

Solutions to Selected Exercises

2.1 Let $h = g \circ f$. Let $A \in \mathcal{H}$. We need to show that $h^{-1}(A) \in \mathcal{F}$. We claim that $h^{-1}(A) = f^{-1}(g^{-1}(A))$. Because g is \mathcal{G}/\mathcal{H} -measurable, $g^{-1}(A) \in \mathcal{G}$ and thus because f is \mathcal{F}/\mathcal{G} -measurable, $f^{-1}(g^{-1}(A))$ is \mathcal{F} -measurable, thus completing the proof, once we show that the claim holds. To show the claim, we show two-sided containment. For showing $h^{-1}(A) \subset f^{-1}(g^{-1}(A))$ let $x \in h^{-1}(A)$. Thus, $h(x) \in A$. By definition, $h(x) = g(f(x)) \in A$. Hence, $f(x) \in g^{-1}(A)$ and thus $x \in f^{-1}(g^{-1}(A))$. For the other direction let $x \in f^{-1}(g^{-1}(A))$. This implies that $f(x) \in g^{-1}(A)$, which implies that $h(x) = g(f(x)) \in A$.

2.3 Since $X(u) \in \mathcal{V}$ for all $u \in \mathcal{U}$ we have $X^{-1}(\mathcal{V}) = \mathcal{U}$. Therefore $\mathcal{U} \in \Sigma_X$. Suppose that $U \in \Sigma_X$, then by definition there exists a $V \in \Sigma$ such that $X^{-1}(V) = U$. Because Σ_X is a σ -algebra we have $V^c \in \Sigma$ and by definition of Σ_X we have $U^c = X^{-1}(V^c) \in \Sigma_X$. Therefore Σ_X is closed under complements. Finally let $(U_i)_i$ be a countable sequence with $U_i \in \Sigma_X$. Then $\bigcup_i U_i = X^{-1}(\bigcup_i X(U_i)) \in \Sigma_X$, which means that Σ_X is closed under countable unions and the proof is completed.

2.5

(a) Let \mathcal{A} be the set of all σ -algebras that contain \mathcal{G} and define

$$\mathcal{F}^* = \bigcap_{\mathcal{F} \in \mathcal{A}} \mathcal{F}.$$

We claim that \mathcal{F}^* is the smallest σ -algebra containing \mathcal{G} . Clearly \mathcal{F}^* contains \mathcal{G} and is contained in all σ -algebras containing \mathcal{G} . Furthermore, by definition it contains exactly those A that are in every σ -algebra that contains \mathcal{G} . It remains to show that \mathcal{F}^* is a σ -algebra. Since $\Omega \in \mathcal{F}$ for all $\mathcal{F} \in \mathcal{A}$ it follows that $\Omega \in \mathcal{F}^*$. Now suppose that $A \in \mathcal{F}^*$. Then $A \in \mathcal{F}$ for all $\mathcal{F} \in \mathcal{A}$ and $A^c \in \mathcal{F}$ for all $\mathcal{F} \in \mathcal{A}$. Therefore $A^c \in \mathcal{F}^*$. Therefore \mathcal{F}^* is closed under complements. Finally, suppose that $(A_i)_i$ is a family in \mathcal{F}^* . Then $(A_i)_i$ are families in \mathcal{F} for all $\mathcal{F} \in \mathcal{A}$ and so $\bigcup_i A_i \in \mathcal{F}$ for all $\mathcal{F} \in \mathcal{A}$ and again we have $\bigcup_i A_i \in \mathcal{F}^*$. Therefore \mathcal{F}^* is a σ -algebra.

(b) Define $\mathcal{H} = \{A : X^{-1}(A) \in \mathcal{F}\}$. Then $\Omega \in \mathcal{H}$ and for $A \in \mathcal{H}$ we have $X^{-1}(A^c) = X^{-1}(A)^c$ so $A^c \in \mathcal{H}$. Furthermore, for $(A_i)_i$ with $A_i \in \mathcal{H}$ we

have

$$X^{-1}\left(\bigcup_i A_i\right) = \bigcup_i X^{-1}(A_i).$$

Therefore \mathcal{H} is a σ -algebra on Ω and by definition $\sigma(\mathcal{G}) \subseteq \mathcal{H}$. Now for any $A \in \mathcal{H}$ we have $f^{-1}(A) \in \mathcal{F}$ by definition. Therefore $f^{-1}(A) \in \mathcal{F}$ for all $A \in \sigma(\mathcal{G})$.

- (c) We need to show that $\mathbb{I}\{A\}^{-1}(B) \in \mathcal{F}$ for all $B \in \mathfrak{B}(\mathbb{R})$. There are four cases. If $\{0, 1\} \in B$, then $\mathbb{I}\{A\}^{-1}(B) = \Omega \in \mathcal{F}$. If $\{1\} \in B$, then $\mathbb{I}\{A\}^{-1}(B) = A \in \mathcal{F}$. If $\{0\} \in B$, then $\mathbb{I}\{A\}^{-1}(B) = A^c \in \mathcal{F}$. Finally, if $\{0, 1\} \cap B = \emptyset$, then $\mathbb{I}\{A\}^{-1}(B) = \emptyset \in \mathcal{F}$. Therefore $\mathbb{I}\{A\}$ is \mathcal{F} -measurable.

2.6 Trivially, $\sigma(X) = \{\emptyset, \mathbb{R}\}$. Hence Y is not $\sigma(X)/\mathfrak{B}(\mathbb{R})$ -measurable because $Y^{-1}([0, 1]) = [0, 1] \notin \sigma(X)$.

2.7 First $\mathbb{P}(\emptyset | B) = \mathbb{P}(\emptyset \cap B) / \mathbb{P}(B) = 0$ and $\mathbb{P}(\Omega | B) = \mathbb{P}(\Omega \cap B) / \mathbb{P}(B) = 1$. Let $(E_i)_i$ be a countable collection of disjoint sets with $E_i \in \mathcal{F}$. Then

$$\begin{aligned} \mathbb{P}\left(\bigcup_i E_i \mid B\right) &= \frac{\mathbb{P}(B \cap \bigcup_i E_i)}{\mathbb{P}(B)} = \frac{\mathbb{P}(\bigcup_i (B \cap E_i))}{\mathbb{P}(B)} \\ &= \sum_i \frac{\mathbb{P}(B \cap E_i)}{\mathbb{P}(B)} = \sum_i \mathbb{P}(E_i | B). \end{aligned}$$

Therefore $\mathbb{P}(\cdot | B)$ satisfies the countable additivity property and the proof is complete.

2.8 Using the definition of conditional probability and the assumption that $\mathbb{P}(A) > 0$ and $\mathbb{P}(B) > 0$ we have:

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(B | A) \mathbb{P}(A)}{\mathbb{P}(B)}.$$

2.9 For part (a),

$$\mathbb{P}(X_1 < 2 | X_2 \text{ is even}) = \frac{\mathbb{P}(X_1 < 2 \text{ and } X_2 \text{ is even})}{\mathbb{P}(X_2 \text{ is even})} = \frac{3/(6^2)}{18/(6^2)} = \frac{1}{6} = \mathbb{P}(X_1 < 2).$$

Therefore $\{X_1 < 2\}$ is independent from $\{X_2 \text{ is even}\}$. For part (b) note that $\sigma(X_1) = \{C \times [6] : C \in 2^{[6]}\}$ and $\sigma(X_2) = \{[6] \times C : C \in 2^{[6]}\}$. It follows that for $|A \cap B| = |A||B|/6^2$ and so

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{|A \cap B|/6^2}{|B|/6^2} = \frac{|A|}{6^2} = \mathbb{P}(A).$$

Therefore A and B are independent.

2.10

- (a) Let $A \in \mathcal{F}$. Then $\mathbb{P}(A \cap \Omega) = \mathbb{P}(A) = \mathbb{P}(A)\mathbb{P}(\Omega)$ and $\mathbb{P}(A \cap \emptyset) = 0 = \mathbb{P}(\emptyset)\mathbb{P}(A)$. Intuitively, Ω and \emptyset happen surely/never respectively, so the occurrence or not of any other event cannot alter their likelihood.
- (b) Let $A \in \mathcal{F}$ satisfy $\mathbb{P}(A) = 1$ and $B \in \mathcal{F}$ be arbitrary. Then $\mathbb{P}(B \cap A^c) \leq \mathbb{P}(A^c) = 0$. Therefore $\mathbb{P}(A \cap B) = \mathbb{P}(A \cap B) + \mathbb{P}(A^c \cap B) = \mathbb{P}(B) = \mathbb{P}(A)\mathbb{P}(B)$. When $\mathbb{P}(A) = 0$ we have $\mathbb{P}(A \cap B) \leq \mathbb{P}(A) = 0 = \mathbb{P}(A)\mathbb{P}(B)$.
- (c) If A and A^c are independent, then $0 = \mathbb{P}(\emptyset) = \mathbb{P}(A \cap A^c) = \mathbb{P}(A)\mathbb{P}(A^c) = \mathbb{P}(A)(1 - \mathbb{P}(A))$. Therefore $\mathbb{P}(A) \in \{0, 1\}$. This makes sense because the knowledge of A provides the knowledge of A^c , so the two events can only be independent if one occurs with probability zero.
- (d) If A is independent of itself, then $\mathbb{P}(A \cap A) = \mathbb{P}(A)^2$. Therefore $\mathbb{P}(A) \in \{0, 1\}$ as before. The intuition is the same as the previous part.
- (e) $\Omega = \{(1, 1), (1, 2), (2, 1), (2, 2)\}$ and $\mathcal{F} = 2^\Omega$.

Ω and A for all $A \in \mathcal{F}$ (16 pairs)

\emptyset and A for all $A \in \mathcal{F} - \Omega$ (15 pairs)

$\{(1, 0), (1, 1)\}$ and $\{(0, 0), (1, 0)\}$

$\{(1, 0), (1, 1)\}$ and $\{(0, 1), (1, 1)\}$

$\{(0, 0), (0, 1)\}$ and $\{(0, 0), (1, 0)\}$

$\{(0, 0), (0, 1)\}$ and $\{(0, 1), (1, 1)\}$

- (f) $\mathbb{P}(X_1 \leq 2, X_1 = X_2) = \mathbb{P}(X_1 = X_2 = 1) + \mathbb{P}(X_1 = X_2 = 2) = 2/9 = \mathbb{P}(X_1 \leq 2)\mathbb{P}(X_1 = X_2)$ because $\mathbb{P}(X_1 \leq 2) = 2/3$ and $\mathbb{P}(X_1 = X_2) = 1/3$.
- (g) If A and B are independent, then $|A \cap B|/n = \mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B) = |A||B|/n^2$. Rearranging shows that $n|A \cap B| = |A||B|$. All steps can be reversed showing the reverse direction.
- (h) Assume that n is prime. By the previous part, $n|A \cap B| = |A||B|$ must hold if A and B are independent of each other. If $|A \cap B| = 0$, the events will be trivial. Hence, assume $|A \cap B| > 0$. Since n is prime, it follows then that n must be either a factor of $|A|$ or a factor of $|B|$. Without loss of generality, assume that it is a factor of $|A|$. This implies $n \leq |A|$. But $|A| \leq n$ also holds, hence $|A| = n$, i.e., A is a trivial event.
- (i) Let X_1 and X_2 be independent Rademacher random variables and $X_3 = X_1X_2$. Clearly these random variables are not mutually independent since X_3 takes multiple values with nonzero probability and is fully determined by X_1 and X_2 . And yet X_3 and X_i are independent for $i \in \{1, 2\}$, which ensures that pairwise independence holds.
- (j) No. Let $\Omega = [6]$ and $\mathcal{F} = 2^\Omega$ and P be the uniform measure. Define events $A = \{1, 3, 4\}$ and $B = \{1, 3, 5\}$ and $C = \{3, 4, 5, 6\}$. Then A and B are clearly dependent and yet

$$\mathbb{P}(A \cap B \cap C) = \mathbb{P}(A \cap B | C)\mathbb{P}(C) = \mathbb{P}(A)\mathbb{P}(B)\mathbb{P}(C).$$

2.11

- (a) $\sigma(X) = (\Omega, \emptyset)$ is trivial. Let Y be another random variable, then X and Y are independent if and only if for all $A \in \sigma(X)$ and $B \in \sigma(Y)$ it holds that $\mathbb{P}(A \cap B) = \mathbb{P}(B)$, which is trivial when $A \in \{\Omega, \emptyset\}$.
- (b) Let A be an event with $\mathbb{P}(A) > 0$. Then $\mathbb{P}(X = x | A) = \mathbb{P}(X = x \cap A) / \mathbb{P}(A) = 1 = \mathbb{P}(X = x)$. Similarly, $\mathbb{P}(X \neq x | A) = 0 = \mathbb{P}(X \neq x)$. Therefore X is independent of all events, including those generated by Y .
- (c) Suppose that A and B are independent. Then $\mathbb{P}(A^c | B) = 1 - \mathbb{P}(A | B) = 1 - \mathbb{P}(A) = \mathbb{P}(A^c)$. Therefore A^c and B are independent and by the same argument so are A^c and B^c as well as A and B^c . The ‘if’ direction follows by noting that $\sigma(X) = \{\Omega, A, A^c, \emptyset\}$ and $\sigma(Y) = \{\Omega, B, B^c, \emptyset\}$ and recalling that every event is independent of Ω or the empty set. For the ‘only if’ note that independence of X and Y means that any pair of events taken from $\sigma(X) \times \sigma(Y)$ are independent, which by the above includes the pair A, B .
- (d) Let $(A_i)_i$ be a countable family of events and $X_i(\omega) = \mathbb{I}\{\omega \in A_i\}$ be the indicator of the i th event. When the random variables/events are pairwise independent, then the above argument goes through unchanged for each pair. In the case of mutual independence the ‘only if’ is again the same. For the ‘if’, suppose that (A_i) are mutually independent. Therefore for any finite subset $K \subset \mathbb{N}$ we have

$$\mathbb{P}\left(\bigcap_{i \in K} A_i\right) = \prod_{i \in K} \mathbb{P}(A_i)$$

The same argument as the previous part shows that for any disjoint finite sets $K, J \subset \mathbb{N}$ we have

$$\mathbb{P}\left(\bigcup_{i \in K} A_i \cup \bigcup_{i \in J} A_i^c\right) = \prod_{i \in K} \mathbb{P}(A_i) \prod_{i \in J} \mathbb{P}(A_i^c).$$

Therefore for any finite set $K \subset \mathbb{N}$ and $(V_i)_{i \in K}$ with $V_i \in \sigma(X_i) = \{\Omega, \emptyset, A_i, A_i^c\}$ it holds that

$$\mathbb{P}\left(\bigcap_{i \in K} V_i\right) = \prod_{i \in K} \mathbb{P}(V_i),$$

which completes the proof that $(X_i)_i$ are mutually independent.

2.12

- (a) Let $A \subset \mathbb{R}$ be an open set. By definition, since f is continuous it holds that $f^{-1}(A)$ is open. But the Borel σ -algebra is generated by all open sets and so $f^{-1}(A) \in \mathfrak{B}(\mathbb{R})$ as required.
- (b) Since $|\cdot| : \mathbb{R} \rightarrow \mathbb{R}$ is continuous and by definition a random variable X on measurable space (Ω, \mathcal{F}) is $\mathcal{F}/\mathfrak{B}(\mathbb{R})$ -measurable it follows by the previous part that $|X|$ is $\mathcal{F}/\mathfrak{B}(\mathbb{R})$ -measurable and therefore a random variable.

(c) Recall that $(X)^+ = \max\{0, X\}$ and $(X)^- = -\min\{0, X\}$. Therefore $(|X|)^+ = |X| = (X)^+ + (X)^-$ and $(|X|)^- = 0$. Recall that $\mathbb{E}[X] = \mathbb{E}[(X)^+] - \mathbb{E}[(X)^-]$ exists if and only if both expectations are defined. Therefore if X is integrable, then $|X|$ is integrable. Now suppose that $|X|$ is integrable, then X is integrable by the dominated convergence theorem.

2.14 Assume without (much) loss of generality that $X_i \geq 0$ for all i . The general case follows by considering positive and negative parts, as usual. First we claim that for any n it holds that

$$\mathbb{E} \left[\sum_{i=1}^n X_i \right] = \sum_{i=1}^n \mathbb{E}[X_i].$$

To show this, note the definition that

$$\begin{aligned} \mathbb{E} \left[\sum_{i=1}^n X_i \right] &= \sup \left\{ \int_{\Omega} h \, d\mathbb{P} : h \text{ is simple and } 0 \leq h \leq \sum_{i=1}^n X_i \right\} \\ &= \sum_{i=1}^n \sup \left\{ \int_{\Omega} h \, d\mathbb{P} : h \text{ is simple and } 0 \leq h \leq X_i \right\}. \end{aligned}$$

Next let $S_n = \sum_{i=1}^n X_i$ and note that by the monotone convergence theorem we have $\lim_{n \rightarrow \infty} \mathbb{E}[S_n] = \mathbb{E}[X]$, which means that

$$\mathbb{E} \left[\sum_{i=1}^{\infty} X_i \right] = \lim_{n \rightarrow \infty} \mathbb{E}[S_n] = \lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{E}[X_i].$$

2.15 Suppose that $X(\omega) = \sum_{i=1}^n \alpha_i \mathbb{I}\{\omega \in A_i\}$ is simple and $c > 0$. Then cX is also simple and

$$\mathbb{E}[cX] = \sum_{i=1}^n c\alpha_i \mathbb{I}\{\omega \in A_i\} = c \sum_{i=1}^n \alpha_i \mathbb{I}\{\omega \in A_i\} = c\mathbb{E}[X].$$

Now suppose that X is positive (but maybe not simple) and $c > 0$, then cX is also positive and

$$\begin{aligned} \mathbb{E}[cX] &= \sup \{ \mathbb{E}[h] : h \text{ is simple and } h \leq cX \} \\ &= \sup \{ \mathbb{E}[ch] : h \text{ is simple and } h \leq X \} \\ &= \sup \{ c\mathbb{E}[h] : h \text{ is simple and } h \leq X \} \\ &= c\mathbb{E}[X]. \end{aligned}$$

Finally for arbitrary random variables and $c > 0$ we have

$$\mathbb{E}[cX] = \mathbb{E}[(cX)^+] - \mathbb{E}[(cX)^-] = c\mathbb{E}[(X)^+] - c\mathbb{E}[(X)^-] = c\mathbb{E}[X].$$

For negative c simply note that $(cX)^+ = -c(X)^-$ and $(cX)^- = -c(X)^+$ and repeat the above argument.

2.16 Suppose $X = \sum_{i=1}^N \alpha_i \mathbb{I}\{A_i\}$ and $Y = \sum_{i=1}^N \beta_i \mathbb{I}\{B_i\}$ are simple functions. Then

$$\begin{aligned} \mathbb{E}[XY] &= \mathbb{E}\left[\sum_{i=1}^N \alpha_i \mathbb{I}\{A_i\} \sum_{i=1}^N \beta_i \mathbb{I}\{B_i\}\right] \\ &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \beta_j \mathbb{P}(A_i \cap B_j) \\ &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \beta_j \mathbb{P}(A_i) \mathbb{P}(A_j) \\ &= \mathbb{E}[X]\mathbb{E}[Y]. \end{aligned}$$

Now suppose that X and Y are arbitrary nonnegative independent random variables. Then

$$\begin{aligned} \mathbb{E}[XY] &= \sup\{\mathbb{E}[h] : h \text{ is simple and } h \leq XY\} \\ &= \sup\{\mathbb{E}[hg] : h \in \sigma(X), g \in \sigma(Y) \text{ are simple and } h \leq X, g \leq Y\} \\ &= \sup\{\mathbb{E}[h]\mathbb{E}[g] : h \in \sigma(X), g \in \sigma(Y) \text{ are simple and } h \leq X, g \leq Y\} \\ &= \sup\{\mathbb{E}[h] : h \in \sigma(X) \text{ is simple and } h \leq X\} \sup\{\mathbb{E}[h] : h \in \sigma(Y) \text{ is simple and } h \leq Y\} \\ &= \mathbb{E}[X]\mathbb{E}[Y]. \end{aligned}$$

Finally, for arbitrary random variables we have via the previous display and the linearity of expectation that

$$\begin{aligned} \mathbb{E}[XY] &= \mathbb{E}[(X)^+ - (X)^-][(Y)^+ - (Y)^-] \\ &= \mathbb{E}[(X)^+(Y)^+] - \mathbb{E}[(X)^+(Y)^-] - \mathbb{E}[(X)^-(Y)^+] + \mathbb{E}[(X)^-(Y)^-] \\ &= \mathbb{E}[(X)^+]\mathbb{E}[(Y)^+] - \mathbb{E}[(X)^+]\mathbb{E}[(Y)^-] - \mathbb{E}[(X)^-]\mathbb{E}[(Y)^+] + \mathbb{E}[(X)^-]\mathbb{E}[(Y)^-] \\ &= \mathbb{E}[(X)^+ - (X)^-]\mathbb{E}[(Y)^+ - (Y)^-]. \end{aligned}$$

2.17 Let X be a standard Rademacher random variable and $Y = X$. Then $\mathbb{E}[X]\mathbb{E}[Y] = 0$ and $\mathbb{E}[XY] = 1$.

2.18 Using the fact that $\int_0^a 1 dx = a$ for $a \geq 0$ and the nonnegativity of X we have

$$X(\omega) = \int_0^\infty \mathbb{I}\{[0, X(\omega)]\}(x) dx.$$

Then by Fubini's theorem we have

$$\begin{aligned} \mathbb{E}[X] &= \mathbb{E}\left[\int_0^\infty \mathbb{I}\{[0, X(\omega)]\}(x) dx\right] \\ &= \int_0^\infty \mathbb{E}[\mathbb{I}\{[0, X(\omega)]\}(x)] dx \\ &= \int_0^\infty \mathbb{P}(X(\omega) \geq x) dx. \end{aligned}$$

3.1 For simplicity, we remove 1 from $[0, 1]$, so we will set $\mathcal{S} = ([0, 1), \mathfrak{B}([0, 1)))$ and $\mathbb{P} \doteq \lambda$ be the uniform (Lebesgue) measure on $[0, 1)$. With some extra bookkeeping the result can also be shown to work when 1 is not removed.

- (a) We have $F_1(x) = \mathbb{I}\{x \in [1/2, 1)\}$, $F_2(x) = \mathbb{I}\{x \in [1/4, 2/4) \cup [3/4, 4/4)\}$, \dots , and for $t \geq 1$, $F_t(x) = \mathbb{I}\{x \in U_t\}$ where $U_t = \cup_{1 \leq s \leq 2^{t-1}} [(2s-1)/2^t, 2s/2^t)$. Since $U_t \in \mathfrak{B}([0, 1))$, F_t are random variables.
- (b) We have $\mathbb{P}(U_t) = \lambda(U_t) = \sum_{s=1}^{2^{t-1}} (1/2^t) = 1/2$.
- (c) For an index set $K \subset \mathbb{N}^+$ let $F_K = (F_k)_{k \in K}$ (that is, $F_K : [0, 1) \rightarrow \{0, 1\}^K$). Note that for $A \in \mathfrak{B}(\mathbb{R}^K)$,

$$\mathbb{P}(F_K \in A) = \mathbb{P}(F_K \in A \cap \{0, 1\}^K). \quad (.26)$$

Further, if $A = \times_{k \in K} A_k$, $A_k \subset \{0, 1\}$,

$$\mathbb{P}(F_K \in A) = \mathbb{P}(F_k \in A_k, k \in K) = \prod_{k \in K} \mathbb{P}(F_k \in A_k). \quad (.27)$$

To see the last equality, first note that if for any $k \in K$, $A_k = \emptyset$ then $\{F_k \in A_k, k \in K\} = \emptyset$ and the equality clearly holds. Hence, assume this does not happen (that is, for any $k \in K$, $A_k \neq \emptyset$). Second, note that $\{F_k \in A_k, k \in K\} = \{F_k \in A_k, k \in A_k, A_k \neq \{0, 1\}\}$. Hence, without loss of generality assume that for all $k \in K$, $A_k = \{a_k\}$ for some $a_k \in \{0, 1\}$. Now, define $U^{(a)} = U$ if $a = 1$ and $U^{(a)} = U^c (= [0, 1) \setminus U)$ when $a = 0$. Then, $\{x : F_k(x) \in A_k\} = \cup_{a \in A_k} U_k^{(a)} = U_k^{(a_k)}$. This implies that $\{F_k \in A_k, k \in K\} = \cap_{k \in K} U_k^{(a_k)}$. Now, a tedious but elementary calculation shows that $\lambda(\cap_{k \in K} U_k^{(a_k)}) = \prod_{k \in K} \lambda(U_k^{(a_k)})$.

Finally, to show that $(F_t)_{t \geq 1}$ is an independent sequence, we need to show that any finite subcollection of F_1, F_2, \dots are independent. For this, it is sufficient to show that for any disjoint index sets $I, J \subset \mathbb{N}^+$, $A \in \mathfrak{B}(\mathbb{R}^I)$, $B \in \mathfrak{B}(\mathbb{R}^J)$, $\mathbb{P}(F_I \in A, F_J \in B) = \mathbb{P}(F_I \in A) \mathbb{P}(F_J \in B)$. By (.26), without loss of generality, we can assume that $A \subset \{0, 1\}^I$ and $B \subset \{0, 1\}^J$. By Caratheodory's theorem, Theorem 2.2, it suffices to consider the case when A and B are rectangles: $A = \times_{i \in I} A_i$ and $B = \times_{j \in J} B_j$. The result then follows from (.27).

- (d) It follows directly from the definition of independence, that any subsequence of an independent sequence is also an independent sequence. That $\mathbb{P}(X_{m,t} = 0) = \mathbb{P}(X_{m,t} = 1) = 0.5$ follows from Part (b).
- (e) We need to show that $X_t = \sum_{i \geq 1} X_{m,t} 2^{-i}$ is uniformly distributed. In fact, X_t is identically distributed to $Y = \sum_{i \geq 1} F_i 2^{-i}$, because both are a function of independent sequence of uniformly distributed random variables. However, $Y(x) = x$, the identity map, hence, it is uniformly distributed (i.e, for any $U \in \mathfrak{B}([0, 1))$, $\mathbb{P}(Y \in U) = \lambda(U)$).

3.5 Let $A \in \mathcal{G}$ and suppose that $X(\omega) = \mathbb{I}_A(\omega)$. Then $\int_{\mathcal{X}} X(x) K(\omega, dx) = K(\omega, A)$, which is \mathcal{F} -measurable by the definition of a probability kernel.

The result extends to simple functions by linearity. For nonnegative X let $X_n \uparrow X$ be a monotone increasing sequence of simple functions converging pointwise to X [Kallenberg, 2002, Lemma 1.11]. Then $U_n(\omega) = \int_{\mathcal{X}} X_n(x)K(\omega, dx)$ is \mathcal{F} -measurable. Monotone convergence ensures that $\lim_{n \rightarrow \infty} U_n(\omega) = \int_{\mathcal{X}} \lim_{n \rightarrow \infty} X_n(x)K(\omega, dx) = \int_{\mathcal{X}} X(x)K(\omega, dx) = U(\omega)$. Hence $\lim_{n \rightarrow \infty} U_n(\omega) = U(\omega)$ and a point-wise convergent sequence of measurable functions is measurable, it follows that U is \mathcal{F} -measurable. The result for arbitrary X follows by decomposing into positive and negative parts.

4.6

(a) The statement is true. Let i be a suboptimal arm. By Lemma 4.2 we have

$$0 = \lim_{n \rightarrow \infty} \frac{R_n(\pi, \nu)}{n} = \limsup_{n \rightarrow \infty} \sum_{i=1}^K \frac{\mathbb{E}[T_i(n)]\Delta_i}{n} \geq \limsup_{n \rightarrow \infty} \frac{\mathbb{E}[T_i(n)]}{n} \Delta_i.$$

Hence $\limsup_{n \rightarrow \infty} \mathbb{E}[T_i(n)]/n \leq 0 \leq \liminf_{n \rightarrow \infty} \mathbb{E}[T_i(n)]/n$ and so $\lim_{n \rightarrow \infty} \mathbb{E}[T_i(n)]/n = 0$ for suboptimal arms i . Since $\sum_{i=1}^K \mathbb{E}[T_i(n)]/n = 1$ it follows that $\lim_{n \rightarrow \infty} \sum_{i:\Delta_i=0} \mathbb{E}[T_i]/n = 1$.

(b) The statement is false. Consider a two-armed bandit for which the second arm is suboptimal and an algorithm that chooses the second arm in rounds $t \in \{1, 2, 4, 8, 16, \dots\}$.

4.7

(a) For $i \in \{1, 2\}$ let μ_i be the mean of P_i . The optimal policy chooses $A_t \in \operatorname{argmax}_i \mu_i$.

(b) Let $p_t(x_1, \dots, x_{t-1})$ be a function such that

$$\mathbb{P}(T_1(n) \geq t \mid S_{t-1}, T_1(n) \geq t-1)$$

5.17 We use the Cramer-Chernoff method:

$$\begin{aligned} \mathbb{P}\left(\sum_{t=1}^n (X_t - \mu_t - \alpha_t) \geq \frac{1}{\lambda} \log\left(\frac{1}{\delta}\right)\right) &= \mathbb{P}\left(\exp\left(\lambda \sum_{t=1}^n (X_t - \mu_t - \alpha_t)\right) \geq \frac{1}{\delta}\right) \\ &\leq \delta \mathbb{E}\left[\exp\left(\lambda \sum_{t=1}^n (X_t - \mu_t - \alpha_t)\right)\right]. \end{aligned}$$

All that remains is to show that the term inside the expectation is a supermartingale. Using the fact that $\exp(x) \leq 1 + x + x^2$ for $x \leq 1$ and $1 + x \leq \exp(x)$ for all $x \in \mathbb{R}$ we have

$$\begin{aligned} \mathbb{E}_{t-1}[\exp(\lambda(X_t - \mu_t - \alpha_t))] &= \exp(-\lambda\alpha_t) \mathbb{E}_{t-1}[\exp(\lambda(X_t - \mu_t))] \\ &\leq \exp(-\lambda\alpha_t) (1 + \lambda^2 \mathbb{E}_{t-1}[(X_t - \mu_t)^2]) \\ &\leq \exp(\lambda(\lambda \mathbb{E}_{t-1}[(X_t - \mu_t)^2] - \alpha_t)) \leq 1. \end{aligned}$$

Therefore $\exp(\lambda \sum_{t=1}^n (X_t - \mu_t - \alpha_t))$ is a supermartingale, which completes the proof.

5.18 By assumption $\mathbb{P}(X_t \leq x) \leq x$, which means that for $\lambda < 1$,

$$\begin{aligned} \mathbb{E}[\exp(\lambda \log(1/X_t))] &= \int_0^\infty \mathbb{P}(\exp(\lambda \log(1/X_t)) \geq x) dx \\ &= 1 + \int_1^\infty \mathbb{P}(X_t \leq x^{-1/\lambda}) dx \leq 1 + \int_1^\infty x^{-1/\lambda} dx = \frac{1}{1-\lambda}. \end{aligned}$$

Applying the Cramer-Chernoff method,

$$\begin{aligned} \mathbb{P}\left(\sum_{t=1}^n \log(1/X_t) \geq \varepsilon\right) &= \mathbb{P}\left(\exp\left(\lambda \sum_{t=1}^n \log(1/X_t)\right) \geq \exp(\lambda \varepsilon)\right) \\ &\leq \exp(-\lambda \varepsilon) \mathbb{E}\left[\exp\left(\lambda \sum_{t=1}^n \log(1/X_t)\right)\right] \leq \left(\frac{1}{1-\lambda}\right)^n \exp(-\lambda \varepsilon). \end{aligned}$$

Choosing $\lambda = (\varepsilon - n)/\varepsilon$ completes the claim.

5.20 Let \mathcal{P} be the set of measures on $([0, 1], \mathfrak{B}([0, 1]))$ and for $q \in \mathcal{P}$ let μ_q be its mean. The theorem will be established by induction over n . The claim is immediate when $x > n$ or $n = 1$. Assume that $n \geq 2$ and $x \in (1, n]$ and the theorem holds for $n - 1$. Then

$$\begin{aligned} \mathbb{P}\left(\sum_{t=1}^n \mathbb{E}[X_t | \mathcal{F}_{t-1}]\right) &= \mathbb{E}\left[\mathbb{P}\left(\sum_{t=2}^n \mathbb{E}[X_t | \mathcal{F}_{t-1}] \geq x - \mathbb{E}[X_1 | \mathcal{F}_0] \mid \mathcal{F}_0\right)\right] \\ &\leq \mathbb{E}\left[f_{n-1}\left(\frac{x - \mathbb{E}[X_1 | \mathcal{F}_0]}{1 - X_1}\right)\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[f_{n-1}\left(\frac{x - \mathbb{E}[X_1 | \mathcal{F}_0]}{1 - X_1}\right) \mid \mathcal{F}_0\right]\right] \\ &\leq \sup_{q \in \mathcal{P}} \int_0^1 f_{n-1}\left(\frac{x - \mu_q}{1 - y}\right) dq(y), \end{aligned}$$

where the first inequality follows from the inductive hypothesis and the fact that $\sum_{t=2}^n X_t/(1 - X_1) \leq 1$ almost surely. The result is completed by proving that for all $q \in \mathcal{P}$,

$$F_n(q) \doteq \int_0^1 f_{n-1}\left(\frac{x - \mu_q}{1 - y}\right) dq(y) \leq f_n(x). \quad (.28)$$

Let $q \in \mathcal{P}$ have mean μ and $y_0 = \max(0, 1 - x + \mu)$. In Lemma .4 below it is shown that

$$f_{n-1}\left(\frac{x - \mu}{1 - y}\right) \leq \frac{1 - y}{1 - y_0} f_{n-1}\left(\frac{x - \mu}{1 - y_0}\right),$$

which after integrating implies that

$$F_n(q) \leq \frac{1 - \mu}{1 - y_0} f_{n-1}\left(\frac{x - \mu}{1 - y_0}\right).$$

Considering two cases. First, when $y_0 = 0$ the display shows that $F_n(q) \leq$

$(1 - \mu)f_{n-1}(x - \mu)$. On the other hand, if $y_0 > 0$ then $x - 1 < \mu \leq 1$ and $F_n(q) \leq (1 - \mu)/(x - \mu) \leq (1 - (x - 1))f_{n-1}(x - (x - 1))$. Combining the two cases we have

$$F_n(q) \leq \sup_{\mu \in [0,1]} (1 - \mu)f_{n-1}(1 - \mu) = f_n(x).$$

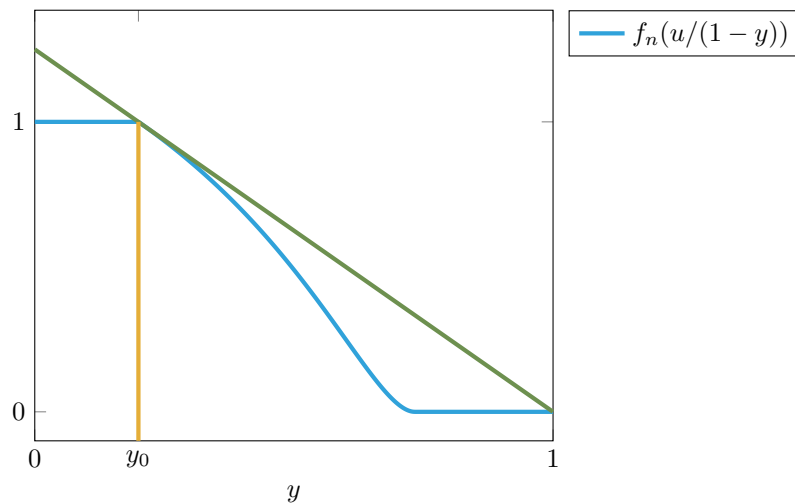
LEMMA .4 Suppose that $n \geq 1$ and $u \in (0, n]$ and $y_0 = \max(0, 1 - u)$. Then

$$f_n\left(\frac{u}{1-y}\right) \leq \frac{1-y}{1-y_0} f_{n-1}\left(\frac{u}{1-y_0}\right) \quad \text{for all } y \in [0, 1].$$

Proof The lemma is equivalent to the claim that the line connecting $(y_0, f_n(u/(1 - y_0)))$ and $(1, 0)$ lies above $f_n(u/(1 - y))$ for all $y \in [0, 1]$ (see figure below). This is immediate for $n = 1$ when $f_n(u/(1 - y)) = \mathbb{I}\{y \leq 1 - u\}$. For larger n basic calculus shows that $f_n(u/(1 - y))$ is concave as a function of y on $[1 - u, 1 - u/n]$ and

$$\left. \frac{\partial}{\partial y} f_n(u/(1 - y)) \right|_{y=1-u} = -1/u.$$

Since $f_n(1) = 1$ this means that the line connecting $(1 - u, 1)$ and $(1, 0)$ lies above $f_n(u/(1 - y))$. This completes the proof when $y_0 = 1 - u$. Otherwise $y_0 \in [1 - u, 1 - u/n]$ and the result follows by concavity of $f_n(u/(1 - y))$ on this interval. \square



8.1 Following the hint, $F \leq \exp(-a)/(1 - \exp(-a))$ where $a = \varepsilon^2/2$. Reordering $\exp(-a)/(1 - \exp(-a)) \leq 1/a$ gives $1 + a \leq \exp(a)$ which is well known (and

easy to prove). Then

$$\begin{aligned} \sum_{t=1}^n \frac{1}{f(t)} &\leq \sum_{t=1}^{20} \frac{1}{f(t)} + \int_{20}^{\infty} \frac{dt}{f(t)} \leq \sum_{t=1}^{20} \frac{1}{f(t)} + \int_{20}^{\infty} \frac{dt}{t \log(t)^2} \\ &= \sum_{t=1}^{20} \frac{1}{f(t)} + \frac{1}{\log(20)} \leq \frac{5}{2}. \end{aligned}$$

9.1 Clearly (M_t) is \mathbb{F} -adapted. Then by Jensen's inequality and convexity of the exponential function,

$$\begin{aligned} \mathbb{E}[M_t | \mathcal{F}_{t-1}] &= \exp\left(\lambda \sum_{s=1}^t X_s\right) \mathbb{E}[\exp(\lambda X_t) | \mathcal{F}_{t-1}] \\ &\geq \exp\left(\lambda \sum_{s=1}^t X_s\right) \exp(\lambda \mathbb{E}[X_t | \mathcal{F}_{t-1}]) \\ &= \exp\left(\lambda \sum_{s=1}^{t-1} X_s\right) \quad \text{a.s.} \end{aligned}$$

Hence M_t is a \mathbb{F} -supermartingale.

10.1 Let g be as in the hint. We have

$$g'(x) = x \left(\frac{1}{(p+x)(1-(p+x))} - 4 \right).$$

Clearly, $g'(0) = 0$. Further, since $q(1-q) \leq 1/4$ for any $q \in [0, 1]$, $g'(x) \geq 0$ for $x > 0$ and $g'(x) \leq 0$ for $x < 0$. Hence, g is increasing for positive x and decreasing for negative x . Thus, $x = 0$ is a minimizer of g . Here, $g(0) = 0$, and so $g(x) \geq 0$ over $[-p, 1-p]$.

10.3 We have $g(\lambda, \mu) = -\lambda\mu + \log(1 + \mu(e^\lambda - 1))$. Taking derivatives, $\frac{d}{d\mu}g(\lambda, \mu) = -\lambda + \frac{e^\lambda - 1}{1 + \mu(e^\lambda - 1)}$ and $\frac{d^2}{d\mu^2}g(\lambda, \mu) = -\frac{(e^\lambda - 1)^2}{(1 + \mu(e^\lambda - 1))^2} \leq 0$, showing that $g(\lambda, \cdot)$ is concave as suggested in the hint. Now, let $S_p = \sum_{t=1}^p (X_t - \mu_t)$, $p \in [n]$ and let $S_0 = 0$. Then for $p \in [n]$,

$$\mathbb{E}[\exp(\lambda S_p)] = \mathbb{E}[\exp(\lambda S_{p-1}) \mathbb{E}[\exp(\lambda(X_p - \mu_p)) | \mathcal{F}_{p-1}]],$$

and, by note Item 2, $\mathbb{E}[\exp(\lambda(X_p - \mu_p)) | \mathcal{F}_{p-1}] \leq \exp(g(\lambda, \mu_p))$. Hence, using that μ_n is not random,

$$\mathbb{E}[\exp(\lambda S_p)] \leq \mathbb{E}[\exp(\lambda S_{p-1}) \exp(g(\lambda, \mu_p))].$$

Chaining this inequalities, using that $S_0 = 0$ together with that $g(\lambda, \cdot)$ is concave, we get

$$\mathbb{E}[\exp(\lambda S_n)] \leq \left(\exp\left(\frac{1}{n} \sum_{t=1}^n g(\lambda, \mu_t)\right) \right)^n \leq \exp(n g(\lambda, \mu)).$$

Thus,

$$\begin{aligned}\mathbb{P}(\hat{\mu} - \mu \geq \varepsilon) &= \mathbb{P}(\exp(\lambda S_n) \geq \exp(\lambda n\varepsilon)) \\ &\leq \mathbb{E}[\exp(\lambda S_n)] \exp(-\lambda n\varepsilon) \\ &\leq (\mu \exp(\lambda(1 - \mu - \varepsilon)) + (1 - \mu) \exp(-\lambda(\mu + \varepsilon)))^n.\end{aligned}$$

From this point, repeat the proof of Lemma 10.2 word by word.

10.4 Abbreviate $p_\theta(x) = \frac{dP_\theta}{dh}(x)$.

- (a) Clearly $p_\theta(x) \geq 0$. By definition, for $B \in \mathfrak{B}(\mathbb{R})$, $P_\theta(B) = \int_B p_\theta(x) dh(x)$. Hence $P_\theta(B) \geq 0$. Furthermore,

$$P_\theta(\mathbb{R}) = \int_{\mathbb{R}} \exp(\theta T(x) - A(\theta)) dh(x) = \exp(-A(\theta)) \int_{\mathbb{R}} \exp(\theta T(x)) dh(x) = 1.$$

Additivity is immediate since $\int_B f dh + \int_C f dh = \int_{B \cup C} f dh$ for disjoint B, C .

- (b) Using the chain rule and passing the derivative under the integral yields the result:

$$\begin{aligned}A'(\theta) &= \frac{\frac{d}{d\theta} \int_{\mathbb{R}} \exp(\theta T(x)) dh(x)}{\int_{\mathbb{R}} \exp(\theta T(x)) dh(x)} \\ &= \frac{\int_{\mathbb{R}} T(x) \exp(\theta T(x)) dh(x)}{\int_{\mathbb{R}} \exp(\theta T(x)) dh(x)} \\ &= \int_{\mathbb{R}} T(x) \exp(\theta T(x) - A(x)) dh(x) \\ &= \int_{\mathbb{R}} T(x) p_\theta(x) dh(x) \\ &= \mathbb{E}_\theta[T].\end{aligned}$$

In order to justify the exchange of integral and derivative use the identity that for all sufficiently small $\varepsilon > 0$ and all $a > 0$,

$$a \leq \frac{\exp(a\varepsilon) + \exp(-a\varepsilon)}{\varepsilon}.$$

Hence for $\theta \in \text{int}(\text{dom}(A))$ there exists a neighborhood N of θ such that for all $\psi \in N$,

$$|T(x)| \exp(\psi T(x)) \leq \phi(x) = \frac{\exp((\theta + \varepsilon)T(x)) + \exp((\theta - \varepsilon)T(x))}{\varepsilon}.$$

Since $\phi(x)$ is integrable for sufficiently small ε it follows by the dominated convergence theorem that the derivative and integral can be exchanged.

(c) When $X \sim P_\theta$ we have

$$\begin{aligned}\mathbb{E}[\exp(\lambda T(X))] &= \int_{\mathbb{R}} \exp(\lambda T(x)) \frac{dP_\theta}{dh}(x) dh(x) \\ &= \int_{\mathbb{R}} \exp(\lambda T(x)) \exp(\theta T(x) - A(\theta)) dh(x) \\ &= \exp(-A(\theta)) \int_{\mathbb{R}} \exp((\theta + \lambda)T(x)) dh(x) \\ &= \exp(-A(\theta)) \exp(A(\theta + \lambda)) \\ &= \exp(A(\theta + \lambda) - A(\theta)).\end{aligned}$$

(d) This is another straightforward calculation:

$$\begin{aligned}d(\theta, \theta') &= \int_{\mathbb{R}} (\theta T(x) - A(\theta) - \theta' T(x) + A(\theta')) \exp(\theta T(x) - A(\theta)) dh(x) \\ &= A(\theta') - A(\theta) + (\theta - \theta') \int_{\mathbb{R}} T(x) \exp(\theta T(x) - A(\theta)) dh(x) \\ &= A(\theta') - A(\theta) - (\theta' - \theta) A'(\theta).\end{aligned}$$

Curiously, this is the Bregman divergence between θ' and θ induced by the convex function A . See Section 26.3 for definitions. Bregman divergence

(e) The Crammer-Chernoff method is the solution. Let $\lambda = n(\theta' - \theta)$. Then

$$\begin{aligned}\mathbb{P}_\theta(\hat{t} \geq \mathbb{E}_{\theta'}[T]) &= \mathbb{P}_\theta(\exp(\lambda \hat{t}) \geq \exp(\lambda \mathbb{E}_{\theta'}[T])) \\ &\leq \mathbb{E}_\theta[\exp(\lambda \hat{t})] \exp(-\lambda \mathbb{E}_{\theta'}[T]) \\ &= \prod_{t=1}^n \mathbb{E}_\theta[\exp(\lambda T(X_t)/n)] \exp(-\lambda A'(\theta')) \\ &= \exp(n(A(\theta + \lambda/n) - A(\theta)) - \lambda A'(\theta')) \\ &= \exp(n(A(\theta') - A(\theta) - (\theta' - \theta)A'(\theta))) \\ &= \exp(-nd(\theta', \theta)).\end{aligned}$$

A symmetric calculation shows that for $\theta' < \theta$,

$$\mathbb{P}_\theta(\hat{t} \leq \mathbb{E}_{\theta'}[T]) \leq \exp(-nd(\theta', \theta)).$$

(f) Let $h = \delta_0 + \delta_1$ be the sum of two Dirac measures and $T(x) = x$. Then $p_\theta(0) = 1/(1 + \exp(\theta))$ and $p_\theta(1) = 1/(1 + \exp(-\theta))$ and the domain of A is \mathbb{R} .

(g) Let h be Gaussian with mean 0 and variance 1 and $T(x) = x$. Then

$$\begin{aligned}A(\theta) &= \log \left(\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \exp(x\theta - x^2/2) dx \right) \\ &= \log \left(\frac{1}{\sqrt{2\pi}} \exp(\theta^2/2) \int_{\mathbb{R}} \exp(-(x - \theta)^2/2) dx \right) \\ &= \theta^2/2.\end{aligned}$$

Hence $p_\theta(x) = \exp(\theta x - \theta^2/2)$ and

$$\begin{aligned} \int_A p_\theta(x) dh(x) &= \frac{1}{\sqrt{2\pi}} \int_A \exp(\theta x - \theta^2/2) \exp(-x^2/2) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_A \exp((\theta - x)^2/2) dx, \end{aligned}$$

which indeed corresponds to a Gaussian with mean θ .

10.5 When the exponential family is in canonical form the mean of P_θ is $\mu(\theta) = \mathbb{E}_\theta[T] = A'(\theta)$. Since A is strictly convex by the assumption that \mathcal{M} is nonsingular it follows that $\mu(\theta)$ is monotone increasing and hence invertible. Let $\mu_{\text{sup}} = \sup_{\theta \in \Theta} \mu(\theta)$ and $\mu_{\text{inf}} = \inf_{\theta \in \Theta} \mu(\theta)$ and define

$$\hat{\theta}(x) = \begin{cases} \sup \Theta & \text{if } x \geq \mu_{\text{sup}} \\ \inf \Theta & \text{if } x \leq \mu_{\text{inf}} \\ \mu^{-1}(x) & \text{otherwise.} \end{cases}$$

The function $\hat{\theta}$ is the bridge between the empirical mean and the maximum likelihood estimator of θ . Precisely, let X_1, \dots, X_n be independent and identically distributed from P_θ and $\hat{\mu}_n = \frac{1}{n} \sum_{t=1}^n X_t$. Then provided that $\hat{\theta}_n = \hat{\theta}(\hat{\mu}_n) \in \Theta$, then $\hat{\theta}_n$ is the maximum likelihood estimator of θ ,

$$\hat{\theta}_n = \operatorname{argmax}_{\theta \in \Theta} \prod_{t=1}^n \frac{dP_\theta}{dh}(X_t).$$

There is an irritating edge case that $\hat{\mu}_n$ does not lie in the range of $\mu : \Theta \rightarrow \mathbb{R}$. When this occurs there is no maximum likelihood estimator.

Part I: Algorithm

Then define $\underline{d}(x, y) = \mathbb{I}\{x \leq y\} \lim_{z \downarrow x} d(z, y)$ and $\bar{d}(x, y) = \mathbb{I}\{x \geq y\} \lim_{z \uparrow x} d(z, y)$. The algorithm chooses $A_t = t$ for the first K rounds and subsequently $A_t = \operatorname{argmax}_i U_i(t)$ where

$$U_i(t) = \sup \left\{ \tilde{\theta} \in \Theta : \bar{d}(\hat{\theta}_i(t-1), \tilde{\theta}) \leq \frac{\log(f(T_i(t-1)))}{T_i(t-1)} \right\}.$$

Part II: Concentration

Given a fixed $\theta \in \Theta$ and independent random variables X_1, \dots, X_n sampled from \mathbb{P}_θ and $\hat{t}_s = \frac{1}{s} \sum_{u=1}^s T(X_u)$ and $\hat{\theta}_s = \hat{\theta}(\hat{t}_s)$. Let $\tilde{\theta} \in \Theta$ be such that $\underline{d}(\tilde{\theta}, \theta) = \varepsilon > 0$. Then

$$\mathbb{P}(\underline{d}(\hat{\theta}_s, \theta) \geq \varepsilon) \leq \mathbb{P}(\hat{t} \geq t(\tilde{\theta})) \leq \exp(-sd(\tilde{\theta}, \theta)) = \exp(-s\varepsilon), \quad (.29)$$

where the second inequality follows from Part (e) of Exercise 10.4. Similarly

$$\mathbb{P}(\bar{d}(\hat{\theta}_s, \theta) \geq \varepsilon) \leq \exp(-s\varepsilon). \quad (.30)$$

Define random variable τ by

$$\tau = \min \left\{ t : \underline{d}(\hat{\theta}_s, \theta - \varepsilon) < \frac{\log(f(t))}{s} \text{ for all } s \in [n] \right\}.$$

In order to bound the expectation of τ we need a connection between $\underline{d}(\hat{\theta}_s, \theta - \varepsilon)$ and $\underline{d}(\hat{\theta}_s, \theta)$. Let $x \leq y - \varepsilon$ and $g(z) = d(x, z)$. Then

$$\begin{aligned} g(y) &= g(y - \varepsilon) + \int_{y-\varepsilon}^y g'(z) dz \\ &= g(y - \varepsilon) + \int_{y-\varepsilon}^y (z - x) A''(z) dz \\ &\geq g(y - \varepsilon) + \inf_{z \in [y-\varepsilon, y]} A''(z) \int_{y-\varepsilon}^y (z - x) dz \\ &= g(y - \varepsilon) + \frac{1}{2} \inf_{z \in [y-\varepsilon, y]} A''(z) \varepsilon (2y - 2x - \varepsilon) \\ &\geq g(y - \varepsilon) + \frac{\varepsilon^2 \inf_{z \in [y-\varepsilon, y]} A''(z)}{2}. \end{aligned}$$

Note that $\inf_{z \in [y-\varepsilon, y]} A''(z) > 0$ is guaranteed because A'' is continuous and $[y - \varepsilon, y]$ is compact and because \mathcal{M} was assumed to be nonsingular. Using this, the expectation of τ is bounded by

$$\begin{aligned} \mathbb{E}[\tau] &= \sum_{t=1}^n \mathbb{P}(\tau \geq t) \\ &\leq \sum_{t=1}^n \sum_{s=1}^n \mathbb{P} \left(\underline{d}(\hat{\theta}_s, \theta - \varepsilon) \geq \frac{\log(f(t))}{s} \right) \\ &\leq \sum_{t=1}^n \sum_{s=1}^n \mathbb{P} \left(\underline{d}(\hat{\theta}_s, \theta) \geq \frac{\varepsilon^2 \inf_{z \in [\theta-\varepsilon, \theta]} A''(z)}{2} + \frac{\log(f(t))}{s} \right) \\ &\leq \sum_{t=1}^n \sum_{s=1}^n \frac{\exp(-s \inf_{z \in [\theta-\varepsilon, \theta]} A''(z) \varepsilon^2 / 2)}{f(t)} \\ &= O(1), \end{aligned} \tag{.31}$$

where the last inequality follows from Eq. (.29) and the final inequality is the same calculation as in the proof of Lemma 10.3. Next let

$$\kappa = \min \{ s \geq 1 : \hat{\theta}_u - \theta < \varepsilon \text{ for all } u \geq s \}.$$

The expectation of κ is easily bounded using Eq. (.30),

$$\mathbb{E}[\kappa] \leq \sum_{s=1}^n \sum_{u=s}^{\infty} (\exp(-ud(\theta + \varepsilon, \theta))) = O(1), \tag{.32}$$

where we used the fact that \mathcal{M} is non-singular to ensure strict positivity of the divergences.

Part III: Bounding $\mathbb{E}[T_i(n)]$

For each arm i let $\hat{\theta}_{is} = \hat{\theta}(\hat{\mu}_{is})$. Now fix a suboptimal arm i and let $\varepsilon < (\theta_1 - \theta_i)/2$ and

$$\tau = \min \left\{ t : \underline{d}(\hat{\theta}_s, \theta - \varepsilon) < \frac{\log(f(t))}{s} \text{ for all } s \in [n] \right\}.$$

Then define

$$\kappa = \min\{s \geq 1 : \hat{\theta}_{iu} < \theta_i + \varepsilon \text{ for all } u \geq s\}.$$

Then by Eq. (.31) and Eq. (.32), $\mathbb{E}[\tau] = O(1)$ and $\mathbb{E}[\kappa] = O(1)$. Suppose that $t \geq \tau$ and $T_i(t-1) \geq \kappa$ and $A_t = i$. Then $U_i(t) \geq U_1(t) \geq \theta_1 - \varepsilon$ and hence

$$d(\theta_i + \varepsilon, \theta_1 - \varepsilon) < \frac{\log(f(n))}{T_i(t-1)}.$$

From this we conclude that

$$T_i(n) \leq 1 + \tau + \kappa + \frac{\log(f(n))}{d(\theta_i + \varepsilon, \theta_1 - \varepsilon)}.$$

Taking expectations and limits shows that

$$\limsup_{n \rightarrow \infty} \frac{\mathbb{E}[T_i(n)]}{\log(n)} \leq \frac{1}{d(\theta_i + \varepsilon, \theta_1 - \varepsilon)}.$$

Since the above holds for all sufficiently small $\varepsilon > 0$ and the divergence d is continuous it follows that

$$\limsup_{n \rightarrow \infty} \frac{\mathbb{E}[T_i(n)]}{\log(n)} \leq \frac{1}{d(\theta_i, \theta_1)}$$

for all suboptimal arms i . The result follows from the fundamental regret decomposition lemma (Lemma 4.2).

10.6 For simplicity we assume the first arm is uniquely optimal. Define

$$\bar{d}(x, y) = \mathbb{I}\{x \geq y\} \lim_{z \uparrow x} d(z, y), \quad \underline{d}(x, y) = \mathbb{I}\{x \leq y\} \lim_{z \downarrow x} d(x, y).$$

Let $\mu(\theta) = \int_{\mathbb{R}} x dP_{\theta}(x)$ and $t(\theta) = \mathbb{E}_{\theta}[T] = A'(\theta)$ and $\mathcal{T} = \{t(\theta) : \theta \in \Theta\}$. Define $\hat{\theta} : \mathbb{R} \rightarrow \text{cl}(\Theta)$ by

$$\hat{\theta}(x) = \begin{cases} t^{-1}(x), & \text{if } x \in \mathcal{T}; \\ \sup \Theta, & \text{if } x \geq \sup \mathcal{T}; \\ \inf \Theta, & \text{if } x \leq \inf \mathcal{T}. \end{cases}$$

The algorithm is a generalization of KL-UCB. Let

$$\hat{t}_i(t) = \frac{1}{T_i(t)} \sum_{s=1}^t \mathbb{I}\{A_s = i\} T(X_s) \quad \text{and} \quad \hat{\theta}_i(t) = \hat{\theta}(\hat{t}_i(t)),$$

which is the empirical estimator of the sufficient statistic. Like UCB, the algorithm plays $A_t = t$ for $t \in [K]$ and subsequently $A_t = \operatorname{argmax}_i U_i(t)$, where

$$U_i(t) = \sup \left\{ \mu(\theta) : d(\hat{\theta}_i(t-1), \theta) \leq \frac{\log(f(T_i(t-1))f(t))}{T_i(t-1)} \right\}$$

and ties in the argmax are broken by choosing the arm with the largest number of plays.

Part I: Concentration

Given a fixed $\theta \in \Theta$ and independent random variables X_1, \dots, X_n sampled from \mathbb{P}_θ and $\hat{t}_s = \frac{1}{s} \sum_{u=1}^s T(X_u)$ and $\hat{\theta}_s = \hat{\theta}(\hat{t}_s)$. Let $\tilde{\theta} \in \Theta$ be such that $\bar{d}(\tilde{\theta}, \theta) = \varepsilon > 0$. Then

$$\mathbb{P} \left(\bar{d}(\hat{\theta}_s, \theta) \geq \varepsilon \right) \leq \mathbb{P} \left(\hat{t} \geq t(\tilde{\theta}) \right) \leq \exp(-sd(\tilde{\theta}, \theta)) = \exp(-s\varepsilon). \quad (.33)$$

Using an identical argument,

$$\mathbb{P} \left(\underline{d}(\hat{\theta}_s, \theta) \geq \varepsilon \right) \leq \exp(-s\varepsilon). \quad (.34)$$

Define random variable τ by

$$\tau = \min \left\{ t : d(\hat{\theta}_s, \theta) < \frac{\log(f(s)f(t))}{s} \text{ for all } s \in [n] \right\}.$$

Then the expectation of τ is bounded by

$$\begin{aligned} \mathbb{E}[\tau] &= \sum_{t=1}^n \mathbb{P}(\tau \geq t) \\ &\leq \sum_{t=1}^n \sum_{s=1}^n \mathbb{P} \left(d(\hat{\theta}_s, \theta + \varepsilon) \geq \frac{\log(f(s)f(t))}{s} \right) \\ &\leq \sum_{t=1}^n \sum_{s=1}^n \frac{2}{f(s)f(t)} \\ &= O(1), \end{aligned} \quad (.35)$$

where the final equality follows from Eqs. (.33) and (.34) and the same calculation as in the proof of Lemma 10.3. Next let

$$\kappa = \min \{ s \geq 1 : |\hat{\theta}_u - \theta| < \varepsilon \text{ for all } u \geq s \}.$$

The expectation of κ is easily bounded Eq. (.33) and Eq. (.34):

$$\mathbb{E}[\kappa] \leq \sum_{s=1}^n \sum_{u=s}^{\infty} (\exp(-ud(\theta + \varepsilon, \theta)) + \exp(-ud(\theta - \varepsilon, \theta))) = O(1), \quad (.36)$$

where we used the fact that \mathcal{M} is non-singular to ensure strict positivity of the divergences.

Part II: Bounding $\mathbb{E}[T_i(n)]$

Choose $\varepsilon > 0$ sufficiently small that for all suboptimal arms i ,

$$\sup_{\phi \in [\theta_i - \varepsilon, \theta_i + \varepsilon]} \mu(\phi) < \mu^*$$

and define

$$d_{i,\min}(\varepsilon) = \min_{\phi \in \Theta} \{d(\theta_i + x, \phi) : \mu(\phi) = \mu^*, x = \pm\varepsilon\}$$

$$d_{i,\inf}(\varepsilon) = \inf_{\phi \in \Theta} \{d(\theta_i + x, \phi) : \mu(\phi) > \mu^*, x = \pm\varepsilon\}.$$

Let $\hat{\theta}_{i_s}$ be the empirical estimate of θ_i based on the first s samples of arm i , which means that $\hat{\theta}_i(t) = \hat{\theta}_{i_{T_i(t)}}$. Let τ be the smallest t such that

$$d(\hat{\theta}_{1_s}, \theta_1) < \frac{\log(f(s)f(t))}{s} \quad \text{for all } s \in [n],$$

which means that $U_1(t) \geq \mu^*$ for all $t \geq \tau$. For suboptimal arms i let κ_i be the random variable

$$\kappa_i = \min\{s : |\hat{\theta}_{i_u} - \theta_i| < \varepsilon \text{ for all } u \geq s\}.$$

Now suppose that $t \geq \tau$ and $T_i(t-1) \geq \kappa_i$ and $A_t = i$. Then $U_i(t) \geq U_1(t) \geq \mu^*$, which implies that

$$d_{i,\min}(\varepsilon) \geq \frac{\log(f(T_i(t-1))f(t))}{T_i(t-1)}.$$

This means that

$$\sum_{i>1} T_i(t) \leq \tau + \sum_{i>1} \left(1 + \max\left\{\kappa_i, \frac{\log(t^4)}{d_{i,\min}(\varepsilon)}\right\}\right). \quad (.37)$$

Then let

$$\Lambda = \max\left\{t : T_1(t-1) \leq \max_{i>1} T_i(t-1)\right\},$$

which by Eq. (.37) and Eq. (.35) and Eq. (.36) satisfies $\mathbb{E}[\Lambda] = O(1)$. Suppose now that $t \geq \Lambda$. Then $T_1(t-1) > \max_{i>1} T_i(t-1)$ and by the definition of the algorithm $A_t = i$ implies that $U_i(t) > \mu^*$ and so

$$T_i(t-1) \leq \frac{\log(T_i(t-1)^2 f(t))}{d_{i,\inf}(\varepsilon)}.$$

Hence

$$T_i(n) \leq 1 + \Lambda + \frac{\log(f(T_i(n))f(t))}{d_{i,\inf}(\varepsilon)}.$$

Since $\mathbb{E}[\Lambda] = O(1)$ we conclude that

$$\limsup_{n \rightarrow \infty} \frac{\mathbb{E}[T_i(n)]}{\log(n)} \leq \frac{1}{d_{i,\inf}(\varepsilon)}.$$

The result because $\lim_{\varepsilon \rightarrow 0} d_{i,\text{inf}}(\varepsilon) = 0$ follows from the fundamental regret decomposition (Lemma 4.2).



In the analysis of KL-UCB for canonical exponential families the asymptotic rate is a good indicator of the finite-time regret in the sense that the $o(\log(n))$ term hidden by the asymptotics has roughly the same leading constant as the dominant term. By contrast, the analysis here indicates that

$$\mathbb{E}[T_i(n)] \approx \frac{\log(n)}{d_{i,\text{inf}}} + \frac{1}{d_{i,\text{min}}},$$

where $d_{i,\text{min}} = d_{i,\text{min}}(0)$. Although the latter term is negligible asymptotically, it may be the dominant term for all reasonable n .

11.2 Let π be a deterministic policy. Therefore A_t is a function of $x_{1A_1}, \dots, x_{t-1, A_{t-1}}$. We define $\nu = (x_{ti})$ inductively by

$$x_{ti} = \begin{cases} 0 & \text{if } A_t = i \\ 1 & \text{otherwise.} \end{cases}$$

Clearly the policy will collect zero reward and yet

$$\max_{i \in [K]} \sum_{t=1}^n x_{ti} \geq \frac{1}{K} \sum_{t=1}^n \sum_{i=1}^K x_{ti} = \frac{n(K-1)}{K}.$$

Therefore the regret is at least $n(K-1)/K = n(1-1/K)$ as required.

11.5

- (a) Clearly the result is true for $t = 1$ when $P_{t1} = P_{t2} = 1/2$. Assume that $P_{t2} = q_{T_2(t)}(x)$ and consider $P_{t+1,2}$. By assumption $t \leq n/2$, which means that $\hat{Y}_{t1} = 0$ and $\hat{Y}_{t2} = xA_{t2}/P_{t2}$. Hence if $A_t = 1$, then $P_{t+1,2} = P_{t2} = q_{T_2(t)}(x) = q_{T_2(t+1)}(x)$. On the other hand, if $A_t = 2$ then

$$P_{t+1,2} = \frac{P_{t2} \exp(-\eta x/P_{t2})}{1 - P_{t2} + \exp(-\eta x/P_{t2})}.$$

And the result follows by induction.

- (b) Use induction over t and basic calculus.
(c) First prove that $q_3(1/2) \leq q_4(1/4)$, which is an elementary (and annoying) calculus exercise. Conclude that $q_{u-1}(1/2) \leq q_u(1/4)$ for all $u \geq 3$. Since $s \geq 4$ and using its definition we have $q_s(1/2) \leq 1/(8n)$ and $q_s(1/4) \geq q_{s-1}(1/2) \geq 1/(8n)$. Now use the intermediate value theorem and the continuity of $q_s(\cdot)$.
(d) Since $q_u(x) \leq 1/2$ for all u and $x \in [1/4, 1/2]$ it follows that

$$q_{u+1}(x) \leq \exp(-\eta x/q_u(x)).$$

And hence $q_{s-1}(x) \geq \frac{\eta}{4} / \log(8n)$. Furthermore, $q_{u+v}(x) \leq \exp(-\eta xv/q_u(x))$

for any $v \in \mathbb{N}^+$. Let $U_\varepsilon = \{u \in \mathbb{N} : q_u(x) \in [\varepsilon, 2\varepsilon]\}$. Then the previous claim shows that

$$|U_\varepsilon| \leq 1 + \frac{2\varepsilon}{\eta x} \log\left(\frac{1}{\varepsilon}\right) \leq 1 + \frac{8\varepsilon}{\eta} \log\left(\frac{1}{\varepsilon}\right).$$

Let $\varepsilon_k = 2^{-k-1}$ and $k_{\max} = \min\{k : \varepsilon_k \leq q_{s-1}(x)\}$. Then

$$\sum_{u=0}^{s-1} \frac{1}{q_u(s)} \leq \sum_{k=1}^{k_{\max}} \frac{1}{\varepsilon_k} |U_\varepsilon| \leq \sum_{k=1}^{k_{\max}} \frac{1}{\varepsilon_k} \left(1 + \frac{8\varepsilon_k}{\eta} \log\left(\frac{1}{\varepsilon_k}\right)\right).$$

Straightforward calculation now shows there exists a universal constant $c > 0$ such that for large enough n ,

$$\sum_{u=0}^{s-1} \frac{1}{q_u(s)} \leq \frac{c \log(n)^2}{\eta} \leq \frac{n}{8},$$

where the last inequality follows by the assumption that $\eta \geq n^{-p}$ for some $p \in (0, 1)$.

- (e) Let $\tau_0 = 0$ and $\tau_u = \min\{t : N(t) = u\}$, which is an almost surely bounded stopping time. Then by assumption $\tau_u - \tau_{u-1}$ is geometrically distributed with success probability $q_{u-1}(x)$. Then using Markov's inequality,

$$\begin{aligned} \mathbb{P}\left(N\left(2 \sum_{u=0}^{s-1}\right) \geq s\right) &= \mathbb{P}\left(\tau_s \leq 2 \sum_{u=0}^{s-1} 1/q_u(x)\right) \\ &= \mathbb{P}\left(\sum_{u=1}^s \tau_u - \tau_{u-1} \leq 2 \sum_{u=0}^{s-1} 1/q_u(x)\right) \\ &= 1 - \mathbb{P}\left(\sum_{u=1}^s \tau_u - \tau_{u-1} > 2 \sum_{u=0}^{s-1} 1/q_u(x)\right) \\ &\geq 1 - \frac{\mathbb{E}\left[\sum_{u=1}^s \tau_u - \tau_{u-1}\right]}{2 \sum_{u=0}^{s-1} 1/q_u(x)} \\ &= \frac{1}{2}. \end{aligned}$$

- (f) By part (d), $\sum_{u=0}^{s-1} 1/q_u(x) \leq n/8$. The result follows by noting that $T_2(t)$ is a counting process with $T_2(1) = 0$ and $T_2(t+1) - T_2(t) \in \{0, 1\}$ and $\mathbb{P}(T_2(t+1) - T_2(t) | T_2(t)) = q_{T_2(t)}(x)$. Then apply part (e).
- (g) The definition of P_t means that on the event E ,

$$P_{t2} \leq \exp\left(\eta \left(\sum_{s=n/2+1}^t \hat{Y}_{s1} - 2n\right)\right).$$

This ensures that $P_{t1} = 1 - P_{t2} \geq 1/2$ as long as $\sum_{s=n/2+1}^t \hat{Y}_{s1} \geq 2n$ and consequentially $\hat{Y}_{t1} \leq 2$ for all $t \in [n]$. Using the above display again shows that $P_{t2} \leq \exp(-\eta n)$ and a union bound completes the step.

- (h) Suppose that $A_t = 2$ and $T_2(t-1) = s$, then $\hat{Y}_{t2} = x/P_{t2} = 8nx \geq 2n$. Hence $\mathbb{P}(E) \geq \mathbb{P}(T_t(n/2) \geq s+1)$. The result is completed by noticing that

$$\begin{aligned} \mathbb{P}(T_t(n/2) \geq s+1) &\geq \mathbb{P}(T_t(n/4) \geq s) \left(1 - \left(1 - \frac{1}{8n}\right)^{n/4}\right) \\ &\geq \frac{1}{2} \left(1 - \left(1 - \frac{1}{8n}\right)^{n/4}\right) \\ &\geq \frac{1}{2} (1 - \exp(-1/32)), \end{aligned}$$

where the first inequality follows from the definition of s and the fact that $p_s(x) = 1/(8n)$. The second inequality is (f). The last follows since $(1-a)^b = \exp(b \log(1-a)) \leq \exp(-ab)$.

- (i) On the event that $A_t = 1$ for all $t > n/2$ the regret satisfies $\hat{R}_n \geq n/2 - nx/2 \geq n/4$. The result follows by parts (g) and (h).
- (j) The regret \hat{R}_n can be negative and negative linear regret can cancel the negative positive regret. In fact, the expected regret for this problem appears to be negative.
- (k) The value of $x = 0.49979435\dots$ suffices.

11.6 The result follows from a long sequence of equalities:

$$\begin{aligned} \mathbb{P}\left(\log a_i + G_i \geq \max_{k \in [K]} \log a_k + G_k\right) &= \mathbb{E}\left[\prod_{k \neq i} \mathbb{P}(\log a_k + G_k \leq \log a_i + G_i)\right] \\ &= \mathbb{E}\left[\prod_{k \neq i} \exp\left(-\frac{a_k}{a_i} \exp(-G_i)\right)\right] \\ &= \mathbb{E}\left[U_i^{\sum_{k \neq i} \frac{a_k}{a_i}}\right] \\ &= \frac{1}{1 + \sum_{k \neq i} \frac{a_k}{a_i}} \\ &= \frac{a_i}{\sum_{k=1}^K a_k}. \end{aligned}$$

12.2 We proceed in five steps.

Step 1: Decomposition

Using that $\sum_a P_{ta} = 1$ and some algebra we get

$$\begin{aligned} & \sum_{t=1}^n \sum_{a=1}^K P_{ta} (Z_{ta} - Z_{tA^*}) \\ &= \underbrace{\sum_{t=1}^n \sum_{a=1}^K P_{ta} (\tilde{Z}_{ta} - \tilde{Z}_{tA^*})}_{(A)} + \underbrace{\sum_{t=1}^n \sum_{a=1}^K P_{ta} (Z_{ta} - \tilde{Z}_{ta})}_{(B)} + \underbrace{\sum_{t=1}^n (\tilde{Z}_{tA^*} - Z_{tA^*})}_{(C)}. \end{aligned}$$

Step 2: Bounding (A)

By assumption (c) we have $\beta_{ta} \geq 0$, which by assumption (a) means that $\eta \tilde{Z}_{ta} \leq \eta \hat{Z}_{ta} \leq \eta |\hat{Z}_{ta}| \leq 1$ for all a . A straightforward modification of the analysis in the last chapter shows that (A) is bounded by

$$\begin{aligned} (A) &\leq \frac{\log(K)}{\eta} + \eta \sum_{t=1}^n \sum_{a=1}^K P_{ta} \tilde{Z}_{ta}^2 \\ &= \frac{\log(K)}{\eta} + \eta \sum_{t=1}^n \sum_{a=1}^K P_{ta} (\hat{Z}_{ta}^2 + \beta_{ta}^2) - 2\eta \sum_{t=1}^n \sum_{a=1}^K P_{ta} \hat{Z}_{ta} \beta_{ta} \\ &\leq \frac{\log(K)}{\eta} + \eta \sum_{t=1}^n \sum_{a=1}^K P_{ta} \hat{Z}_{ta}^2 + 3 \sum_{t=1}^n \sum_{a=1}^K P_{ta} \beta_{ta}, \end{aligned}$$

where in the last two line we used the assumptions that $\eta \beta_{ta} \leq 1$ and $\eta |\hat{Z}_{ta}| \leq 1$.

Step 3: Bounding (B)

For (B) we have

$$(B) = \sum_{t=1}^n \sum_{a=1}^K P_{ta} (Z_{ta} - \tilde{Z}_{ta}) = \sum_{t=1}^n \sum_{a=1}^K P_{ta} (Z_{ta} - \hat{Z}_{ta} + \beta_{ta}).$$

We prepare to use Exercise 5.17. By assumptions (c) and (d) respectively we have $\eta \mathbb{E}_{t-1}[\hat{Z}_{ta}^2] \leq \beta_{ta}$ and $\mathbb{E}_{t-1}[\hat{Z}_{ta}] = Z_{ta}$. By Jensen's inequality,

$$\eta \mathbb{E}_{t-1} \left[\left(\sum_{a=1}^K P_{ta} (Z_{ta} - \hat{Z}_{ta}) \right)^2 \right] \leq \eta \sum_{a=1}^K P_{ta} \mathbb{E}_{t-1}[\hat{Z}_{ta}^2] \leq \sum_{a=1}^K P_{ta} \beta_{ta}.$$

Therefore by Exercise 5.17, with probability at least $1 - \delta$

$$(B) \leq 2 \sum_{t=1}^n \sum_{a=1}^K P_{ta} \beta_{ta} + \frac{\log(1/\delta)}{\eta}.$$

Step 4: Bounding (C)

For (C) we have

$$(C) = \sum_{t=1}^n (\tilde{Z}_{tA^*} - Z_{tA^*}) = \sum_{t=1}^n \left(\hat{Z}_{tA^*} - Z_{tA^*} - \beta_{tA^*} \right).$$

Because A^* is random we cannot directly apply Exercise 5.17, but need a union bound over all actions. Let a be fixed. Then by Exercise 5.17 and the assumption that $\eta|\hat{Z}_{ta}| \leq 1$ and $\mathbb{E}_{t-1}[\hat{Z}_{ta}] = Z_{ta}$ and $\eta\mathbb{E}_{t-1}[\hat{Z}_{ta}^2] \leq \beta_{ta}$, with probability at least $1 - \delta$.

$$\sum_{t=1}^n \left(\hat{Z}_{ta} - Z_{ta} - \beta_{ta} \right) \leq \frac{\log(1/\delta)}{\eta}.$$

Therefore by a union bound we have with probability at most $1 - K\delta$,

$$(C) \leq \frac{\log(1/\delta)}{\eta}.$$

Step 5: Putting it together

Combining the bounds on (A), (B) and (C) in the last three steps with the decomposition in the first step shows that with probability at least $1 - (K + 1)\delta$,

$$R_n \leq \frac{3\log(1/\delta)}{\eta} + \eta \sum_{t=1}^n \sum_{a=1}^K P_{ta} \hat{Z}_{ta}^2 + 5 \sum_{t=1}^n \sum_{a=1}^K P_{ta} \beta_{ta}.$$

where we used the assumption that $\delta \leq 1/K$.

14.1 Let $\mu = P - Q$, which is a signed measure on (Ω, \mathcal{F}) . By the Hahn decomposition theorem there exist disjoint sets $A, B \subset \Omega$ such that $A \cup B = \Omega$ and $\mu(E) \geq 0$ for all measurable $E \subseteq A$ and $\mu(E) \leq 0$ for all measurable $E \subseteq B$. Then

$$\begin{aligned} \int_{\Omega} X dP - \int_{\Omega} X dQ &= \int_A X d\mu + \int_B X d\mu \\ &\leq b\mu(A) + a\mu(B) \\ &= (b - a)\mu(A) \\ &\leq (b - a)\delta(P, Q), \end{aligned}$$

where we used the fact that $\mu(B) = P(B) - Q(B) = Q(A) - P(A) = -\mu(A)$.

14.7 Dobrushin's theorem says that for any ω ,

$$D(P(\cdot | \omega), Q(\cdot | \omega)) = \sup_{(A_i)} \sum_i P(A_i | \omega) \log \left(\frac{P(A_i | \omega)}{Q(A_i | \omega)} \right),$$

where the supremum is taken over all finite partitions of \mathbb{R} with rational-valued end-points. By the definition of a Markov kernel it follows that the quantity inside the supremum on the right-hand side is \mathcal{F} -measurable as a function of ω

for any finite partition. Since the supremum is over a countable set, the whole right-hand side is \mathcal{F} -measurable as required.

14.8 First assume that $P \ll Q$. Then let P^t and Q^t be the restrictions of P and Q to $(\mathbb{R}^t, \mathfrak{B}(\mathbb{R}^t))$ given by

$$P^t(A) = P(A \times \Omega^{n-t}) \quad \text{and} \quad Q^t(A) = Q(A \times \Omega^{n-t}).$$

You should check that $P \ll Q$ implies that $P^t \ll Q^t$ and hence there exists a Radon-Nikodym derivative dP^t/dQ^t . Define

$$F(x_t | x_1, \dots, x_{t-1}) = \frac{dP^t}{dQ^t}(x_1, \dots, x_t) \bigg/ \frac{dP^{t-1}}{dQ^{t-1}}(x_1, \dots, x_{t-1}),$$

which is well defined for all $x_1, \dots, x_{t-1} \in \mathbb{R}^{t-1}$ except for a set of P^{t-1} -measure zero. Then for any $A \in \mathfrak{B}(\mathbb{R}^{t-1})$ and $B \in \mathfrak{B}(\mathbb{R})$,

$$\begin{aligned} \int_A \int_B F(x_t | \omega) Q_t(dx_t | \omega) P^{t-1}(d\omega) &= \int_A \int_B \frac{dP^t}{dQ^t}(x_t, \omega) Q_t(dx_t | \omega) Q^{t-1}(d\omega) \\ &= \int_{A \times B} \frac{dP^t}{dQ^t} dQ^t \\ &= P(A \times B). \end{aligned}$$

A monotone class argument shows that $F(x_t | \omega)$ is P^{t-1} -almost surely the Radon-Nikodym derivative of $P_t(\cdot | \omega)$ with respect to $Q_t(\cdot | \omega)$. Hence

$$\begin{aligned} D(P, Q) &= \mathbb{E}_P \left[\log \left(\frac{dP}{dQ} \right) \right] \\ &= \sum_{t=1}^n \mathbb{E}_P [\log (F(X_t | X_1, \dots, X_{t-1}))] \\ &= \sum_{t=1}^n \mathbb{E}_P [D(P_t(\cdot | X_1, \dots, X_{t-1}), Q_t(\cdot | X_1, \dots, X_{t-1}))]. \end{aligned}$$

Now suppose that $P \not\ll Q$. Then by definition $D(P, Q) = \infty$. We need to show this implies there exists a $t \in [n]$ such that $D(P_t(\cdot | \omega), Q_t(\cdot | \omega)) = \infty$ with nonzero probability. Proving the contrapositive, let

$$U_t = \{\omega : D(P_t(\cdot | \omega), Q_t(\cdot | \omega)) < \infty\}$$

and assume that $P(U_t) = 1$ for all t . Then $U = \bigcap_{t=1}^n U_t$ satisfies $P(U) = 1$. On U_t let $F(x_t | x_1, \dots, x_{t-1}) = dP_t(\cdot | x_1, \dots, x_{t-1})/dQ_t(\cdot | x_1, \dots, x_{t-1})(x_t)$ and otherwise let $F(x_t | x_1, \dots, x_{t-1}) = 0$. Iterating applications of Fubini's theorem shows that for any $(A_t)_{t=1}^n$ with $A_t \in \mathfrak{B}(\mathbb{R})$ it holds that

$$\int_{A_1 \times \dots \times A_n} \prod_{t=1}^n F(x_t | x_1, \dots, x_{t-1}) Q(dx_1, \dots, dx_n) = P(A_1 \times \dots \times A_n).$$

Hence $\prod_{t=1}^n F(x_t | x_1, \dots, x_{t-1})$ behaves like the Radon-Nikodym derivative of

P with respect to Q on rectangles. Another monotone class argument extends this to all measurable sets and the existence of dP/dQ guarantees that $P \ll Q$.

18.1

(a) By Jensen's inequality,

$$\begin{aligned} \sum_{c \in \mathcal{C}} \sqrt{\sum_{t=1}^n \mathbb{I}\{c_t = c\}} &= |\mathcal{C}| \sum_{c \in \mathcal{C}} \frac{1}{|\mathcal{C}|} \sqrt{\sum_{t=1}^n \mathbb{I}\{c_t = c\}} \\ &\leq |\mathcal{C}| \sqrt{\sum_{c \in \mathcal{C}} \frac{1}{|\mathcal{C}|} \sum_{t=1}^n \mathbb{I}\{c_t = c\}} \\ &= \sqrt{|\mathcal{C}|n}, \end{aligned}$$

where the inequality follows from Jensen's inequality and the concavity of $\sqrt{\cdot}$ and the last equality follows since $\sum_{c \in \mathcal{C}} \sum_{t=1}^n \mathbb{I}\{c_t = c\} = n$.

(b) When each context occurs $n/|\mathcal{C}|$ times we have

$$\sum_{c \in \mathcal{C}} \sqrt{\sum_{t=1}^n \mathbb{I}\{c_t = c\}} = \sqrt{n|\mathcal{C}|},$$

which matches the upper bound proven in the previous part.

18.6

- (a) This follows because $X_i \leq \max_j X_j$ for any family of random variables $(X_i)_i$. Hence $\mathbb{E}[X_i] \leq \mathbb{E}[\max_j X_j]$.
- (b) Modify the proof by proving a bound on

$$\mathbb{E} \left[\sum_{t=1}^n E_{m^*}^{(t)} x_t - \sum_{t=1}^n X_t \right]$$

for an arbitrarily fixed m^* . By the definition of learning experts, $\mathbb{E}_t[\hat{X}_t] = x_t$ and so Eq. (18.10) also remains valid. Note this would not be true in general if $E^{(t)}$ were allowed to depend on A_t . The rest follows the same way as in the oblivious case.

18.7 The inequality $E_n^* \leq nK$ is trivial (since $\max_m E_{mi}^{(t)} \leq 1$). To prove $E_n^* \leq nM$, let $m_{t,i}^* = \operatorname{argmax}_m E_{mi}^{(t)}$. Then, $E_n^* = \sum_t \sum_i \sum_m E_{m,i}^{(t)} \mathbb{I}\{m = m_{t,i}^*\} \leq \sum_t \sum_m \sum_i E_{m,i}^{(t)} = nM$, where the last step used that $\sum_i E_{m,i}^{(t)} = 1$.

19.2 Let \mathcal{T} be the set of rounds t when $\|x_t\|_{V_{t-1}^{-1}} \geq 1$ and $G_t = V_0 +$

$\sum_{s=1}^t \mathbb{I}_{\mathcal{T}}(s) x_t x_t^\top$. Then

$$\begin{aligned} \left(\frac{d\lambda + |\mathcal{T}|L^2}{d} \right)^d &\geq \left(\frac{\text{trace}(G_n)}{d} \right)^d \\ &\geq \det(G_n) \\ &= \det(V_0) \prod_{t \in \mathcal{T}} (1 + \|x_t\|_{G_{t-1}^{-1}}^2) \\ &\geq \det(V_0) \prod_{t \in \mathcal{T}} (1 + \|x_t\|_{V_{t-1}^{-1}}^2) \\ &\geq \lambda^d 2^{|\mathcal{T}|}. \end{aligned}$$

Rearranging and taking the logarithm shows that

$$|\mathcal{T}| \leq \frac{d}{\log(2)} \log \left(1 + \frac{|\mathcal{T}|L^2}{d\lambda} \right).$$

Abbreviate $a = d/\log(2)$ and $b = L^2/d\lambda$, which are both positive. Then

$$a \log(1 + b(3a \log(1 + ab))) \leq a \log(1 + 3a^2 b^2) \leq a \log(1 + ab)^3 = 3a \log(1 + ab).$$

Since $x - a \log(1 + bx)$ is decreasing for $x \geq 3a \log(1 + ab)$ it follows that

$$|\mathcal{T}| \leq 3a \log(1 + ab) = \frac{3d}{\log(2)} \log \left(1 + \frac{L^2}{\lambda \log(2)} \right).$$

20.1

- (a) If $C \subset \mathcal{A}$ is an ε -covering then it is also an ε' -covering with any $\varepsilon' \geq \varepsilon$. Hence, $\varepsilon \rightarrow N(\varepsilon)$ is a decreasing function of ε .
- (b) The inequality $M(2\varepsilon) \leq N(\varepsilon)$ amounts to showing that any 2ε packing has a cardinality at most the cardinality of any ε covering. Assume this does not hold, that is, there is a 2ε packing $P \subset \mathcal{A}$ and an ε -covering $C \subset \mathcal{A}$ such that $|P| \geq |C| + 1$. By the pigeonhole principle, there is $c \in C$ such that there are distinct $x, y \in P$ such that $x, y \in B(c, \varepsilon)$. Then $\|x - y\| \leq \|x - c\| + \|c - y\| \leq 2\varepsilon$, which contradicts that P is a 2ε -packing.

If $M(\varepsilon) = \infty$, the inequality $N(\varepsilon) \leq M(\varepsilon)$ is trivially true. Otherwise take a maximum ε -packing P of \mathcal{A} . This packing is automatically an ε -covering as well (otherwise P would not be a maximum packing), hence, the result.

- (c) We show the inequalities going left to right. For the first inequality, if $N \doteq N(\varepsilon) = \infty$ then there is nothing to be shown. Otherwise let C be a minimum cardinality ε -cover of \mathcal{A} . Then from the definition of cover and the additivity of volume, $\text{vol}(\mathcal{A}) \leq \sum_{x \in \mathcal{A}'} \text{vol}(B(x, \varepsilon)) = N\varepsilon^d \text{vol}(B)$. Reordering gives the inequality.

The next inequality, namely that $N(\varepsilon) \leq M(\varepsilon)$ has already been shown.

Consider now the inequality bounding $M \doteq M(\varepsilon)$. Let P be a maximum cardinality ε -packing of \mathcal{A} . Then, for any $x, y \in P$ distinct, $B(x, \varepsilon/2) \cap B(y, \varepsilon/2) = \emptyset$. Further, for $x \in P$, $B(x, \varepsilon/2) \subset \mathcal{A} + \frac{\varepsilon}{2}B$ and thus

$\cup_{x \in P} B(x, \varepsilon/2) \subset \mathcal{A} + \frac{\varepsilon}{2}B$, hence, by the additivity of volume, $M \text{vol}(\frac{\varepsilon}{2}B) \leq \text{vol}(\mathcal{A} + \frac{\varepsilon}{2}B)$.

For the next inequality note that $\varepsilon B \subset \mathcal{A}$ immediately implies that $\mathcal{A} + \frac{\varepsilon}{2}B \subset \mathcal{A} + \frac{1}{2}\mathcal{A}$ (check the containment using the definitions), while the convexity of \mathcal{A} implies that $\mathcal{A} + \frac{1}{2}\mathcal{A} \subset \frac{3}{2}\mathcal{A}$. For this second claim let $u \in \mathcal{A} + \frac{1}{2}\mathcal{A}$. Then $u = x + \frac{1}{2}y$ for some $x, y \in \mathcal{A}$. By the convexity of \mathcal{A} , $\frac{2}{3}u = \frac{2}{3}x + \frac{1}{3}y \in \mathcal{A}$ and hence $u = \frac{3}{2}(\frac{2}{3}u) \in \frac{3}{2}\mathcal{A}$. For the final inequality note that for measurable X and $c > 0$ we have $\text{vol}(cX) = c^d \text{vol}(X)$. This is true because cX is the image of X under the linear mapping represented by a diagonal matrix with c on the diagonal and this matrix has determinant c^d .

- (d) Let \mathcal{A} be bounded, and say, $\mathcal{A} \subset rB$ for some $r > 0$. Then $\text{vol}(\mathcal{A} + \varepsilon/2B) \leq \text{vol}(rB + \varepsilon/2B) = \text{vol}((r + \varepsilon/2)B) < +\infty$, hence, the previous part gives that $N(\varepsilon) \leq M(\varepsilon) < +\infty$. Now assume that $N(\varepsilon) < \infty$ and let C be a minimum cover of \mathcal{A} . Then $\mathcal{A} \subset \cup_{x \in C} B(x, \varepsilon) \subset \cup_{x \in C} (\|x\| + \varepsilon)B \subset \max_{x \in C} (\|x\| + \varepsilon)B$ hence, \mathcal{A} is bounded.

20.4 Proving that \bar{M}_t is \mathcal{F}_t -measurable is actually not trivial. It follows because $M_t(\cdot)$ is measurable and by the ‘sections’ lemma [Kallenberg, 2002, Lem 1.26]. It remains to show that $\mathbb{E}[\bar{M}_t | \mathcal{F}_{t-1}] \leq \bar{M}_{t-1}$ almost surely. Proceeding by contradiction, suppose that $\mathbb{P}(\mathbb{E}[\bar{M}_t | \mathcal{F}_{t-1}] - \bar{M}_{t-1} > 0) > 0$. Then there exists an $\varepsilon > 0$ such that the set $A = \{\omega : \mathbb{E}[\bar{M}_t | \mathcal{F}_{t-1}](\omega) - \bar{M}_{t-1}(\omega) > \varepsilon\} \in \mathcal{F}_{t-1}$ satisfies $\mathbb{P}(A) > 0$. Then

$$\begin{aligned} 0 < \int_A (\mathbb{E}[\bar{M}_t | \mathcal{F}_{t-1}] - \bar{M}_{t-1}) d\mathbb{P} &= \int_A (\bar{M}_t - \bar{M}_{t-1}) d\mathbb{P} \\ &= \int_A \int_{\mathbb{R}^d} (M_t(x) - M_{t-1}(x)) dh(x) d\mathbb{P} \\ &= \int_{\mathbb{R}^d} \int_A (M_t(x) - M_{t-1}(x)) d\mathbb{P} dh(x) \\ &\leq 0, \end{aligned}$$

where the first equality follows from the definition of conditional expectation, the second by substituting the definition of \bar{M}_t and the third from Fubini-Tonelli’s theorem. The last follows from Lemma 20.2 and the definition of conditional expectation again. The proof is completed by noting the deep result that $0 \not\leq 0$. In this proof it is necessary to be careful to avoid integrating over conditional $\mathbb{E}[M_t(x) | \mathcal{F}_{t-1}]$, which are only defined for each x almost surely and need not be measurable as a function of x .

20.7 We first show a bound on the right tail of S_t . A symmetric argument suffices for the left tail. Let $Y_s = X_s - \mu_s |X_s|$ and $M_t(\lambda) = \exp(\sum_{s=1}^t (\lambda Y_s - \lambda^2 |X_s|/2))$. Define filtration $\mathcal{G}_1 \subset \dots \subset \mathcal{G}_n$ by $\mathcal{G}_t = \sigma(\mathcal{F}_{t-1}, |X_t|)$. Using the fact that $X_s \in \{-1, 0, 1\}$ we have for any $\lambda > 0$ that

$$\mathbb{E}[\exp(\lambda Y_s - \lambda^2 |X_s|/2) | \mathcal{G}_s] \leq 1.$$

Therefore $M_t(\lambda)$ is a supermartingale for any $\lambda > 0$. The next step is to use the method of mixtures with a uniform distribution on $[0, 2]$. Let $M_t = \int_0^2 M_t(\lambda) d\lambda$. Then Markov's inequality shows that for any \mathcal{G}_t -measurable stopping time τ with $\tau \leq n$ almost surely, stopping time $\mathbb{P}(M_\tau \geq 1/\delta) \leq \delta$. Next we need a bound on M_τ . The following holds whenever $S_t \geq 0$.

$$\begin{aligned} M_t &= \frac{1}{2} \int_0^2 M_t(\lambda) d\lambda \\ &= \frac{1}{2} \sqrt{\frac{\pi}{2N_t}} \left(\operatorname{erf} \left(\frac{S_t}{\sqrt{2N_t}} \right) + \operatorname{erf} \left(\frac{2N_t - S_t}{\sqrt{2N_t}} \right) \right) \exp \left(\frac{S_t^2}{2N_t} \right) \\ &\geq \frac{\operatorname{erf}(\sqrt{2})}{2} \sqrt{\frac{\pi}{2N_t}} \exp \left(\frac{S_t^2}{2N_t} \right). \end{aligned}$$

The bound on the upper tail completed via a stopping time, which shows that

$$\mathbb{P} \left(\text{exists } t \leq n : S_t \geq \sqrt{2N_t \log \left(\frac{2}{\delta \operatorname{erf}(\sqrt{2})} \sqrt{\frac{2N_t}{\pi}} \right)} \text{ and } N_t > 0 \right) \leq \delta.$$

The result follows by symmetry and union bound.

20.8 Following the hint, we show that $\exp(L_t(\theta_*))$ is a martingale. Indeed, letting $\mathcal{F}_t = \sigma(X_1, \dots, X_t)$,

$$\begin{aligned} \mathbb{E}[\exp(L_t(\theta_*)) | \mathcal{F}_{t-1}] &= \mathbb{E} \left[p_{\hat{\theta}_{t-1}}(X_t) / p_\theta(X_t) \right] \exp(L_{t-1}(\theta_*)) \\ &= \exp(L_{t-1}(\theta_*)) \int p_\theta(x) \frac{p_{\hat{\theta}_{t-1}}(x)}{p_\theta(x)} d\mu(x) \\ &= \exp(L_{t-1}(\theta_*)) \int p_{\hat{\theta}_{t-1}}(x) d\mu(x) \\ &= \exp(L_{t-1}(\theta_*)). \end{aligned}$$

Then, applying the Cramer-Chernoff trick,

$$\begin{aligned} \mathbb{P}(L_t(\theta_*) \geq \log(1/\delta) \text{ for some } t \geq 1) &= \mathbb{P} \left(\sup_{t \in \mathbb{N}} \exp(L_t(\theta_*)) \geq 1/\delta \right) \\ &\leq \delta \mathbb{E}[\exp(L_0(\theta_*))] = \delta, \end{aligned}$$

where the inequality is due to Theorem 3.5, the maximal inequality of nonnegative supermartingales.

21.1 Following the hint,

$$\nabla f(\pi)_a = \frac{\operatorname{trace}(\operatorname{adj}(V(\pi)) a a^\top)}{\det(V(\pi))} = \frac{a^\top \operatorname{adj}(V(\pi)) a}{\det(V(\pi))} = a^\top V(\pi)^{-1} a = \|a\|_{V(\pi)^{-1}}^2,$$

where in the third equality we used that $\operatorname{adj}(V(\pi))$ is symmetric since $V(\pi)$ is symmetric, hence, following the hint, $\operatorname{adj}(V(\pi)) / \det(V(\pi)) = V(\pi)^{-1}$.

21.2 By the determinant product rule,

$$\begin{aligned}\log \det(H + tZ) &= \log \det(H^{1/2}(I + tH^{-1/2}ZH^{-1/2})H^{1/2}) \\ &= \log \det(H) + \log \det(I + tH^{-1/2}ZH^{-1/2}) \\ &= \log \det(H) + \sum_i \log(1 + t\lambda_i),\end{aligned}$$

where λ_i are the eigenvalues of $H^{-1/2}ZH^{-1/2}$. Since $\log(1 + t\lambda_i)$ is concave, their sum is also concave, proving that $t \mapsto \log \det(H + tZ)$ is concave.

21.3 Let \mathcal{A} be a compact subset of \mathbb{R}^d and $(\mathcal{A}_n)_n$ be a sequence of finite subsets with $\mathcal{A}_n \subset \mathcal{A}_{n+1}$ and with $\lim_{n \rightarrow \infty} d(\mathcal{A}, \mathcal{A}_n) = 0$ where d is the Hausdorff metric. Then let π_n be a G -optimal design for \mathcal{A}_n with support of size at most $d(d+1)/2$. By compactness there (π_n) has a cluster point π^* for which $g(\pi^*) = d$.

23.6

- (a) This follows from straightforward calculus.
 (b) The result is trivial for $\Lambda < 0$. For $\Lambda \geq 0$ we have

$$\begin{aligned}M_n &= \int_{\mathbb{R}} f(\lambda) \exp\left(\lambda S_n - \frac{\lambda^2 n}{2}\right) d\lambda \\ &\geq \int_{\Lambda}^{\Lambda(1+\varepsilon)} f(\lambda) \exp\left(\lambda S_n - \frac{\lambda^2 n}{2}\right) d\lambda \\ &\geq \varepsilon \Lambda f(\Lambda(1+\varepsilon)) \exp\left(\Lambda(1+\varepsilon)S_n - \frac{\Lambda^2(1+\varepsilon)^2 n}{2}\right) \\ &= \varepsilon \Lambda f(\Lambda(1+\varepsilon)) \exp\left(\frac{(1-\varepsilon^2)S_n^2}{2n}\right).\end{aligned}$$

- (c) Let $n \in \mathbb{N}$. Since M_t is a supermartingale with $M_0 = 1$ it follows that

$$P_n = \mathbb{P}(\text{exists } t \leq n : M_t \geq 1/\delta) \leq \delta.$$

Hence $\mathbb{P}(\text{exists } t : M_t \geq 1/\delta) \leq \delta$. Substituting the result from the previous part and rearranging completes the proof.

- (d) A suitable choice of f is $f(\lambda) = \frac{\mathbb{I}\{\lambda \leq e^{-e}\}}{\lambda \log(\frac{1}{\lambda}) (\log \log(\frac{1}{\lambda}))^2}$.
 (e) Let $\varepsilon_n = \min\{1/2, 1/\log \log(n)\}$ and $\delta \in [0, 1]$ be the largest (random) value such that S_n never exceeds

$$\sqrt{\frac{2n}{1-\varepsilon_n^2} \left(\log\left(\frac{1}{\delta}\right) + \log\left(\frac{1}{\varepsilon_n \Lambda_n f(\Lambda_n(1+\varepsilon_n))}\right) \right)}.$$

By Part (c) we have $\mathbb{P}(\delta > 0) = 1$. Furthermore, $\limsup_{n \rightarrow \infty} S_n/n = 0$ almost surely by the strong law of large numbers, so that $\Lambda_n \rightarrow 0$ almost surely. On the intersection of these almost sure events we have

$$\limsup_{n \rightarrow \infty} \frac{S_n}{\sqrt{2n \log \log(n)}} \leq 1.$$

24.1 Assume without loss of generality that $i = 1$ and let $\theta^{(-1)} \in \Theta^{p-1}$. The objective is to prove that

$$\frac{1}{|\Theta|} \sum_{\theta^{(1)} \in \Theta} R_{n1}(\theta) \geq \frac{\sqrt{kn}}{8}.$$

For $j \in [k]$ let $T_j(n) = \sum_{t=1}^n \mathbb{1}\{B_{t1} = j\}$ be the number of times base action j is played in the first bandit. Define $\psi_0 \in \mathbb{R}^d$ to be the vector with $\psi_0^{(-1)} = \theta^{(-1)}$ and $\psi_0^{(1)} = 0$. For $j \in [k]$ let $\psi_j \in \mathbb{R}^d$ be given by $\psi_j^{(-1)} = \theta^{(-1)}$ and $\psi_j^{(1)} = \Delta e_j$. Abbreviate $\mathbb{P}_j = \mathbb{P}_{\psi_j}$ and $\mathbb{E}_j[\cdot] = \mathbb{E}_{\mathbb{P}_j}[\cdot]$. With this notation, we have

$$\frac{1}{|\Theta|} \sum_{\theta^{(1)} \in \Theta} R_{n1}(\theta) = \frac{1}{k} \sum_{j=1}^k \Delta(n - \mathbb{E}_j[T_j(n)]). \quad (.38)$$

Lemma 15.1 gives that

$$D(\mathbb{P}_0, \mathbb{P}_j) = \frac{1}{2} \mathbb{E}_0 \left[\sum_{t=1}^n \langle A_t, \psi_0 - \psi_j \rangle^2 \right] = \frac{\Delta^2}{2} \mathbb{E}_0 [T_j(n)].$$

Choosing $\Delta = \sqrt{k/n}/2$ and applying Pinsker's inequality yields

$$\begin{aligned} \sum_{j=1}^k \mathbb{E}_j[T_j(n)] &\leq \sum_{j=1}^k \mathbb{E}_0[T_j(n)] + n \sum_{j=1}^k \sqrt{\frac{1}{2} D(\mathbb{P}_0, \mathbb{P}_j)} \\ &= n + n \sum_{j=1}^k \sqrt{\frac{\Delta^2}{4} \mathbb{E}_0[T_j(n)]} \\ &\leq n + n \sqrt{\frac{k\Delta^2}{4} \sum_{j=1}^k \mathbb{E}_0[T_j(n)]} \quad (\text{Cauchy-Schwarz}) \\ &= n + n \sqrt{\frac{k\Delta^2 n}{4}} \\ &\leq 3nk/4. \quad (\text{since } k \geq 2) \end{aligned}$$

Combining the above display with Eq. (.38) completes the proof:

$$\frac{1}{|\Theta|} \sum_{\theta^{(1)} \in \Theta} R_{n1}(\theta) = \frac{1}{k} \sum_{j=1}^k \Delta(n - \mathbb{E}_j[T_j(n)]) \geq \frac{n\Delta}{4} = \frac{1}{8} \sqrt{kn}.$$

25.3 For (a) let $\theta_1 = \Delta$ and $\theta_i = 0$ for $i > 1$ and let $\mathcal{A} = \{e_1, \dots, e_{d-1}\}$. Then adding e_d increases the asymptotic regret. For (b) let $\theta_1 = \Delta$ and $\theta_i = 0$ for $1 < i < d$ and $\theta_d = 1$ and $\mathcal{A} = \{e_1, \dots, e_{d-1}\}$. Then for small values of Δ added e_d decreases the asymptotic regret.

26.1 Let \mathbb{P} be the on the space on which X is defined. Following the hint, let $x_0 = \mathbb{E}[X] \in \mathbb{R}^d$. Then let $a \in \mathbb{R}^d$ and $b \in \mathbb{R}$ be such that $\langle a, x_0 \rangle + b = f(x_0)$ and

$\langle a, x \rangle + b \leq f(x)$ for all $x \in \mathbb{R}^d$. Such a hyperplane is guaranteed to exist by the supporting hyperplane theorem. Then

$$\int f(X) d\mathbb{P} \geq \int (\langle a, X \rangle + b) d\mathbb{P} = \langle a, x_0 \rangle + b = f(x_0) = f(\mathbb{E}[X]).$$

An alternative is of course to follow the ideas next to the picture in the main text. As you may recall that proof is given for the case when X is discrete. To extend the proof to the general case, one can use the ‘standard machinery’ of building up the integral from simple functions, but the resulting proof, originally due to [Needham \[1993\]](#), is much longer than what was given above.

26.7

- (a) Fix $u \in \mathbb{R}^d$. By definition $f^*(u) = \sup_x \langle x, u \rangle - f(x)$. To find this value we solve for x where the derivative of $\langle x, u \rangle - f(x)$ in x is equal to zero. As calculated before, $\nabla f(x) = \log(x)$. Thus, we need to find the solution to $u = \log(x)$, giving $x = \exp(u)$. Plugging this value, we get $f^*(u) = \langle \exp(u), u \rangle - f(\exp(u))$. Now, $f(\exp(u)) = \langle \exp(u), \log(\exp(u)) \rangle - \langle \exp(u), \mathbf{1} \rangle = \langle \exp(u), u \rangle - \langle \exp(u), \mathbf{1} \rangle$. Hence, $f^*(u) = \langle \exp(u), \mathbf{1} \rangle$ and $\nabla f^*(u) = \exp(u)$.
- (b) From our calculation, $\text{dom}(\nabla f^*) = \mathbb{R}^d$.
- (c) $D_{f^*}(u, v) = f^*(u) - f^*(v) - \langle \nabla f^*(v), u - v \rangle = \langle \exp(u) - \exp(v), \mathbf{1} \rangle - \langle \exp(v), u - v \rangle$.
- (d) To check Part (a) of Theorem 26.3 note that $\nabla f(x) = \log(x)$ and $\nabla f^*(u) = \exp(u)$, which are indeed inverses of each other and their respective domains match that of $\text{int}(\text{dom}(f))$ and $\text{int}(\text{dom}(f^*))$, respectively. To check Part (b) of Theorem 26.3, we calculate $D_{f^*}(\nabla f(y), \nabla f(x))$:

$$\begin{aligned} D_{f^*}(\nabla f(y), \nabla f(x)) &= \langle \exp(\log(y)) - \exp(\log(x)), \mathbf{1} \rangle - \langle \exp(\log(x)), \log(y) - \log(x) \rangle \\ &= \langle y - x, \mathbf{1} \rangle - \langle x, \log(y) - \log(x) \rangle, \end{aligned}$$

which is indeed equal to $D_f(x, y)$.

27.3 Let $x \in \mathbb{R}^d$. Then, the Cauchy-Schwarz inequality shows that

$$\|x\|_{A^{-1}}^2 = \langle x, A^{-1}x \rangle \leq \|x\|_{B^{-1}} \|A^{-1}x\|_B \leq \|x\|_{B^{-1}} \|A^{-1}x\|_A = \|x\|_{B^{-1}} \|x\|_{A^{-1}}.$$

Hence $\|x\|_{A^{-1}} \leq \|x\|_{B^{-1}}$ for all x , which completes the claim.

27.5 Part (a): First, we have $\|x + y\| = \sup_{u \in \mathcal{A}^\circ} |\langle x + y, u \rangle| \leq \sup_{u \in \mathcal{A}^\circ} |\langle x, u \rangle| + \sup_{u \in \mathcal{A}^\circ} |\langle y, u \rangle| \leq \sup_{u \in \mathcal{A}^\circ} |\langle x, u \rangle| + \sup_{v \in \mathcal{A}^\circ} |\langle y, v \rangle| = \|x\| + \|y\|$. For $c \in \mathbb{R}$, $x \in \mathbb{R}^d$, $\|cx\| = |c| \sup_{u \in \mathcal{A}^\circ} |\langle x, u \rangle| = |c| \|x\|$. That $\|0\| = 0$ is trivial.

Part (b): This follows immediately from the definitions: If $\mathcal{A}' \subset \mathcal{A}$ and $d(\mathcal{A}', \mathcal{A}) \leq \varepsilon$ then for any $x \in \mathcal{A}$ there exists $y \in \mathcal{A}'$ such that $\|x - y\| \leq \varepsilon$ (because \mathcal{A}' is finite, hence the inf can be replaced by min). It follows that $x \in B(y, \varepsilon)$ and that $\mathcal{A} \subset \cup_{y \in \mathcal{A}'} B(y, \varepsilon)$. For the reverse direction assume that

$\mathcal{A}' \subset \mathcal{A}$ is an ε -cover of \mathcal{A} and take $x \in \mathcal{A}$. Then, there exist $y \in \mathcal{A}'$ such that $x \in B(y, \varepsilon)$, hence $\|x - y\| \leq \varepsilon$. Since this holds for all $x \in \mathcal{A}$, $d(\mathcal{A}', \mathcal{A}) \leq \varepsilon$.

Part (c): We first show that $\mathcal{A}^* \subset B$: Take any $x \in \mathcal{A} \cup -\mathcal{A}$. By the definition of \mathcal{A}° , for any $u \in \mathcal{A}^\circ$, $|\langle x, u \rangle| \leq 1$. Taking the supremum over u gives $\|x\| = \sup_{u \in \mathcal{A}^\circ} |\langle x, u \rangle| \leq 1$, proving that $\mathcal{A}^* \subset B$. Taking the convex hull of both sides, $\text{co}(\mathcal{A}^*) \subset \text{co}(B) = B$, where the equality uses that B is convex, which follows because $\|\cdot\|$ satisfies the triangle inequality. From the above it also follows that $\text{span}(\mathcal{A}) \subset \text{span}(B) = \{z : \|z\| < \infty\}$.

It remains to be shown that $\{z : \|z\| < \infty\} \subset \text{span}(\mathcal{A})$, which is equivalent to that if $z \notin \text{span}(\mathcal{A})$ then $\|z\| = \infty$. Hence, we prove this latter implication. To show this implication, take any vector $z \in \mathbb{R}^d \setminus \text{span}(\mathcal{A})$ and write $z = z^\perp + z^\parallel$ where $z^\perp \in \mathcal{A}^\perp$, $z^\parallel \in \text{span}(\mathcal{A})$. Since $z \notin \text{span}(\mathcal{A})$, $\|z^\perp\|_2 > 0$. Now, $\mathcal{A}^\perp \subset \mathcal{A}^\circ$: Indeed, if $u \in \mathcal{A}^\perp$ then $\langle x, u \rangle = 0$ for any $x \in \mathcal{A}$, hence, $u \in \mathcal{A}^\circ$. Note also that $\mathcal{A}^\circ = -\mathcal{A}^\circ$ thanks to $|\langle x, u \rangle| = |\langle x, -u \rangle|$. Thus, $\|z\| = \sup_{u \in \mathcal{A}^\circ} \langle z, u \rangle$. Based on these, we compute $\|z\| = \sup_{u \in \mathcal{A}^\circ} \langle z, u \rangle \geq \sup_{\lambda > 0} \langle z, \lambda z^\perp \rangle = (\sup_{\lambda > 0} \lambda) \|z^\perp\|_2 = \infty$.

Part (d): We will show that $\log N(\varepsilon) \leq cd \log(1/\varepsilon)$, from which the result follows by Part (b).

First, note that with some modifications the conclusion of Exercise 20.1 remain valid. In particular, dropping \mathcal{A} from the notation, one part of this exercise claimed that $M(2\varepsilon) \leq N(\varepsilon) \leq M(\varepsilon)$, where recall that $M(\varepsilon)$ is the ε -packing number of \mathcal{A} (size of the largest ε -packing of \mathcal{A}). This part remains valid without any change, but we only need $N(\varepsilon) \leq M(\varepsilon)$ from here. Then, the next part of the exercise claimed that the inequalities

$$N(\varepsilon) \leq M(\varepsilon) \leq \frac{\text{vol}(\mathcal{A} + \frac{\varepsilon}{2}B)}{\text{vol}(\frac{\varepsilon}{2}B)} \stackrel{(*)}{\leq} \frac{\text{vol}(\frac{3}{2}\mathcal{A})}{\text{vol}(\frac{\varepsilon}{2}B)} \leq \left(\frac{3}{\varepsilon}\right)^d \frac{\text{vol}(\mathcal{A})}{\text{vol}(B)}$$

hold. Here, the first bound on $M(\varepsilon)$ remains valid without any change in the argument, though since $B \subset \text{span}(\mathcal{A})$, the volume vol is taken to be the p -dimensional Lebesgue measure on $\text{span}(\mathcal{A})$ (recall that $p = \dim(\text{span}(\mathcal{A}))$).

The next inequality, marked by (*), required that $\varepsilon B \subset \mathcal{A}$ and that \mathcal{A} is convex. We slightly modify this inequality. First, we use $\mathcal{A} + \frac{\varepsilon}{2}B \subset \text{co}(\mathcal{A}^*) + \frac{\varepsilon}{2}B$. Now, if $\varepsilon B \subset \text{co}(\mathcal{A}^*)$ then $\text{co}(\mathcal{A}^*) + \frac{\varepsilon}{2}B \subset \text{co}(\mathcal{A}^*) + \frac{1}{2}\text{co}(\mathcal{A}^*) \subset \frac{3}{2}\text{co}(\mathcal{A}^*)$ from which it follows that

$$\frac{\text{vol}(\mathcal{A} + \frac{\varepsilon}{2}B)}{\text{vol}(\frac{\varepsilon}{2}B)} \leq \frac{\text{vol}(\frac{3}{2}\text{co}(\mathcal{A}^*))}{\text{vol}(\frac{\varepsilon}{2}B)}.$$

To finish, note that for $S \subset \text{span}(\mathcal{A})$, $c > 0$, we have $\text{vol}(cS) = c^p \text{vol}(S)$.

So it remains to show that for ε small enough, $\varepsilon B \subset \text{co}(\mathcal{A}^*)$. This can be shown in two steps: Defining $B'_2 = \{z \in \text{span}(\mathcal{A}) : \|z\|_2 \leq 1\}$ as the intersection of the 2-norm unit ball with $\text{span}(\mathcal{A})$, it is enough to show that for some $\lambda_1, \lambda_2 > 0$,

$$B \subset \lambda_1 B'_2 \quad \text{and} \quad \lambda_2 B'_2 \subset \text{co}(\mathcal{A}^*). \quad (39)$$

To show the second inclusion let $r = \inf_{x \in \partial \text{co}(\mathcal{A}^*)} \|x\|_2$ where ∂ is meant in the

topology of $\text{span}(\mathcal{A})$. Clearly, $rB'_2 \subset \text{co}(\mathcal{A}^*)$. If $r > 0$, we are done. Otherwise, note that since $x \mapsto \|x\|_2$ is continuous, there exists $x^* \in \partial \text{co}(\mathcal{A}^*)$ such that $0 = r = \|x^*\|_2$, i.e., $0 \in \partial \text{co}(\mathcal{A}^*)$. Since $x^* \in \partial \text{co}(\mathcal{A}^*)$, there exists a closed halfspace H of \mathbb{R}^d such that $\text{co}(\mathcal{A}^*) \subset H \cap \text{span}(\mathcal{A}) \neq \text{span}(\mathcal{A})$ (i.e., $H \cap \text{span}(\mathcal{A})$ is a closed halfspace of $\text{span}(\mathcal{A})$) and such that $x^* = 0$ is on the boundary of H . By definition, $\mathcal{A}^* = -\mathcal{A}^*$. It follows that $\text{co}(\mathcal{A}^*) = -\text{co}(\mathcal{A}^*)$ and thus, together with $-\text{co}(\mathcal{A}^*) \subset -H \cap \text{span}(\mathcal{A})$, we get $\text{co}(\mathcal{A}^*) \subset -H \cap \text{span}(\mathcal{A})$. Thus means that $\text{co}(\mathcal{A}^*) \subset H \cap (-H) \cap \text{span}(\mathcal{A})$, implying that $\text{span}(\mathcal{A}) = \text{span}(\mathcal{A}^*) = \text{span}(\text{co}(\mathcal{A}^*)) \subset H \cap (-H) \cap \text{span}(\mathcal{A})$, which is a contradiction since $H \cap (-H)$ is an at most $\dim(\text{span}(\mathcal{A})) - 1$ dimensional hyperplane of $\text{span}(\mathcal{A})$, hence $\dim(H \cap (-H) \cap \text{span}(\mathcal{A})) < \dim(\text{span}(\mathcal{A}))$, finishing the proof of the second inclusion.

Consider now the first inclusion in (39). We claim that if $v \in \text{span}(\mathcal{A})$ and $\|v\|_2 = 1$ then $cv \in \mathcal{A}^\circ$ whenever $c \leq 1/(\max_{x \in \mathcal{A}} \|x\|_2)$. If this claim holds then taking $c = 1/(\max_{x \in \mathcal{A}} \|x\|_2)$ we compute

$$\|v\| = \sup_{u \in \mathcal{A}^\circ} |\langle v, u \rangle| \geq c \langle v, v \rangle = c.$$

Let $x \in B$. Then $x \in \text{span}(\mathcal{A})$. Let $v = x/\|x\|_2$ so that $\|v\|_2 = 1$ and $\|v\| = \|x\|/\|x\|_2$. By the previous inequality $\|v\| \geq c$. Putting things together, $c\|x\|_2 \leq \|x\| \leq 1$. Hence, $cB \subset B'_2$ and $c > 0$ follows because A is bounded.

It remains to prove the claim. For this let $v \in \text{span}(\mathcal{A})$, $\|v\|_2 = 1$, $0 \leq c \leq 1/(\max_{x \in \mathcal{A}} \|x\|_2)$. We want to show that $cv \in \mathcal{A}^\circ$. For this to hold we need to show that for any $x \in \mathcal{A}$, $|\langle x, cv \rangle| \leq 1$. Take some $x \in \mathcal{A}$. Then, $|\langle x, cv \rangle| = c|\langle x, v \rangle| \leq c \max_{y \in \mathcal{A}} |\langle y, v \rangle| \leq c \max_{y \in \mathcal{A}} \|y\|_2 \leq 1$, where the one but last inequality follows because $\|v\|_2 = 1$ and the last inequality follows by the choice of c .

An alternative solution would have been to show the equivalence of $\|\cdot\|$ to $\|\cdot\|_2$ when both are restricted to $\text{span}(\mathcal{A})$ (actually, this is most of the work in the above proof) and then argue that for equivalent norms covering and packing numbers change only by little and then call the result of Exercise 20.1.

27.10 Let $\alpha = \sup_{x, y \in \mathcal{K}} \langle x - y, u \rangle$ and $\mathcal{K}_t = \{x \in \mathcal{K} : \langle x^* - x, u \rangle = t\}$ and $f(t) = \text{vol}(\mathcal{K}_t)$. The volume here is calculated with respect the $(d-1)$ -dimensional volume form, since this is the dimension of the slice \mathcal{K}_t (see figure). Now,

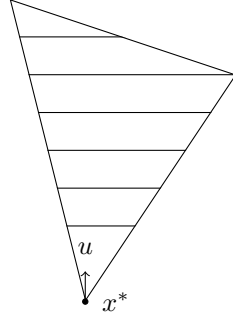
$$\frac{\text{vol}(\mathcal{K}_t)}{\int_{\mathcal{K}} \exp(-\langle x^* - x, u \rangle)} dx = \frac{\int_0^{\alpha/\|u\|} f(t) dt}{\int_0^{\alpha/\|u\|} \exp(-t\|u\|) f(t) dt} = \frac{\int_0^1 f\left(\frac{\alpha s}{\|u\|}\right) ds}{\int_0^1 f\left(\frac{\alpha s}{\|u\|}\right) \exp(-\alpha s) ds}.$$

Convexity ensures that $f(t) \geq (tf(t_0))^{d-1}$ for any $t_0 \geq 0$. The worst-case occurs when this is an equality and so we may take $f(x) = x^{d-1}$. Then,

$$\frac{\text{vol}(\mathcal{K}_t)}{\int_{\mathcal{K}} \exp(-\langle x^* - x, u \rangle)} dx \leq \frac{\int_0^1 s^{d-1} ds}{\int_0^1 s^{d-1} \exp(-\alpha s) ds} \leq e \max(1, \alpha^d),$$

where the last inequality follows by treating the cases with $\alpha < 1$ and $\alpha \geq 1$

differently. In the former case $\exp(-\alpha s) \geq 1/e$ for $s \in [0, 1]$ while in the latter the integral in the denominator is truncated to $s \in [0, 1/\alpha]$ for which $\exp(-\alpha s) \geq 1/e$ and direct calculation yields the result.



The whole triangle is \mathcal{K} with x^* at the bottom corner. The thin lines represent \mathcal{K}_t for different values of t , which $(d-1)$ -dimensional subsets of \mathcal{K} that lie in affine spaces with normal vector u .

28.5

(a) For the first relation we have $\nabla F^*(x)_i = \exp(x)_i$ and $F^*(x) = \sum_{i=1}^K \exp(x_i)$. Then $D_F(P_t, \tilde{P}_{t+1}) = D_{F^*}(\nabla F^*(\tilde{P}_{t+1}), \nabla F^*(P_t)) = \sum_{i=1}^K P_{ti} (\exp(-\eta \hat{Y}_{ti}) - 1 + \eta \hat{Y}_{ti})$. The second relation follows from the inequality $\exp(x) \leq 1 + x + x^2/2$ for $x \leq 0$.

(b) Using part (a) we have

$$\begin{aligned} \frac{1}{\eta} \mathbb{E} \left[\sum_{t=1}^n D_F(P_t, \tilde{P}_{t+1}) \right] &\leq \frac{\eta}{2} \mathbb{E} \left[\sum_{t=1}^n \sum_{i=1}^K P_{ti} \hat{Y}_{ti}^2 \right] \\ &\leq \frac{\eta}{2} \mathbb{E} \left[\sum_{t=1}^n \sum_{i=1}^K \frac{\mathbb{I}\{A_t = i\}}{P_{ti}} \right] \\ &= \frac{\eta n K}{2}. \end{aligned}$$

(c) Simple calculus shows that for $p \in \mathcal{P}_{K-1}$, $F(p) \geq -\log(K) - 1$ and $F(p) \leq -1$ is obvious. Therefore $\text{diam}_F(\mathcal{A}) = \max_{p,q \in \mathcal{A}} F(p) - F(q) \leq \log(K)$.

(d) By the previous exercise Exp3 chooses A_t sampled from P_t . Then applying Theorem 28.1 and parts (b) and (c) and choosing $\eta = \sqrt{\log(K)/(2nK)}$ yields the result.

28.6 The first step is the same as the proof of Theorem 28.1.

$$R_n(a) = \sum_{t=1}^n \langle a_t - a, y_t \rangle = \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle + \sum_{t=1}^n \langle a_{t+1} - a, y_t \rangle.$$

Next let $\Phi_t(a) = F(a)/\eta + \sum_{s=1}^t \langle a, y_s \rangle$ so that

$$\begin{aligned}
 \sum_{t=1}^n \langle a_{t+1} - a, y_t \rangle &= \sum_{t=1}^n \langle a_{t+1}, y_t \rangle - \Phi_n(a) + \frac{F(a)}{\eta} \\
 &= \sum_{t=1}^n (\Phi_t(a_{t+1}) - \Phi_{t-1}(a_{t+1})) - \Phi_n(a) + \frac{F(a)}{\eta} \\
 &= -\Phi_0(a_1) + \sum_{t=0}^{n-1} (\Phi_t(a_{t+1}) - \Phi_t(a_{t+2})) + \Phi_n(a_{n+1}) - \Phi_n(a) + \frac{F(a)}{\eta} \\
 &\leq \frac{F(a) - F(a_1)}{\eta} + \sum_{t=0}^{n-1} (\Phi_t(a_{t+1}) - \Phi_t(a_{t+2})). \tag{.40}
 \end{aligned}$$

Now $D_{\Phi_t}(a, b) = \frac{1}{\eta} D_F(a, b)$. Therefore

$$\begin{aligned}
 \Phi_t(a_{t+1}) - \Phi_t(a_{t+2}) &= \langle \nabla \Phi_t(a_{t+1}), a_{t+1} - a_{t+2} \rangle - \frac{1}{\eta} D_F(a_{t+2}, a_{t+1}) \\
 &\leq -\frac{1}{\eta} D_F(a_{t+2}, a_{t+1}).
 \end{aligned}$$

Substituting this into Eq. (.40) completes the proof.

28.7 Abbreviate $D(x, y) = D_F(x, y)$. By the definition of \tilde{a}_{t+1} and the first-order optimality conditions we have $\eta_t y_t = \nabla F(a_t) - \nabla F(\tilde{a}_{t+1})$. Therefore

$$\begin{aligned}
 \langle a_t - a, y_t \rangle &= \frac{1}{\eta_t} \langle a_t - a, \nabla F(a_t) - \nabla F(\tilde{a}_{t+1}) \rangle \\
 &= \frac{1}{\eta_t} (-\langle a - a_t, \nabla F(a_t) \rangle - \langle a_t - \tilde{a}_{t+1}, \nabla F(\tilde{a}_{t+1}) \rangle + \langle a - \tilde{a}_{t+1}, \nabla F(\tilde{a}_{t+1}) \rangle) \\
 &= \frac{1}{\eta_t} (D(a, a_t) - D(a, \tilde{a}_{t+1}) + D(a_t, \tilde{a}_{t+1})).
 \end{aligned}$$

Summing completes the proof. For the second part use the fact that $D(a, \tilde{a}_{t+1}) \geq D(a, a_{t+1})$.

28.8 We use the same argument as the solution to Exercise 28.6. First,

$$R_n(a) = \sum_{t=1}^n \langle a_t - a, y_t \rangle = \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle + \sum_{t=1}^n \langle a_{t+1} - a, y_t \rangle.$$

The next step also mirrors that in Exercise 28.6, but now we have to keep track

of the changing potentials:

$$\begin{aligned}
\sum_{t=1}^n \langle a_{t+1} - a, y_t \rangle &= \sum_{t=1}^n \langle a_{t+1}, y_t \rangle - \Phi_n(a) + F_n(a) \\
&= \sum_{t=1}^n (\Phi_t(a_{t+1}) - \Phi_{t-1}(a_{t+1})) + \sum_{t=1}^n (F_{t-1}(a_{t+1}) - F_t(a_{t+1})) - \Phi_n(a) + F_n(a) \\
&= -\Phi_0(a_1) + \sum_{t=0}^{n-1} (\Phi_t(a_{t+1}) - \Phi_t(a_{t+2})) + \Phi_n(a_{n+1}) - \Phi_n(a) \\
&\quad + F_n(a) + \sum_{t=1}^n (F_{t-1}(a_{t+1}) - F_t(a_{t+1})) \\
&\leq F_n(a) - F_0(a_1) + \sum_{t=0}^{n-1} (\Phi_t(a_{t+1}) - \Phi_t(a_{t+2})) + \sum_{t=1}^n (F_{t-1}(a_{t+1}) - F_t(a_{t+1})) .
\end{aligned}$$

Now $D_{\Phi_t}(a, b) = D_{F_t}(a, b)$. Therefore

$$\begin{aligned}
\Phi_t(a_{t+1}) - \Phi_t(a_{t+2}) &= \langle \nabla \Phi_t(a_{t+1}), a_{t+1} - a_{t+2} \rangle - D_{F_t}(a_{t+2}, a_{t+1}) \\
&\leq -D_{F_t}(a_{t+2}, a_{t+1}),
\end{aligned}$$

which combined with the previous big display completes the proof.

28.9

- (a) Apply your solution to Exercise 26.9.
(b) Since \hat{Y}_t is unbiased we have

$$R_n = \max_{i \in [K]} \mathbb{E} \left[\sum_{t=1}^n (y_{tA_t} - y_{ti}) \right] = \max_{i \in [K]} \mathbb{E} \left[\sum_{t=1}^n \langle P_t - P, \hat{Y}_t \rangle \right].$$

Then apply the result from Exercise 28.8.

- (c) The negentropy satisfies $F(p) \in [-\log(K) - 1, -1]$ for all $p \in \mathcal{P}_{K-1}$. Therefore

$$\begin{aligned}
F_n(P) - F_0(P_1) + \sum_{t=1}^n (F_{t-1}(P_{t+1}) - F_t(P_{t+1})) \\
&= \frac{F(P)}{\eta_n} - \frac{F(P_1)}{\eta_0} - \sum_{t=1}^n F(P_{t+1}) \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \\
&\leq -\frac{1}{\eta_n} + \frac{\log(K) + 1}{\eta_0} + \sum_{t=1}^n (1 + \log(K)) \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \\
&= \frac{\log(K)}{\eta_n}.
\end{aligned}$$

- (d) Consider two cases. First, if $P_{t+1, A_t} \geq P_{tA_t}$, then

$$\langle P_t - P_{t+1}, \hat{Y}_t \rangle = (P_{tA_t} - P_{t+1, A_t}) \hat{Y}_{tA_t} \leq 0.$$

On the other hand, if $P_{t+1,A_t} \leq P_{tA_t}$, then by Theorem 26.5 with $H = \nabla^2 f(q) = \text{diag}(1/q)$ for some $q \in [P_t, P_{t+1}]$ we have

$$\langle P_t - P_{t+1}, \hat{Y}_t \rangle - \frac{D_F(P_t, P_{t+1})}{\eta_t} \leq \frac{\eta_t}{2} \|\hat{Y}_t\|_{H^{-1}}^2,$$

Since $P_{t+1,A_t} \leq P_{tA_t}$ we have

$$\frac{\eta_t}{2} \|\hat{Y}_t\|_{H^{-1}}^2 = \frac{\eta_t q_{A_t} \hat{Y}_{tA_t}^2}{2} \leq \frac{\eta_t P_{tA_t} \hat{Y}_{tA_t}^2}{2} \leq \frac{\eta_t y_{tA_t}^2}{2P_{tA_t}}.$$

Therefore

$$\langle P_t - P_{t+1}, \hat{Y}_t \rangle - \frac{D_F(P_t, P_{t+1})}{\eta_t} \leq \frac{\eta_t y_{tA_t}^2}{2P_{tA_t}} \leq \frac{\eta_t}{2P_{tA_t}}.$$

(e)

$$R_n \leq \mathbb{E} \left[\frac{\log(K)}{\eta_n} + \frac{1}{2} \sum_{t=1}^n \frac{\eta_t}{P_{tA_t}} \right] = \frac{\log(K)}{\eta_n} + \frac{K}{2} \sum_{t=1}^n \eta_t.$$

(f) Choose $\eta_t = \sqrt{\frac{\log(K)}{Kt}}$ and use the fact that $\sum_{t=1}^n \sqrt{1/t} \leq 2\sqrt{n}$.

28.12 We make use of mirror descent on the simplex with negentropy regularization. Let $p_t \in \mathcal{P}_{K-1}$ be the choice of mirror descent in round t and let $q_t = \text{argmax}_q \langle p_t, Gq \rangle$ and $y_t = Gq_t$, which means that $p_1 = (1/K, \dots, 1/K)$ and

$$p_t = \text{argmin}_{p \in \mathcal{P}_{K-1}} \langle p, y_{t-1} \rangle + D_F(p, p_{t-1}),$$

where $F(p) = -\sum_{i=1}^K p_i \log(p_i) - p_i$ is the unnormalized negentropy. Then let $\bar{p}_n = \frac{1}{n} \sum_{t=1}^n p_t$ and $\bar{q}_n = \frac{1}{n} \sum_{t=1}^n q_t$. By the definition of the regret, q_t and straightforward algebra,

$$\begin{aligned} R_n &= \max_p \sum_{t=1}^n \langle p_t - p, Gq_t \rangle \\ &= \sum_{t=1}^n \langle p_t, Gq_t \rangle - \min_p \sum_{t=1}^n \langle p, Gq_t \rangle \\ &\geq \max_q \sum_{t=1}^n \langle p_t, Gq \rangle - \min_p \sum_{t=1}^n \langle p, Gq_t \rangle \\ &= n \max_q \langle \bar{p}_n, Gq \rangle - n \min_p \langle p, G\bar{q}_n \rangle \\ &\geq n \max_q \langle \bar{p}_n, Gq \rangle - n \max_q \min_p \langle p, Gq \rangle. \end{aligned}$$

Dividing both sides by n and taking the limit as n tends to infinity shows that

$$\limsup_{n \rightarrow \infty} \max_q \langle \bar{p}_n, Gq \rangle \leq \max_q \min_p \langle p, Gq \rangle.$$

Note that here we used the result from Exercise 28.11 that $R_n \leq \sqrt{2n \log(K)}$. Since $\bar{p}_n \in \mathcal{P}_{E-1}$, which is compact, it follows that there exists a p^* such that

$$\min_p \max_q \langle p, Gq \rangle \leq \max_q \langle p^*, Gq \rangle \leq \max_q \min_p \langle p, Gq \rangle.$$

The result is complete because $\min_p \max_q \langle p, Gq \rangle \geq \max_q \min_p \langle p, Gq \rangle$ is obvious.

31.1 As suggested, Exp4 is used with each element of $\Gamma_{n,m}$ identified with one expert. Consider an arbitrary enumeration of $\Gamma_{n,m} = \{a^{(1)}, \dots, a^{(G)}\}$ where $G = |\Gamma_{n,m}|$. The predictions of expert $g \in [G]$ for round $t \in [n]$ encoded as a probability distribution over $[K]$ (as required by the prediction-with-expert-advice framework) is $E_{g,k}^t = \mathbb{I}\{a_t^{(g)} = k\}$, $k \in [K]$. The expected regret of Exp4 when used with these experts is

$$R_n^{\text{experts}} = \mathbb{E} \left[\sum_{t=1}^n y_{tA_t} - \min_{g \in [G]} \sum_{t=1}^n E_g^{(t)} y_t \right],$$

where compared to Chapter 18 we switched to losses. By definition,

$$\sum_{t=1}^n E_g^{(t)} y_t = \sum_{t=1}^n y_{t, a_t^{(g)}}$$

and hence

$$R_n^{\text{experts}} = R_{n,m}.$$

Thus, Theorem 18.1 indeed proves (31.1). To prove (31.2) it remains to show that $G = |\Gamma_{n,m}| \leq Cm \log(Kn/m)$. For this note that $G = \sum_{s=1}^m G_{n,s}^*$ where $G_{n,s}^*$ is the number of sequences from $[K]^n$ that switch exactly $s-1$ times. When $m-1 \leq n/2$, a crude upper bound on G is $mG_{n,m}^*$. For $s=1$, $G_{n,s}^* = K$. For $s > 1$, a sequence with $s-1$ switches is determined by the location of the switches, and the identity of the action taken in each segment where the action does not change. The possible switch locations are of the form $(t, t+1)$ with $t = 1, \dots, n-1$. Thus the number of these locations is $n-1$, of which, we need to choose $s-1$. There are $\binom{n-1}{s-1}$ ways of doing this. Since there are s segments and for the first segment we can choose any action and for the others we can choose any other action than the one chosen for the previous segments, there are KK^{s-1} valid ways of assigning actions to segments. Thus, $G_{n,s}^* = KK^{s-1} \binom{n-1}{s-1}$. Define $\Phi_m(n) = \sum_{i=0}^m \binom{n}{i}$. Hence, $G \leq K^m \sum_{s=0}^{m-1} \binom{n-1}{s} = K^m \Phi_{m-1}(n-1) \leq K^m \Phi_m(n)$. Now note that for $n \geq m$, $0 \leq m/n \leq 1$, hence

$$\left(\frac{m}{n}\right)^m \Phi_m(n) \leq \sum_{i=0}^m \left(\frac{m}{n}\right)^i \binom{n}{i} \leq \sum_{i=0}^n \left(\frac{m}{n}\right)^i \binom{n}{i} = \left(1 + \frac{m}{n}\right)^n \leq e^m.$$

Reordering gives $\Phi_m(n) \leq \left(\frac{en}{m}\right)^m$. Hence, $\log(G) \leq m \log(en/m)$. Plugging this into (31.1) gives (31.2).

31.3 Use the construction and analysis in Exercise 11.5 and note that when

$m = 2$ the random version of the regret is nonnegative on the bandit constructed there.

32.2 The argument is half-convincing. The heart of the argument is that under the criterion that at least one item should attract the user, it may be suboptimal to present the list composed of the fittest items. The example with the query ‘jaguar’ is clear: Assume half of the users will mean ‘jaguar’ as the big cat, while the other half will mean it as the car. Presenting items that are relevant for both meanings may have a better chance to satisfy a randomly picked user than going with the top K list, which may happen to support only one of the meanings. This shows that there is indeed an issue with ‘linearizing’ the problem by just considering individual item fitness values.

However, the argument is confusing in other ways. First, it treats conditions (for example, independence of attractiveness) that are sufficient but not necessary to validate the probabilistic ranking principle (PRP) as if they were also necessary. In fact, in click model studied here, the mentioned independence assumption is not needed. To clarify, the strong assumption in the stochastic click model, is that the optimal list is indeed optimal. Under this assumption, the independence assumption is not needed.

Next, that the same document can have different relevance to different users fits even the cascade model, where the vector of attraction values are different each time they are sampled from the model. So this alone would not undermine the PRP.

Finally, the last sentence confuses relevance and ‘usefulness’. Again, in the cascade model, the relevance (attractiveness) of a document (item) does not depend on the relevance of any other document. Yet in the reward in the cascade model is exactly one if and only if at least one document presented is relevant (attractive).

32.6 Following the proof of Theorem 32.1 the first part until Eq. (32.5) we have

$$R_n \leq nK\mathbb{P}(F_n) + \sum_{j=1}^L \sum_{i=1}^{\min\{K,j-1\}} \mathbb{E} \left[\mathbb{I}\{F_n^c\} \sum_{t=1}^n U_{tij} \right].$$

As before the first term is bounded using Lemma 32.2. Then using the first part of the proof of Lemma 32.5 shows that

$$\mathbb{I}\{F_n^c\} \sum_{t=1}^n U_{tij} \leq 1 + \sqrt{2N_{nij} \log \left(\frac{c\sqrt{n}}{\delta} \right)}.$$

Substituting into the previous display and applying Cauchy-Schwarz shows that

$$R_n \leq nK\mathbb{P}(F_n) + KL + \sqrt{2KL\mathbb{E} \left[\sum_{j=1}^L \sum_{i=1}^{\min\{K,j-1\}} N_{nij} \right] \log \left(\frac{c\sqrt{n}}{\delta} \right)}.$$

Writing out the definition of N_{nij} reveals that we need to bound

$$\begin{aligned} \mathbb{E} \left[\sum_{j=1}^L \sum_{i=1}^{\min\{K, j-1\}} N_{nij} \right] &\leq \sum_{t=1}^n \mathbb{E} \left[\mathbb{E}_{t-1} \left[\sum_{d=1}^{M_t} \sum_{j \in \mathcal{P}_{td}} \sum_{i \in \mathcal{P}_{td} \cap [K]} U_{tij} \middle| \mathcal{F}_{t-1} \right] \right] \\ &\leq \sum_{t=1}^n \mathbb{E} \left[\mathbb{E}_{t-1} \left[\sum_{d=1}^{M_t} \sum_{j \in \mathcal{P}_{td}} \sum_{i \in \mathcal{P}_{td} \cap [K]} (C_{ti} + C_{tj}) \middle| \mathcal{F}_{t-1} \right] \right] = (\text{A}). \end{aligned}$$

Expanding the two terms in the inner sum and bounding each separately leads to

$$\begin{aligned} \mathbb{E}_{t-1} \left[\sum_{d=1}^{M_t} \sum_{j \in \mathcal{P}_{td}} \sum_{i \in \mathcal{P}_{td} \cap [K]} C_{ti} \middle| \mathcal{F}_{t-1} \right] &= \mathbb{E}_{t-1} \left[\sum_{d=1}^{M_t} |\mathcal{P}_{td}| \sum_{i \in \mathcal{P}_{td} \cap [K]} C_{ti} \middle| \mathcal{F}_{t-1} \right] \\ &\leq \sum_{d=1}^{M_t} |\mathcal{I}_{td} \cap [K]| |\mathcal{P}_{td} \cap [K]| \leq K^2. \end{aligned}$$

For the second term,

$$\begin{aligned} \mathbb{E}_{t-1} \left[\sum_{d=1}^{M_t} \sum_{j \in \mathcal{P}_{td}} \sum_{i \in \mathcal{P}_{td} \cap [K]} C_{tj} \middle| \mathcal{F}_{t-1} \right] &= \mathbb{E}_{t-1} \left[\sum_{d=1}^{M_t} |\mathcal{P}_{td} \cap [K]| \sum_{j \in \mathcal{P}_{td}} C_{tj} \middle| \mathcal{F}_{t-1} \right] \\ &\leq \sum_{d=1}^{M_t} |\mathcal{P}_{td} \cap [K]| |\mathcal{I}_{td} \cap [K]| \leq K^2. \end{aligned}$$

Hence (A) $\leq nK^2$ and $R_n \leq nK\mathbb{P}(F_n) + KL + \sqrt{4K^3Ln \log\left(\frac{c\sqrt{n}}{\delta}\right)}$ and the result follows from Lemma 32.2.

33.3 Abbreviate $f(\alpha) = \inf_{d \in D} \langle \alpha, d \rangle$, which is clearly positively homogeneous: $f(c\alpha) = cf(\alpha)$ for any $c \geq 0$. Because D is nonempty, $f(\mathbf{0}) = 0$. Hence we can ignore $\alpha = \mathbf{0}$ in both optimization problems and so

$$\begin{aligned} \left(\sup_{\alpha \in \mathcal{P}_{K-1}} f(\alpha) \right)^{-1} L &= \inf_{\alpha \in \mathcal{P}_{K-1}} \frac{L}{f(\alpha)} \\ &= \inf_{\alpha \geq 0: \|\alpha\|_1 > 0} \frac{L\|\alpha\|_1}{f(\alpha)} \\ &= \inf_{\alpha \geq 0: \|\alpha\|_1 > 0} \|L\alpha/f(\alpha)\|_1 \\ &= \inf \{ \|\alpha\|_1 : f(\alpha) \geq L \}, \end{aligned}$$

where we used the positive homogeneity of $f(\alpha)$ and the ℓ_1 norm.

33.4

(a) For each $i > 1$ define

$$\mathcal{E}_i = \{ \tilde{\nu} \in \mathcal{E} : \mu_1(\tilde{\nu}) = \mu_i(\tilde{\nu}) \text{ and } \mu_j(\tilde{\nu}) = \mu_j(\nu) \text{ for } j \notin \{1, i\} \}.$$

You can easily show that

$$\begin{aligned} \inf_{\tilde{\nu} \in \mathcal{E}_{\text{alt}}(\nu)} \sum_{i=1}^K \alpha_i D(\nu_i, \tilde{\nu}_i) &= \min_{i>1} \inf_{\tilde{\nu} \in \mathcal{E}_{\text{alt}}(\nu)} (\alpha_1 D(\nu_1, \tilde{\nu}_1) + \alpha_i D(\nu_i, \tilde{\nu}_i)) \\ &= \min_{i>1} \inf_{\tilde{\mu} \in \mathbb{R}} \left(\frac{\alpha_1(\mu_1(\nu) - \tilde{\mu})^2}{2\sigma_1^2} + \frac{\alpha_i(\mu_i(\nu) - \tilde{\mu})^2}{2\sigma_i^2} \right) \\ &= \frac{1}{2} \min_{i>1} \frac{\alpha_1 \alpha_i \Delta_i^2}{\alpha_1 \sigma_i^2 + \alpha_i \sigma_1^2}. \end{aligned}$$

(b) Let $\alpha_1 = \alpha$ so that $\alpha_2 = 1 - \alpha$. By the previous part

$$\begin{aligned} (c^*(\nu))^{-1} &= \max_{\alpha \in [0,1]} \inf_{\tilde{\nu} \in \mathcal{E}_{\text{alt}}(\nu)} \sum_{i=1}^K \alpha_i D(\nu_i, \tilde{\nu}_i) = \max_{\alpha \in [0,1]} \frac{\alpha_1 \alpha_2 \Delta_2^2}{\alpha_1 \sigma_2^2 + \alpha_2 \sigma_1^2} \\ &= \max_{\alpha \in [0,1]} \frac{\alpha(1-\alpha)\Delta_2^2}{\alpha\sigma_2^2 + (1-\alpha)\sigma_1^2} \\ &= \frac{\Delta_2^2}{2(\sigma_1^2 + \sigma_2^2)}. \end{aligned}$$

(c) By the result in Exercise 33.3 and Part (a) of this exercise,

$$\begin{aligned} c^*(\nu) &= \inf \left\{ \|\alpha\|_1 : \alpha \in [0, \infty)^K, \inf_{\tilde{\nu} \in \mathcal{E}_{\text{alt}}(\nu)} \sum_{i=1}^K \alpha_i D(\nu_i, \tilde{\nu}_i) = 1 \right\} \\ &= \inf \left\{ \|\alpha\|_1 : \alpha \in [0, \infty)^K, \min_{i>1} \frac{\alpha_1 \alpha_i \Delta_i^2}{2\alpha_1 \sigma_i^2 + 2\alpha_i \sigma_1^2} = 1 \right\}. \end{aligned}$$

Let $\alpha_1 = 2a\sigma_1^2/\Delta_{\min}^2$, which by the constraint that $\alpha \geq 0$ must satisfy $a > 1$. Then

$$\begin{aligned} c^*(\nu) &= \inf_{\alpha_1 > 2\sigma_1^2/\Delta_{\min}^2} \alpha_1 + \sum_{i=2}^K \frac{2\alpha_1 \sigma_i^2}{\alpha_1 \Delta_i^2 - 2\sigma_1^2} \\ &\leq \inf_{a>1} \frac{2a\sigma_1^2}{\Delta_{\min}^2} + \frac{a}{a-1} \sum_{i=2}^K \frac{2\sigma_i^2/\Delta_i^2}{a-1} \\ &= \left(\sqrt{\frac{2\sigma_1^2}{\Delta_{\min}^2}} + \sqrt{\sum_{i=2}^K \frac{2\sigma_i^2}{\Delta_i^2}} \right)^2 \\ &= \frac{2\sigma_1^2}{\Delta_{\min}^2} + \sum_{i=2}^K \frac{2\sigma_i^2}{\Delta_i^2} + \frac{4\sigma_1}{\Delta_{\min}} \sqrt{\sum_{i=2}^K \frac{2\sigma_i^2}{\Delta_i^2}}. \end{aligned}$$

(d) From the previous part

$$c^*(\nu) = \inf \left\{ \|\alpha\|_1 : \alpha \in [0, \infty)^K, \min_{i>1} \frac{\alpha_1 \alpha_i \Delta_i^2}{2\alpha_1 \sigma_i^2 + 2\alpha_i \sigma_1^2} = 1 \right\},$$

Let $\alpha_1 = 2a\sigma_1^2/\Delta_{\min}^2$ with $a > 1$. Then

$$\begin{aligned} c^*(\nu) &= \inf_{\alpha_1 > 2\sigma_1^2/\Delta_{\min}^2} \left(\alpha_1 + \sum_{i=2}^K \frac{2\alpha_1\sigma_i^2}{\alpha_1\Delta_i^2 - 2\sigma_1^2} \right) \\ &\leq \inf_{a > 1} \left(\frac{2a\sigma_1^2}{\Delta_{\min}^2} + \frac{a}{a-1} \sum_{i=2}^K \frac{2\sigma_i^2}{\Delta_i^2} \right) \\ &= \left(\sqrt{\frac{2\sigma_1^2}{\Delta_{\min}^2}} + \sqrt{\sum_{i=2}^K \frac{2\sigma_i^2}{\Delta_i^2}} \right)^2 \\ &= \frac{2\sigma_1^2}{\Delta_{\min}^2} + \sum_{i=2}^K \frac{2\sigma_i^2}{\Delta_i^2} + \frac{4\sigma_1}{\Delta_{\min}} \sqrt{\sum_{i=2}^K \frac{2\sigma_i^2}{\Delta_i^2}}. \end{aligned}$$

(e) Notice that the inequality in the previous part is now an equality.

33.6

(a) Let $\nu \in \mathcal{E}$ be an arbitrary Gaussian bandit with $\mu_1(\nu) > \max_{i>1} \mu_i(\nu)$ and assume that

$$\liminf_{n \rightarrow \infty} \frac{-\log(\mathbb{P}_{\nu\pi}(\Delta_{A_{n+1}} > 0))}{\log(n)} > 1 + \varepsilon. \quad (.41)$$

Notice that if Eq. (.41) were not true then we would be done. Then let ν' be a Gaussian bandit in $\mathcal{E}_{\text{alt}}(\nu)$ with $\mu(\nu') = \mu(\nu)$ except that $\mu_i(\nu') = \mu_i(\nu) + \Delta_i(\nu)(1+\delta)$ where $i > 1$ and $\delta = \sqrt{1+\varepsilon} - 1$. By Theorem 14.2 and Lemma 15.1,

$$\begin{aligned} \mathbb{P}_{\nu\pi}(A_{n+1} \neq 1) + \mathbb{P}_{\nu'\pi}(A_{n+1} \neq i) &\geq \frac{1}{2} \exp(-D(\mathbb{P}_{\nu\pi}, \mathbb{P}_{\nu'\pi})) \\ &\geq \frac{1}{2} \exp\left(-\frac{(1+\delta)^2 \Delta_i(\nu)^2 \mathbb{E}_{\nu\pi}[T_i(n)] \log(n)}{2}\right). \end{aligned}$$

Because π is assumed to be asymptotically optimal $\lim_{n \rightarrow \infty} \mathbb{E}_{\nu\pi}[T_i(n)]/\log(n) = 2/\Delta_i(\nu)$ and hence

$$\mathbb{P}_{\nu\pi}(A_{n+1} \neq 1) + \mathbb{P}_{\nu'\pi}(A_{n+1} \neq i) = \Omega\left(\frac{1}{n^{1+\varepsilon}}\right).$$

Using Eq. (.41) shows that

$$\limsup_{n \rightarrow \infty} n^{1+\varepsilon} \mathbb{P}_{\nu'\pi}(A_{n+1} \neq i) > 0,$$

which implies that

$$\liminf_{n \rightarrow \infty} \frac{-\log(\mathbb{P}_{\nu'\pi}(A_{n+1} \neq i))}{\log(n)} \leq (1+\delta)^2 \leq 1 + \varepsilon.$$

(b) No. Consider the algorithm that chooses $A_n = 1 + (n \bmod K)$.

- (c) The same argument as Part (a) shows there exists a $\nu \in \mathcal{E}$ with a unique optimal arm such that

$$\liminf_{n \rightarrow \infty} \frac{-\log(\mathbb{P}_{\nu\pi}(A_{n+1} \notin i^*(\nu)))}{\log(n)} = O(1),$$

which means the probability of selecting a suboptimal arm decays only polynomially with n .

33.7

- (a) Given a bandit $\tilde{\nu} \in \mathcal{E}$ define $\Phi : \mathcal{E} \times \mathcal{P}_{K-1} \rightarrow \mathbb{R}$ by

$$\Phi_{\tilde{\nu}}(\nu, x) = \sum_{i=1}^K x_i (\mu_i(\nu) - \mu_i(\tilde{\nu}))^2,$$

which is continuous for fixed $\tilde{\nu}$ and satisfies

$$\Phi(\nu, x) = \inf_{\tilde{\nu} \in \mathcal{E}_{\text{alt}}(\nu)} \Phi_{\tilde{\nu}}(\nu, x). \quad (.42)$$

Let $\varepsilon > 0$ be sufficiently small so that $d(\nu, \xi) \leq \varepsilon$ implies that $i^*(\xi) = i^*(\nu)$, which exists since it was assumed that ν has a unique optimal arm. Define

$$\mathcal{E}_{\text{alt}} = \left\{ \tilde{\nu} \in \mathcal{E}_{\text{alt}}(\nu) : \min_i \mu_i(\tilde{\nu}) \geq \min_i \mu_i(\nu) - 2\varepsilon, \max_i \mu_i(\tilde{\nu}) \leq \max_i \mu_i(\nu) + 2\varepsilon \right\}$$

Then for any ξ and $\beta \in \mathcal{P}_{K-1}$ with $d(\nu, \xi) \leq \varepsilon$ it holds that

$$\Phi(\xi, \beta) = \inf_{\tilde{\nu} \in \mathcal{E}_{\text{alt}}} \Phi_{\tilde{\nu}}(\xi, \beta).$$

Hence on the ball $B = \{(\xi, \beta) : d(\nu, \xi) < \varepsilon\}$ the function Φ is an infimum over a compact set of continuous functions and hence Φ is continuous on B . That α^* is continuous now follows from the fact in the hint and because $\alpha^*(\nu)$ is unique.

- (b) Define random variable $\Lambda \geq 1$ by

$$\Lambda = \min \left\{ \lambda \geq 1 : d(\hat{\nu}_t, \nu) \leq \sqrt{\frac{2 \log(\Lambda K t(t+1))}{\min_i T_i(t)}} \text{ for all } t \right\}.$$

Then $\mathbb{E}[\log(\Lambda)] \leq 1$. Hence $d(\hat{\nu}_t, \nu) \leq \delta$ for all $t \geq \tau_\delta$ given by

$$\tau_\nu = \min \left\{ t : \sqrt{\frac{2 \log(\Lambda K t(t+1))}{\min_i T_i(t)}} \leq \varepsilon(\delta) \right\}.$$

- (c) Let $w(\varepsilon) = \inf\{x : d(\xi, \nu) \leq x \implies \|\alpha^*(\nu) - \alpha^*(\xi)\| \leq \varepsilon\}$, which by (a) satisfies $w(\varepsilon) > 0$ for all $\varepsilon > 0$. Hence $\mathbb{E}[\tau_\alpha(\varepsilon)] \leq \mathbb{E}[\tau_\nu(w(\varepsilon))] < \infty$.
- (d) By definition of the algorithm $A_t = i$ implies that either $T_i(t-1) \leq \sqrt{t}$ or $A_t = \operatorname{argmax}_i \alpha_i^*(\hat{\nu}_{t-1}) - T_i(t-1)/(t-1)$. Now suppose that

$$t \geq \frac{2}{\varepsilon} (\tau_\alpha(\varepsilon/(2K)) + K(1 + \sqrt{t})).$$

Then the definition of the algorithm implies that

$$T_i(t) \leq \max \left\{ T_i(\tau_\alpha(\varepsilon/(2K))), 1 + t(\alpha_i^*(\nu) + \varepsilon/(2K)), 1 + \sqrt{t} \right\} \leq t\alpha_i^*(\nu) + t\varepsilon.$$

Furthermore, since $\sum_{i=1}^K T_i(t) = t$ it follows that

$$\begin{aligned} T_i(t) &\geq t - \sum_{j \neq i} \max \left\{ T_j(\tau_\alpha(\varepsilon/(2K))), 1 + t(\alpha_j^*(\nu) + \varepsilon/(2K)), 1 + \sqrt{t} \right\} \\ &\geq t\alpha_i^*(\nu) - \tau_\alpha(\varepsilon/(2K)) - K(1 + \sqrt{t}) - t\varepsilon/2 \\ &\geq t\alpha_i^*(\nu) - t\varepsilon. \end{aligned}$$

And the result follows from the previous part, which ensures that

$$\mathbb{E} \left[\frac{2}{\varepsilon} (\tau_\alpha(\varepsilon/(2K)) + K(1 + \sqrt{t})) \right] < \infty.$$

(e) Given $\varepsilon > 0$ let $\tau_\beta(\varepsilon) = \max \{t : t\Phi(\nu, \alpha^*(\nu)) \geq \beta_t(\delta) + \varepsilon t\}$ and

$$u(\varepsilon) = \sup_{(\xi, x)} \{ \Phi(\xi, x) : d(\xi, \nu) \leq \varepsilon, \|x - \alpha^*(\nu)\|_\infty \leq \varepsilon \}.$$

Then for $t \geq \max\{\tau_\nu(\varepsilon), \tau_T(\varepsilon), \tau_\beta(u(\varepsilon))\}$ it holds that

$$t\Phi(\hat{\nu}_t, \alpha^*(\hat{\nu}_t)) \geq t(\Phi(\nu, \alpha^*(\nu)) - u(\varepsilon)) \geq \beta_t(\delta),$$

which implies that

$$\tau_\delta \leq \max\{\tau_\nu(\varepsilon), \tau_T(\varepsilon), \tau_\beta(u(\varepsilon))\} \leq \tau_\nu(\varepsilon) + \tau_T(\varepsilon) + \tau_\beta(u(\varepsilon)).$$

Taking the expectation,

$$\mathbb{E}[\tau_\delta] \leq \mathbb{E}[\tau_\nu(\varepsilon)] + \mathbb{E}[\tau_T(\varepsilon)] + \mathbb{E}[\tau_\beta(u(\varepsilon))].$$

Taking the limit as $\delta \rightarrow 0$ and using the previous parts shows that for any sufficiently small $\varepsilon > 0$,

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\beta(u(\varepsilon))]}{\log(1/\delta)} = \frac{1}{\Phi(\nu, \alpha^*(\nu)) - u(\varepsilon)}.$$

Continuity of Φ at $(\nu, \alpha^*(\nu))$ ensures that $\lim_{\varepsilon \rightarrow 0} u(\varepsilon) = 0$ and the result follows since $c^*(\nu) = 1/\Phi(\nu, \alpha^*(\nu))$.

33.8 For the first part we need to prove

$$H_2(\mu) \leq H_1(\mu) \leq (1 + \log(K))H_2(\mu)$$

where

$$H_1(\mu) = \sum_{i=1}^K \frac{1}{\Delta_i^2}, \quad H_2(\mu) = \max_{i \in [K]} \frac{i}{\Delta_i^2},$$

and where $\Delta_1 = \Delta_2 \leq \dots \leq \Delta_K$. For the first inequality note that $1/\Delta_2^2 \geq \dots \geq 1/\Delta_K^2$. Summing up the first i of these numbers, we see that

$H_1(\mu) \geq \sum_{j=1}^i 1/\Delta_j^2 \geq i/\Delta_i^2$. Taking the max of both sides over i , we get $H_1(\mu) \geq H_2(\mu)$. To get the reverse inequality, note that

$$H_1(\mu) = \frac{1}{2} \frac{2}{\Delta_1^2} + \sum_{i=2}^K \frac{1}{i} \frac{i}{\Delta_i^2} \leq \left(\frac{1}{2} + \sum_{i=2}^K \frac{1}{i} \right) \max_{j \geq 2} \frac{j}{\Delta_j^2} \leq (1 + \log(K)) \max_{j \geq 1} \frac{j}{\Delta_j^2}.$$

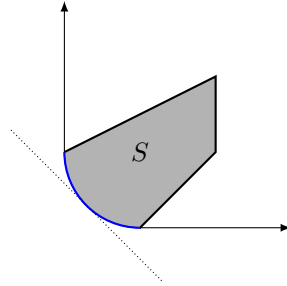
For the second part note that whenever $\Delta_1 = \dots = \Delta_K > 0$ it holds that $H_1(\mu) = H_2(\mu)$. Further, when $\Delta_i = \sqrt{i}c$ with $c > 0$ for $i \geq 2$, $i/\Delta_i^2 = 1/c$, hence $H_1(\mu) = (1/2 + \sum_{i=2}^K \frac{1}{i}) \frac{1}{c} = L \max_i \frac{i}{\Delta_i^2}$.

33.10 We have $\mathbb{P}(\max_{i \in [n]} \mu(X_i) < \mu_\alpha^*) = \mathbb{P}(\mu(X_1) < \mu_\alpha^*)^n \leq (1 - \alpha)^n \leq \delta$. Solving for n gives the required inequality.

34.7 Let \tilde{S} be the extension of S by adding rays in the positive direction: $\tilde{S} = \{x + u : x \in S, u \geq 0\}$. Clearly \tilde{S} remains convex and $\lambda(S) \subseteq \partial \tilde{S}$ is on the boundary (see figure). Let $x \in \lambda(S)$. By the supporting hyperplane theorem the convexity of \tilde{S} there exists a vector $a \in \mathbb{R}^N$ and $b \in \mathbb{R}$ such that $\langle a, x \rangle = b$ and $\langle a, y \rangle \geq b$ for all $y \in \tilde{S}$. Furthermore $a \geq 0$ since $x + e_i \in \tilde{S}$ and so $\langle x + e_i, a \rangle = b + a_i \geq b$. Define $q(\nu_i) = a_i/\|a\|_1$. Then for any policy π ,

$$\sum_{\nu \in \mathcal{E}} q(\nu) \ell(\nu, \pi) = \frac{1}{\|a\|_1} \sum_{i=1}^N a_i \ell(\nu_i, \pi) = \frac{\langle a, \ell(\pi) \rangle}{\|a\|_1} \geq b$$

with equality for the policy π with $\ell(\pi) = x$.



34.8 Let Π be the space of all policies and $\Pi_D = (e_1, \dots, e_N)$ the space of all deterministic policies, which is finite. An arbitrary policy π is uniquely represented by a convex combination of deterministic policies and from now on we identify policies Π with probability vectors in $\mathcal{P}(\Pi_D)$. This identification leads to a natural bijection between Π and $\mathcal{P}(\Pi_D)$ from which Π inherits a metric and topology. Even more straightforwardly, \mathcal{E} is identified with $[0, 1]^K$ and inherits a metric from that space. As metric spaces both Π and \mathcal{E} are compact and closed and the regret $R_n(\nu, \pi)$ is continuous in both arguments. Let $(\nu_k)_{k=1}^\infty$ be a sequence of bandit environments that is dense in \mathcal{E} and $\mathcal{E}_k = \{\nu_1, \dots, \nu_k\}$ and $S_k = \{(R_n(\nu_1, \pi), \dots, R_n(\nu_k, \pi)) : \pi \in \Pi\} \subset \mathbb{R}^k$, which is closed and convex. Now fix a policy π . By the same argument as the previous exercise for each k there exists a prior q_k supported on \mathcal{E}_k and Bayesian optimal policy with respect to

q_k such that $R_n(\nu, \pi_k) \leq R_n(\nu, \pi)$ for all $\nu \in \mathcal{E}_k$. By compactness the sequence $((\pi_k, \text{Supp}(\pi_k)))_k$ contains a subsequence converging to some policy $\pi^* \in \Pi$ and $\text{Supp}(\pi^*) \subseteq [N]$. Since $(\nu_k)_{k=1}^\infty$ is dense in \mathcal{E} it follows from the continuity of the regret that $R_n(\nu, \pi^*) \geq R_n(\nu, \pi)$ for all $\nu \in \mathcal{E}$. It remains to show that π^* is Bayesian optimal with respect to some prior. The space of posteriors is not compact, so we cannot guarantee that $(q_k)_k$ contains a convergent subsequence. Fortunately there is another way. Let k be such that $\text{Supp}(\pi_k) = \text{Supp}(\pi^*)$. Then

$$\begin{aligned} \sum_{\nu \in \mathcal{E}_k} q_k(\nu) R_n(\nu, \pi_k) &= \sum_{i=1}^N \pi_k(i) \sum_{\nu \in \mathcal{E}_k} q_k(\nu) R_n(\nu, e_i) \\ &= \sum_{i=1}^N \pi^*(i) \sum_{\nu \in \mathcal{E}_k} q_k(\nu) R_n(\nu, e_i). \end{aligned}$$

Hence π^* is Bayesian optimal with respect to q_k .

35.1 Let π be the policy of MOSS from Chapter 9, which for any 1-subgaussian bandit ν with rewards in $[0, 1]$ satisfies

$$R_n(\pi, \nu) \leq C \min \left\{ \sqrt{Kn}, \frac{K \log(n)}{\Delta_{\min}(\nu)} \right\},$$

where $\Delta_{\min}(\nu)$ is the smallest positive suboptimality gap. Let \mathcal{E}_n be the set of bandits in \mathcal{E} for which there exists an arm i with $\Delta_i \in (0, n^{-1/4})$. Then, for $C' = CK$,

$$\begin{aligned} \text{BR}_n^*(\mathbb{Q}) &\leq \text{BR}_n(\pi, \mathbb{Q}) \\ &= \int_{\mathcal{E}} R_n(\pi, \nu) d\mathbb{Q}(\nu) \\ &= \int_{\mathcal{E}_n} R_n(\pi, \nu) d\mathbb{Q}(\nu) + \int_{\mathcal{E}_n^c} R_n(\pi, \nu) d\mathbb{Q}(\nu) \\ &\leq C' \sqrt{n} \mathbb{Q}(\mathcal{E}_n) + C' \int_{\mathcal{E}_n^c} n^{1/4} \log(n) d\mathbb{Q}(\nu) \\ &= C' \sqrt{n} \mathbb{Q}(\mathcal{E}_n) + o(\sqrt{n}). \end{aligned}$$

The first part follows since $\cap_n \mathcal{E}_n = \emptyset$ and thus $\lim_{n \rightarrow \infty} \mathbb{Q}(\mathcal{E}_n) = 0$ for any measure \mathbb{Q} . For the second part we describe roughly what needs to be done. The idea is to make use of the minimax lower bound technique in Exercise 15.2, which shows that for a uniform prior concentrated on a finite set of K bandits the regret is $\Omega(\sqrt{Kn})$. The only problems are that (a) the rewards were assumed to be Gaussian and (b) the prior depends on n . The first issue is corrected by replacing the Gaussian distributions with Bernoulli distributions with means close to $1/2$. For the second issue you should compute this prior for $n \in \{1, 2, 4, 8, \dots\}$ and denote them $\mathbb{Q}_1, \mathbb{Q}_2, \dots$. Then let $\mathbb{Q} = \sum_{k=1}^\infty p_k \mathbb{Q}_k$ where $p_k \propto (k \log^2(k))^{-1}$. The result follows easily.

35.2 Recall that $E_t = U_t$ and for $t < n$,

$$E_t = \max\{U_t, \mathbb{E}[E_{t+1} \mid \mathcal{F}_t]\}.$$

Integrability of $(U_t)_{t=1}^n$ ensures that $(E_t)_{t=1}^n$ are integrable. By definition $E_t \geq \mathbb{E}[E_{t+1} \mid \mathcal{F}_t]$. Hence $(E_t)_{t=1}^n$ is a supermartingale adapted to \mathbb{F} . Hence for any stopping time $\kappa \in \mathfrak{R}_1^n$ the optional stopping theorem says that

$$\mathbb{E}[U_\kappa] \leq \mathbb{E}[E_\kappa] \leq E_1.$$

On the other hand, for τ satisfying the requirements of the lemma the process $M_t = E_{t \wedge \tau}$ is a martingale and hence $\mathbb{E}[U_\tau] = \mathbb{E}[M_\tau] = M_1 = E_1$.

35.4 Clearly $v_\gamma(x)$ is monotone nonincreasing in γ . Let

$$g = \sup_{\tau \geq 2} \frac{\mathbb{E}_x \left[\sum_{t=1}^{\tau-1} \alpha^t r(S_t) \right]}{\mathbb{E}_x \left[\sum_{t=1}^{\tau-1} \alpha^t \right]}.$$

For any $\varepsilon > 0$ there exists a stopping time $\tau \geq 2$ such that

$$\mathbb{E}_x \left[\sum_{t=1}^{\tau-1} \alpha^t r(S_t) \right] > (g - \varepsilon) \mathbb{E}_x \left[\sum_{t=1}^{\tau-1} \alpha^t \right],$$

which implies that

$$\mathbb{E}_x \left[\sum_{t=1}^{\tau-1} \alpha^t (r(S_t) - g + \varepsilon) \right] > 0.$$

And hence $v_{\gamma-\varepsilon}(x) > 0$ and so $\inf \{\gamma \in \mathbb{R} : v_\gamma(x) = 0\} > g - \varepsilon$. Since this holds for any ε it follows that $\inf \{\gamma \in \mathbb{R} : v_\gamma(x) = 0\} \geq g$. Similarly for any stopping time τ and $\varepsilon > 0$ we have

$$\mathbb{E}_x \left[\sum_{t=1}^{\tau-1} \alpha^t (r(S_t) - g - \varepsilon) \right] < 0,$$

which implies that $\sup \{\gamma \in \mathbb{R} : v_\gamma(x) = 0\} \leq g$ and the result is complete.

36.3 We have

$$\begin{aligned} \sum_{t \in \mathcal{T}} \mathbb{I}\{A_t = i\} &\leq \sum_{t=1}^n \sum_{s=1}^n \mathbb{I}\{T_i(t) = s, T_i(t-1) = s-1, G_i(T_i(t-1)) > 1/n\} \\ &= \sum_{s=1}^n \mathbb{I}\{G_i(s-1) > 1/n\} \sum_{t=1}^n \mathbb{I}\{T_i(t) = s, T_i(t-1) = s-1\} \\ &= \sum_{s=1}^n \mathbb{I}\{G_i(s-1) > 1/n\}, \end{aligned}$$

where the first equality uses that when $A_t = i$, $T_i(t) = s$ and $T_i(t-1) = s-1$ for some $s \in [n]$ and that $t \in \mathcal{T}$ implies $G_i(T_i(t-1)) > 1/n$. The next equality is by algebra, and the last follows because for any $s \in [n]$, there is at most one time point $t \in [n]$ such that $T_i(t) = s$ and $T_i(t-1) = s-1$.

For the next inequality, note that

$$\begin{aligned} \mathbb{E} \left[\sum_{t \notin \mathcal{T}} \mathbb{I} \{E_i^c(t)\} \right] &= \sum_t \mathbb{E}[\mathbb{E}[\mathbb{I} \{E_i^c(t), G_i(T_i(t-1)) \leq 1/n\} | \mathcal{F}_{t-1}]] \\ &= \sum_t \mathbb{E}[\mathbb{I} \{G_i(T_i(t-1)) \leq 1/n\} G_i(T_i(t-1))] \\ &\leq \sum_t \mathbb{E}[\mathbb{I} \{G_i(T_i(t-1)) \leq 1/n\} 1/n] \\ &= \mathbb{E} \left[\sum_{t \notin \mathcal{T}} 1/n \right], \end{aligned}$$

where the second equality used that $\mathbb{I} \{G_i(T_i(t-1)) \leq 1/n\}$ is \mathcal{F}_{t-1} -measurable and $\mathbb{E}[\mathbb{I} \{E_i^c(t)\} | \mathcal{F}_{t-1}] = 1 - \mathbb{P}(\theta_i(t) \leq \mu_1 - \varepsilon | \mathcal{F}_{t-1}) = G_i(T_i(t-1))$.

36.4

- (a) Let $f(y) = \sqrt{s/(2\pi)} \exp(-sy^2/2)$ be the probability density function of a centered Gaussian with variance $1/s$ and $F(y) = \int_{-\infty}^y f(x)dx$ be its cumulative distribution function. Then

$$\begin{aligned} Q_{1s} &= \int_{\mathbb{R}} f(y + \varepsilon)F(y)/(1 - F(y))dy \\ &\leq \frac{1}{2} \int_0^\infty f(y + \varepsilon)/(1 - F(y)) + 2 \int_{-\infty}^0 f(y + \varepsilon)F(y)dy. \end{aligned} \quad (.43)$$

For the first term in Eq. (.43) we use the following bound on $1 - F(y)$ for $y \geq 0$:

$$1 - F(y) \geq \frac{\exp(-sy^2/2)}{y\sqrt{s} + \sqrt{sy^2 + 4}}.$$

Hence

$$\begin{aligned} \frac{1}{2} \int_0^\infty \frac{f(y + \varepsilon)}{1 - F(y)} dy &\leq \frac{1}{2} \int_0^\infty f(y + \varepsilon) \exp(sy^2/2) (y\sqrt{s} + \sqrt{sy^2 + 4}) dy \\ &\leq \exp(-s\varepsilon^2/2) \int_0^\infty \exp(-sy\varepsilon) (y\sqrt{s} + 1) \sqrt{\frac{s}{2\pi}} dy \\ &= \frac{1 + \varepsilon\sqrt{s}}{\varepsilon^2 s \sqrt{2\pi}} \exp(-s\varepsilon^2/2). \end{aligned}$$

For the second term in Eq. (.43),

$$2 \int_{-\infty}^0 f(y + \varepsilon)F(y)dy \leq 2 \int_{\mathbb{R}} f(y + \varepsilon)F(y)dy \leq 2 \exp(-s\varepsilon^2).$$

Summing from $s = 1$ to ∞ shows that

$$\sum_{s=1}^{\infty} \left(2 \exp(-s\varepsilon^2) + \frac{1 + \varepsilon\sqrt{s}}{\varepsilon^2 s \sqrt{2\pi}} \exp(-s\varepsilon^2/2) \right) \leq \frac{c}{\varepsilon^2} \log \left(\frac{1}{\varepsilon} \right),$$

where the last line follows from a Mathematica slog.

- (b) Let $\hat{\mu}_{is}$ be the empirical mean of arm i after s observations. Then $Q_{is} \leq 1/n$ if

$$\hat{\mu}_{is} + \sqrt{\frac{2 \log(n)}{s}} \leq \mu_1 - \varepsilon.$$

Hence for $s \geq u = \frac{2 \log(n)}{(\Delta_i - \varepsilon)}$ we have

$$\begin{aligned} \mathbb{P}(Q_{is} > 1/n) &\leq \mathbb{P}\left(\hat{\mu}_{is} + \sqrt{\frac{2 \log(n)}{s}} > \mu_1 - \varepsilon\right) \\ &= \mathbb{P}\left(\hat{\mu}_{is} - \mu_i > \Delta_i - \varepsilon - \sqrt{\frac{2 \log(n)}{s}}\right) \\ &\leq \exp\left(-\frac{s \left(\Delta_i - \varepsilon - \sqrt{\frac{2 \log(n)}{s}}\right)^2}{2}\right). \end{aligned}$$

Summing,

$$\begin{aligned} \sum_{s=1}^n \mathbb{P}(Q_{is} > 1/n) &\leq u + \sum_{s=[u]}^n \exp\left(-\frac{s \left(\Delta_i - \varepsilon - \sqrt{\frac{2 \log(n)}{s}}\right)^2}{2}\right) \\ &\leq 1 + \frac{2}{(\Delta_i - \varepsilon)^2} (\log(n) + \sqrt{\pi \log(n)} + 1), \end{aligned}$$

where the last inequality follows by bounding the sum by an integral as in the proof of Lemma 8.1.

- 36.7** Let $\Delta_t = X_{tA^*} - X_{tA_t}$ and $I_t = I_{\mathbb{P}_{t-1}}(A^*; A_t, X_t)$. Then information directed sampling chooses $A_t \sim \pi_t$ where π_t minimizes $\Gamma_t = \mathbb{E}_{t-1}[\Delta_t]^2 / I_t$.

$$\begin{aligned} \text{BR}_n &= \mathbb{E} \left[\sum_{t=1}^n (X_{tA^*} - X_{tA_t}) \right] = \mathbb{E} \left[\sum_{t=1}^n \Delta_t \right] = \mathbb{E} \left[\sum_{t=1}^n \sqrt{\Gamma_t I_t} \right] \\ &\leq \mathbb{E} \left[\sqrt{n \sum_{t=1}^n \Gamma_t I_t} \right] \leq \sqrt{n \mathbb{E} \left[\sum_{t=1}^n \Gamma_t I_t \right]}. \end{aligned}$$

Since π_t is chosen to minimize Γ_t it follows by Lemma 36.3 that $\Gamma_t \leq K/2$ almost surely. Hence

$$\text{BR}_n \leq \sqrt{\frac{nK}{2} \mathbb{E} \left[\sum_{t=1}^n I_t \right]} \leq \sqrt{\frac{nK \log(K)}{2}}.$$

37.8

- (a) $A_n \mathbf{1} = \frac{1}{n} \sum_{t=0}^{n-1} P^t \mathbf{1} = \mathbf{1}$, which means that A_n is right stochastic.

- (b) This follows immediately from the definitions.
- (c) Let (B_n) and (C_n) be convergent subsequences of (A_n) with $\lim_{n \rightarrow \infty} B_n = B$ and $\lim_{n \rightarrow \infty} C_n = C$. It suffices to show that $B = C$. From Part (b), $B_m + \frac{1}{n_m}(P^{n_m} - I) = B_m P = P B_m$. Taking the limit as m tends to infinity and using the fact that P^n is $[0, 1]$ -valued we see that $B = B P = P B$. Similarly, $C = C P = P C$ and it follows that $B = B P^i = P^i B$ and $C = C P^i = P^i C$ hold for any $i \geq 0$. Hence $B = B C_m = C_m B$ and $C = C B_m = B_m C$ for any $m \geq 1$. Taking limit as m tends to infinity shows that $B = B C = C B$ and $C = C B = B C$, which together imply that $B = C$.
- (d) We have already seen in the proof of Part (c) that $P^* = P^* P = P P^*$. From this, it follows that $P^* = P^* P^i$ for any $i \geq 0$, which implies that $P^* = P^* A_n$ holds for any $n \geq 1$. Taking limit shows that $P^* = P^* P^*$.
- (e) Let $\nu^\top = e_1^\top P^*$. Then by the previous part $\nu^\top P = e_1^\top P^* P = e_1^\top (P^* P) = e_1^\top P^* = \nu^\top$.
- (f) Abbreviate $P = Q_t$. First notice that for each $a \in \mathcal{A}$ and $b \in [K]$ there exists an $i \in \mathbb{N}^+$ such that $P_{ba}^i > 0$. Similarly, for any $c \in [K]$ and $b \notin \mathcal{A}$ we have $P_{cb}^* = 0$. Then since $P^* = P^* P^i$ it follows that for $a \in \mathcal{A}$ we have $P_{ba}^* > 0$ for all $b \in [K]$. Suppose ν is a stationary distribution and $b \notin \mathcal{A}$. Then $\nu_b = \sum_a P_{ab}^* \nu_a = 0$ and hence the stationary distribution is unique outside of \mathcal{A} . Now suppose that μ is another stationary distribution and $\mu \neq \nu$. Let $a = \operatorname{argmax}_{\nu_a - \mu_a}$, which means that $\nu_a > \mu_a$ and so $a \in \mathcal{A}$. Since both are stationary we have $\nu_a - \mu_a = \sum_b P_{ba}^* (\nu_b - \mu_b)$ and hence $\nu_b = \mu_b$ whenever $P_{ba}^* > 0$. But this means that $\|\nu\|_1 > \|\mu\|_1$, which is a contradiction.

37.9 First we show that P_t is indeed a probability vector. By assumption \tilde{P}_t is the stationary distribution, which is a probability distribution. Let $\bar{P}_t = \text{REDISTRIBUTE}(\tilde{P}_t)$ so that

$$P_t = (1 - \gamma)\bar{P}_t + \frac{\gamma}{K}\mathbf{1},$$

which means we need to show that \bar{P}_t is a probability distribution. Since \bar{P}_t is obtained by the iterative procedure given in the REDISTRIBUTE function it is sufficient to show that the vector q tracked by this algorithm is indeed a distribution. The claim is that each loop of the REDISTRIBUTE function does not break this property. The first observation is that the algorithm always moves mass from actions in \mathcal{A} to actions in \mathcal{D} . All that must be shown is that $\bar{P}_{ta} \geq 0$ for all $a \in \mathcal{A}$. To see this note first that if $a \in \mathcal{A}$ is one of the choices of the algorithm in Line 11, then $\rho c_a q_a \leq p_a / (2K)$ and so

$$\bar{P}_{ta} \geq \tilde{P}_{ta} / 2 \quad \text{for all } a \in \mathcal{A} \geq 0. \quad (.44)$$

Part (a): Since $\gamma \leq 1/2$ this follows from Eq. (.44).

Part (b): First we show that $\sum_{a \in [K]} (\bar{P}_{ta} - \tilde{P}_{ta}) \ell_a = 0$. It suffices to show that the redistribution in each inner loop of the algorithm does not change this value,

which is true because

$$\begin{aligned} (c_a q_a + c_b q_b) \ell_d &= (c_a q_a + c_b q_b) (\alpha \ell_a + (1 - \alpha) \ell_b) \\ &= \frac{q_a q_b}{\alpha q_b + (1 - \alpha) q_a} (\alpha \ell_a + (1 - \alpha) \ell_b) \\ &= \rho c_a q_a \ell_a + \rho c_b q_b \ell_b. \end{aligned}$$

Then using the definition of P_t we have

$$\begin{aligned} \left| \sum_{a \in [K]} (P_{ta} - \tilde{P}_{ta}) \langle \ell_a, u \rangle \right| &= \left| \sum_{a \in [K]} (P_{ta} - \bar{P}_{ta}) \langle \ell_a, u \rangle \right| \\ &= \gamma \left| \sum_{a \in [K]} \left(\frac{1}{K} - \bar{P}_{ta} \right) \langle \ell_a, u \rangle \right| \\ &\leq \gamma, \end{aligned}$$

where we used the assumption that $\ell_a \in [0, 1]^E$ for all actions and $u \in \mathcal{P}_{E-1}$ so that $\langle \ell_a, u \rangle \in [0, 1]$.

(c): There are three cases: Either $b = k$ or $b = a$ or b is degenerate. If $b = k$, then the result is immediate from Part (a). If $b = a$, then, Part (a) combined with (37.10) implies that $P_{tb} = P_{ta} \geq \tilde{P}_{ta}/4 \geq \tilde{P}_{tk} Q_{tka}/4 \geq \tilde{P}_{tk} Q_{tka}/(4K)$. Finally, if b is degenerate, then by the definition of the rebalancing algorithm we have

$$\bar{P}_{tb} \geq \frac{\min(\tilde{P}_{tk}, \tilde{P}_{ta})}{2K} \geq \frac{\min(\tilde{P}_{tk}, \tilde{P}_{tk} Q_{tka})}{2K} = \frac{\tilde{P}_{tk} Q_{tka}}{2K}$$

and the result follows from Eq. (.44).

(d): This is trivial from the definition of P_t .

(e): Let $a, b \in \mathcal{A}$ be the Pareto optimal actions that share their probability with d in REDISTRIBUTE. Thus, $d \in \mathcal{N}_{ab}$. Since d is a duplicate of $k \in \mathcal{A}$, and \mathcal{A} is duplicate-free, it must be that k is either equal to a or b . With no loss of generality assume that $k = a$. Then, $\ell_d = \ell_a$ and it follows that $\alpha = 1$ and so $c_a = 1$ and $c_b = 0$. This means that $\bar{P}_{td} = \tilde{P}_{ta}/(2K)$ and using Eq. (.44) again yields the result.

37.10

(a) This follows immediately from Part (e) of Lemma 37.3.

(b) There are two cases: Either b is a duplicate of a , or it is not a duplicate action of either k or a . If b is a duplicate of a , Part (e) of Lemma 37.3 gives that $P_{tb} \geq \tilde{P}_{ta}/(4K)$. Hence,

$$\frac{Q_{tka} \tilde{P}_{tk}}{P_{tb}} \leq 4K \frac{Q_{tka} \tilde{P}_{tk}}{\tilde{P}_{ta}} \leq 4K \frac{\sum_k Q_{tka} \tilde{P}_{tk}}{\tilde{P}_{ta}} = 4K \frac{\tilde{P}_{ta}}{\tilde{P}_{ta}} = 4K,$$

where the second to last equality follows from the definition of \tilde{P}_t . In the remaining case, Part (c) of Lemma 37.3 gives the result immediately.

- (c) That S and S_a are disjoint is clear. Next let a and a' be distinct elements of $\mathcal{N}_k \cap \mathcal{A}$ and suppose that $b \in \mathcal{N}_{ak} \cap \mathcal{N}_{a'k}$. By Lemma 37.1(b), either $b = k$ or $C_b = C_a \cap C_k = C_{a'} \cap C_k$. This is impossible, however, since k , a and a' are all distinct Pareto optimal actions.

37.12

- (a) We only sketch the argument. Let $\mathcal{V} = \{a, b\} \times [F]$ be the vertices of a graph and let $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$ be the set of edges where $((a, f), (b, f')) \in \mathcal{E}$ if there exists an $i \in [E]$ such that $\Phi_{ai} = f$ and $\Phi_{bi} = f'$, which makes $(\mathcal{V}, \mathcal{E})$ a bipartite graph. We assume without loss of generality this graph is fully connected, since if not, then the following argument can be applied to each connected component. Let $f, f' \in [F]$ and use the fact there exists a path between (a, f) and (a, f') to show that $v(a, f) - v(a, f') \leq 2F$. By translating $v(a, f)$ in one direction and $v(b, f)$ in the other we may choose v so that $\max_f |v(a, f)| \leq F$. Conclude that $\max_f |v(b, f)| \leq F + 1$ using the fact that $v(b, f) + v(a, f') = \ell_{ai} - \ell_{bi}$ for some $f' \in [F]$ and $i \in [E]$.
- (b) Let $\ell_1 = (1, 0, 0, 0, \dots, 0, 0)$ and $\ell_2 = (0, 1, 0, 0, \dots, 0, 0)$. For all other actions $a > 2$ let $\ell_a = \mathbf{1}$. Hence $\ell_1 - \ell_2 = (1, -1, 0, 0, \dots, 0, 0)$. Next we define the feedback matrix by

$$\Phi = \begin{pmatrix} 1 & 2 & 2 & 1 & 2 & 1 & 2 & & \\ 1 & 1 & 2 & 2 & 1 & 2 & 1 & \ddots & \\ 1 & 1 & 1 & 2 & 2 & 1 & 2 & \ddots & \\ 1 & 1 & 1 & 1 & 2 & 2 & 1 & \ddots & \\ 1 & 1 & 1 & 1 & 1 & 2 & 2 & \ddots & \\ 1 & 1 & 1 & 1 & 1 & 1 & 2 & \ddots & \\ & & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \end{pmatrix}.$$

For this choice we have $v(1, 2) = v(1, 1) - 2$ and $v(2, 2) = v(2, 1) + 2$ and $v(3, 2) = v(3, 1) - 2$. For $k > 3$ we have $|v(k, 2) - v(k, 1)| = |v(k-2, 2) - v(k-2, 1)| + |v(k-1, 2) - v(k-1, 1)|$. Hence $|v(k, 2) - v(k, 1)|$ grows exponentially by the usual techniques for evaluating Fibonacci sequences.

38.2 The solution to Part (b) is immediate from Part (a), so we only show the solution to Part (a). Abbreviate \mathbb{P}_μ^{π} to \mathbb{P} and $\mathbb{P}_\mu^{\pi'}$ to \mathbb{P}' . Let $\pi' = (\pi'_1, \pi'_2, \dots)$ be the Markov policy to be constructed: For each $t \geq 1$, $s \in \mathcal{S}$, $\pi'_t(\cdot|s)$ is a distribution over \mathcal{A} .

Fix $(s, a) \in \mathcal{S} \times \mathcal{A}$ and consider first $t = 1$. We want $\mathbb{P}'(S_1 = s, A_1 = a) = \mathbb{P}(S_1 = s, A_1 = a)$. By the definition of \mathbb{P}' , \mathbb{P} and by the definition of conditional probabilities, we have $\mathbb{P}'(S_1 = s, A_1 = a) = \mathbb{P}'(A_1 = a|S_1 = s)\mathbb{P}'(S_1 = s) = \pi'_1(a|s)\mu(s) = \mathbb{P}(S_1 = a, A_1 = a) = \mathbb{P}(A_1 = a|S_1 = s)\mu(s)$. Thus, defining $\pi'_1(a|s) = \mathbb{P}(A_1 = a|S_1 = s)$ we see that the desired equality holds for $t = 1$. Now,

for $t > 1$ first notice that from $\mathbb{P}'(A_{t-1} = a, S_{t-1} = s) = \mathbb{P}'(A_{t-1} = a, S_{t-1} = s)$, $(s, a) \in \mathcal{S} \times \mathcal{A}$ it follows by summing these equations over a that $\mathbb{P}'(S_{t-1} = s) = \mathbb{P}'(S_{t-1} = s)$. Hence, the same calculation as for $t = 1$ applies, showing that $\pi'_t(a|s) = \mathbb{P}(A_t = a|S_t = s)$ will work.

38.4 We show that $D(M) < \infty$ implies that M is strongly connected and that $D(M) = \infty$ implies that M is not strongly connected. Assume first that $D(M) < \infty$. Take any $s, s' \in \mathcal{S}$. Assume first that $s \neq s'$. By definition, there is a policy whose expected travel time from state s to s' is finite. Take this policy. It follows that this policy reaches state s' from state s with positive probability, because otherwise the expected travel time would be infinite. Formally, if T is the random travel time of a policy whose expected travel time between s and s' is finite, $\{T = \infty\}$ is the event that the policy does not reach state s' . Now, for any $n \in \mathbb{N}$, $T > n\mathbb{1}\{T = \infty\}$. Taking expectations and reordering gives $\mathbb{E}[T]/n > \mathbb{P}(T = \infty)$. Letting $n \rightarrow \infty$, we see that $\mathbb{P}(T = \infty) = 0$ (thus we see that the policy reaches state s' in fact with probability one). It remains to consider the case when $s = s'$. If the MDP has a single state, it is strongly connected by definition. Otherwise, there exist a state $s'' \in \mathcal{S}$ that is distinct from $s = s'$. Since $D(M)$ is finite, there is a policy that reaches s' from s with positive probability and another one that reaches s again from s' with positive probability. Compose these two policies the obvious way to find the policy that travels from s to s with positive probability.

Assume now that $D(M) = \infty$, while M is strongly connected (proof by contradiction). Since M is strongly connected, for any s, s' , there is a policy that has a positive probability of reaching s' from s . But this means that the uniformly random policy (the policy which chooses uniformly at random between the actions at any state) has also a positive probability of reaching any state from any other state. We claim that the expected travel time of this policy is finite between any pairs of states. Indeed, this follows by noticing that the states under this policy form a time-homogenous Markov chain whose transition probability matrix is irreducible and the hitting times in some a Markov chain, which coincide with the expected travel times in the MDP for the said policy, are finite. [link](#)

38.5 We follow the advice of the hint. For the second part, note that the minimum in the definition of $d^*(\mu_0, U)$ is attained when n_k is maximized for small indices until $|U|$ is exhausted. In particular, if $(n_k)_{0 \leq k \leq m}$ denotes the optimal solution ($n_k = 0$ for $k > m$) then $n_0 = A^0, \dots, n_{m-1} = A^k$, $0 \leq n_m = |U| - \sum_{k=0}^{m-1} A^k (= |U| - \frac{A^m - 1}{A - 1}) < A^m$. Hence, $|U| < A^m + \frac{A^m - 1}{A - 1} \leq 2A^m$,

implying that $m \geq \log_A(|U|/2)$. Thus,

$$\begin{aligned}
 d^*(\mu_0, U) &= \sum_{k=0}^{m-1} k A^k + m n_m \\
 &= m|U| + \sum_{k=0}^{m-1} (k-m)A^k \\
 &\stackrel{(a)}{=} m|U| + \frac{m}{A-1} - \frac{A^{m+1} - A}{(A-1)^2} \\
 &\geq |U| \left(m-1 - \frac{1}{A-1} \right) \\
 &\geq |U|(\log_A(|U|) - 3),
 \end{aligned}$$

where step (a) follows since $|U| < \frac{A^m - 1}{A-1}$. Choosing $U = \mathcal{S}$, we see that the expected minimum time to reach a random state in \mathcal{S} is lower bounded by $\log_A(\mathcal{S}) - 3$. The expected minimum time to reach an arbitrary state in \mathcal{S} must also be above this quantity, proving the desired result.

38.8

- (a) $A_n \mathbf{1} = \frac{1}{n} \sum_{t=0}^{n-1} P^t \mathbf{1} = \mathbf{1}$, which means that A_n is right stochastic.
- (b) This follows immediately from the definitions.
- (c) Let (B_n) and (C_n) be convergent subsequences of (A_n) with $\lim_{n \rightarrow \infty} B_n = B$ and $\lim_{n \rightarrow \infty} C_n = C$. It suffices to show that $B = C$. From Part (b), $B_m + \frac{1}{n_m}(P^{n_m} - I) = B_m P = P B_m$. Taking the limit as m tends to infinity and using the fact that P^n is $[0, 1]$ -valued we see that $B = B P = P B$. Similarly, $C = C P = P C$ and it follows that $B = B P^i = P^i B$ and $C = C P^i = P^i C$ hold for any $i \geq 0$. Hence $B = B C_m = C_m B$ and $C = C B_m = B_m C$ for any $m \geq 1$. Taking limit as m tends to infinity shows that $B = B C = C B$ and $C = C B = B C$, which together imply that $B = C$.
- (d) We have already seen in the proof of Part (c) that $P^* = P^* P = P P^*$. From this, it follows that $P^* = P^* P^i$ for any $i \geq 0$, which implies that $P^* = P^* A_n$ holds for any $n \geq 1$. Taking limit shows that $P^* = P^* P^*$.
- (e) Let $B = P - P^*$. By algebra $I - B^i = (I - B)(I + B + \cdots + B^{i-1})$. Summing over $i = 1, \dots, n$ and dividing by n and using the fact that $B^i = P^i - P^*$ for all $i \geq 1$,

$$I - \frac{1}{n} \sum_{i=1}^n P^i + P^* = (I - B)H_n, \quad (.45)$$

where $H_n = \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{i-1} (P - P^*)^k$. The limit of the left-hand side of (.45) exists and is equal to the identity matrix I . Hence the limit of the right-hand side also exists and in particular the limit of H_n must exist. Denoting this by H_∞ we find that $I = (I - B)H_\infty$ and thus $I - B$ is invertible and its inverse H is equal to H_∞ .

(f) Let $D_n = \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{i-1} (P^k - P^*)$. Then

$$B^k = \begin{cases} I, & \text{if } k = 0; \\ P^k - P^*, & \text{otherwise.} \end{cases}$$

Using this we calculate $H_n - D_n = \frac{1}{n} \sum_{i=1}^n (P - P^*)^0 - \frac{1}{n} \sum_{i=1}^n (P^0 - P^*) = I - I + P^* = P^*$. Hence $H - \lim_{n \rightarrow \infty} D_n = P^*$. From the definition of D we have $D = \lim_{n \rightarrow \infty} D_n$.

(g) This follows immediately because

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{i-1} P^k (r - \rho) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{i-1} (P^k - P^*) r = Dr.$$

(h) One way to prove this is to note that by the previous part $v = Dr = (H - P^*)r$, hence $r = H^{-1}(v + \rho)$. Now, $H^{-1}(v + \rho) = (I - P + P^*)(v + \rho) = v - Pv + P^*v + \rho - P\rho + P^*\rho = v - Pv + \rho$, where we used that $P^*v = P^*Dr$ and $P^*D = P^*(H - P^*) = P^*H - P^* = (P^* - P^*H^{-1})H = 0$ and that $P\rho = PP^*r = P^*r = P^*\rho = P^*P^*r$.

Alternatively, the following direct argument also works. In this argument we only use that v is well-defined. Let $v_n = \sum_{t=0}^{n-1} P^t(r - \rho)$, $\bar{v}_n = \frac{1}{n} \sum_{i=1}^n v_i$. Note that $\lim_{n \rightarrow \infty} \bar{v}_n = v$. Then, $v_{k+1} = Pv_k + (r - \rho)$. Taking the average of these over $k = 1, \dots, n$ we get

$$\frac{1}{n} ((n+1)\bar{v}_{n+1} - v_1) = P\bar{v}_n + (r - \rho).$$

Taking the limit of both sides proves that $v = Pv + r - \rho$, which, after reordering gives $v + \rho = r + Pv$.

38.9

(a) First note that $|\max_x f(x) - \max_y g(y)| \leq \max_x |f(x) - g(x)|$. Then for $v, w \in \mathbb{R}^S$,

$$\begin{aligned} \|T_\gamma v - T_\gamma w\|_\infty &\leq \max_{s \in S} \max_{a \in \mathcal{A}} \gamma |\langle P_a(s), v - w \rangle| \\ &\leq \max_{s \in S} \max_{a \in \mathcal{A}} \gamma \|P_a(s)\|_1 \|v - w\|_\infty \\ &= \gamma \|v - w\|_\infty. \end{aligned}$$

Hence T is a contraction with respect to the supremum norm as required.

(b) This follows immediately from the Banach fixed point theorem, which also guarantees the uniqueness of a value function v satisfying $v = T_\gamma v$.

(c) Recall that the greedy policy is $\pi(s) = \operatorname{argmax}_a r_a(s) + \gamma \langle P_a(s), v \rangle$. Then

$$v(s) = \max_{a \in \mathcal{A}} r_a(s) + \gamma \langle P_a(s), v \rangle = r_\pi(s) + \gamma \langle P_\pi(s), v \rangle.$$

(d) We have $v = r + \gamma P_\pi v$. Solving for v completes the result.

- (e) If π is a memoryless policy, it is trivial to see that $v_\gamma^\pi = r_\pi + \gamma P_\pi v_\gamma^\pi$. Let π^* be the greedy policy with respect to v , the unique solution of $v = T_\gamma v$. By the previous part of this exercise, it follows that $v_\gamma^{\pi^*} = v$. By Exercise 38.2, it suffices to show that for any Markov policy π , $v_\gamma^\pi \leq v$. If π_t is the memoryless policy used in time step t when following π , $v_\gamma^\pi = \sum_{t=1}^{\infty} \gamma^{t-1} P_\pi^{(t-1)} r_{\pi_t}$, where $P^{(0)} = I$ and for $t \geq 1$, $P^{(t)} = P_{\pi_1} \dots P_{\pi_t}$. For $n \geq 1$, let $v_{\gamma,n}^\pi = \sum_{t=1}^n \gamma^{t-1} P_\pi^{(t-1)} r_{\pi_t}$. It is easy to see that $v_{\gamma,1}^\pi = r_{\pi_1} \leq T\mathbf{0}$. Assume that for some $n \geq 1$,

$$v_{\gamma,n}^\pi \leq T^n \mathbf{0}. \quad (.46)$$

Notice that $f \leq g$ implies $Tf \leq Tg$. Hence, $Tv_{\gamma,n}^\pi \leq T^{n+1}\mathbf{0}$. Further, $v_{\gamma,n+1}^\pi = r_{\pi_{n+1}} + \gamma P_{\pi_{n+1}} v_{\gamma,n}^\pi \leq Tv_{\gamma,n}^\pi$. This shows that (.46) holds for all $n \geq 1$. Letting $n \rightarrow \infty$, the right-hand side converges to v , while the left-hand side converges to v_γ^π . Hence, $v_\gamma^\pi \leq v$.

38.10

- (a) Let $0 \leq \gamma < 1$. Algebra gives

$$P_\gamma^* P = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P^t P = \frac{P_\gamma^* - (1 - \gamma)I}{\gamma}.$$

Hence $\gamma P_\gamma^* P = P_\gamma^* - (1 - \gamma)I$. It is easy to check that P_γ^* is right stochastic. By the compactness of the space of right stochastic matrices, $(P_\gamma^*)_\gamma$ has at least one cluster point A as $\gamma \rightarrow 1^-$. It follows that $AP = A$, which implies that $AP^* = A$. Now, $(P_\gamma^*)^{-1} P^* = \frac{(I - \gamma P)P^*}{1 - \gamma} = P^*$, which implies that $P^* = AP^* = A$. Since this holds for any cluster point we conclude that $\lim_{\gamma \rightarrow 1^-} P_\gamma^* = P^*$.

- (b) Since $I - P + P^*$ is invertible and $PP^* = P^* = P^*P^*$, the required claim is equivalent to

$$\lim_{\gamma \rightarrow 1^-} (I - P + P^*) \left(\frac{P_\gamma^* - P^*}{1 - \gamma} \right) = I - P^*.$$

Rewriting $P_\gamma^* = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P^t$ shows that

$$\begin{aligned} (I - P + P^*) \left(\frac{P_\gamma^* - P^*}{1 - \gamma} \right) &= (I - P + P^*) \left(\sum_{t=0}^{\infty} \gamma^t P^t - \frac{P^*}{1 - \gamma} \right) \\ &= \sum_{t=0}^{\infty} (I - P) \gamma^t P^t = \frac{1}{\gamma} (I - P^*). \end{aligned}$$

The result is completed by taking the limit as γ tends to one from below and using Part (a).

38.11

- (a) Let (γ_n) be an arbitrary monotone increasing sequence with $\gamma_n < 1$ for all n and $\lim_{n \rightarrow \infty} \gamma_n = 1$. Let v_n be the fixed point of T_{γ_n} and π_n be the greedy policy with respect to v_n . Since greedy policies are always deterministic and there are only finitely many deterministic policies it follows there exists a subsequence $n_1 < n_2 < \dots$ and policy π such that $\pi_{n_k} = \pi$ for all k .
- (b) For arbitrary π and $0 \leq \gamma < 1$, let v_γ^π be the value function of policy π in the γ -discounted MDP, v^π the value function of π and ρ^π its gain. Let D_π be the deviation matrix underlying P_π . Define f_γ^π by

$$v_\gamma^\pi = \frac{\rho^\pi}{1-\gamma} + v_\pi + f_\gamma^\pi. \quad (.47)$$

By Part (b) of Exercise 38.10, and because $\rho^\pi = P_\pi^* r_\pi$ and $v_\pi = D_\pi r_\pi$, it holds that $\|f_\gamma^\pi\|_\infty \rightarrow 0$ as $\gamma \rightarrow 1$.

Fix now π to be the policy whose existence is guaranteed by the previous part. By Part (e) of Exercise 38.9, π is γ_n -discount optimal for all $n \geq 1$. Suppose that ρ^π is not a constant. In any case, ρ^π is piecewise constant on the recurrent classes of the Markov chain with transition probabilities P^π . Let $\rho_*^\pi = \max_{s \in \mathcal{S}} \rho^\pi(s)$. Let $R \subset \mathcal{S}$ be the recurrent class in this Markov chain where ρ^π is the largest and take a policy π' that is identical to π over R , while π' is set up such that it gets to R with probability one. Such a π' exist because the MDP is strongly connected. Fix any $s \in \mathcal{S} \setminus R$. We claim that there exists some $\gamma^* \in (0, 1)$ such that for all $\gamma \geq \gamma^*$,

$$v_\gamma^{\pi'}(s) > v_\gamma^\pi(s). \quad (.48)$$

If this was true and n is large enough so that $\gamma_n \geq \gamma^*$ then, since π is γ_n -discount optimal, $v_{\gamma_n}^\pi(s) \geq v_{\gamma_n}^{\pi'}(s) > v_{\gamma_n}^\pi(s)$, which is a contradiction.

Hence, it remains to show (.48). By the construction of π' , $\rho^{\pi'}(s) = \rho_*^\pi > \rho^\pi(s)$. From (.47),

$$v_\gamma^{\pi'}(s) = \frac{\rho^{\pi'}(s)}{1-\gamma} + v_{\pi'}(s) + f_\gamma^{\pi'}(s) > \frac{\rho^\pi(s)}{1-\gamma} + v_\pi(s) + f_\gamma^\pi(s) = v_\gamma^\pi(s),$$

where the inequality follows by taking $\gamma \geq \gamma^*$ for some $\gamma^* \in (0, 1)$. The existence of any appropriate γ^* follows because $f_\gamma^{\pi'}(s), f_\gamma^\pi(s) \rightarrow 0$, while $1/(1-\gamma) \rightarrow \infty$ as $\gamma \rightarrow 1$.

- (c) Since v is the value function of π and ρ is its gain, by (38.4) we have $\rho \mathbf{1} + v = r_\pi + P_\pi v$. Let π' be an arbitrary stationary policy and v_n be as before. Let $f_n = f_{\gamma_n}$ where f_γ is defined by (.47). Note that $\|f_n\|_\infty \rightarrow 0$ as $n \rightarrow \infty$. Then,

$$\begin{aligned} 0 &\geq r_{\pi'} + (\gamma_n P_{\pi'} - I)v_n \\ &= r_{\pi'} + (\gamma_n P_{\pi'} - I) \left(\frac{\rho \mathbf{1}}{1-\gamma_n} + v + f_n \right) \\ &= r_{\pi'} - \rho \mathbf{1} + \gamma_n P_{\pi'} v - v + (\gamma_n P_{\pi'} - I)f_n. \end{aligned}$$

Note that when $\pi' = \pi$, the first inequality becomes an equality. Taking the limit as n tends to infinity and rearranging shows that

$$\rho \mathbf{1} + v \geq r_{\pi'} + P_{\pi'} v$$

and

$$\rho \mathbf{1} + v = r_{\pi} + P_{\pi} v. \quad (.49)$$

Since π' was arbitrary, $\rho + v(s) \geq \max_a r_a(s) + \langle P_a(s), v \rangle$ holds for all $s \in \mathcal{S}$. This combined with (.49) shows that the pair (ρ, v) satisfies the Bellman optimality equation as required.

38.12 Clearly the optimal policy is to take action STAY in any state and this policy has gain $\rho^* = 0$. Pick any solution (ρ, v) to the Bellman optimality equations. Therefore $\rho = \rho^* = 0$ by Theorem 38.1. The Bellman optimality equation for state 1 is $v(1) = \max(v(1), -1 + v(2))$, which is equivalent to $v(1) \geq -1 + v(2)$. Similarly, the Bellman optimality equation for state 2 is equivalent to $v(2) \geq -1 + v(1)$. Thus the set of solutions is a subset of

$$\{(\rho, v) \in \mathbb{R} \times \mathbb{R}^2 : \rho = 0, v(1) - 1 \leq v(2) \leq v(1) + 1\}.$$

The same argument shows that any element of this set is a solution to the optimality equations.

38.13 Let $T : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$ be defined by $(Tv)(s) = \max_a r_a(s) - \rho^* + \langle P_a(s), v \rangle$ so the Bellman optimality equation Eq. (38.5) can be written in the compact form $v = Tv$. Let $v \in \mathbb{R}^{\mathcal{S}}$ be a solution to Eq. (38.5). The proof follows from the definition of the diameter and by showing that for any states $s_1, s_2 \in \mathcal{S}$ and memoryless policy π it holds that

$$v(s_2) \leq v(s_1) + (\rho^* - \min_{s,a} r_a(s)) \mathbb{E}^{\pi}[\tau_{s_2} \mid S_1 = s_1].$$

The remainder of the proof is devoted to proving this result for fixed $s_1, s_2 \in \mathcal{S}$ and memoryless policy π . Abbreviate $\tau = \tau_{s_2}$ and let $\mathbb{E}[\cdot]$ denote the expectation with respect to the measure induced by the interaction of π and the MDP conditioned on $S_1 = s_1$. Since the result is trivial when $\mathbb{E}[\tau] = \infty$, for the remainder we assume that $\mathbb{E}[\tau] < \infty$. Define operator $\bar{T} : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$ by

$$(\bar{T}u)(s) = \begin{cases} \min_{s,a} r_a(s) - \rho^* + \langle P_{\pi}(s), u \rangle, & \text{if } s \neq s_2; \\ v(s_2), & \text{otherwise.} \end{cases}$$

Since $r_{\pi}(s) - \rho^* \geq \min_{s,a} r_a(s) - \rho^*$ and $Tv = v$ it follows that $(\bar{T}v)(s) \leq (Tv)(s) = v(s)$. Notice that for $u \leq w$ it holds that $\bar{T}u \leq \bar{T}w$. Then by induction we have $\bar{T}^n v \leq v$ for all $n \in \mathbb{N}^+$. By unrolling the recurrence we have

$$v(s_1) \geq (\bar{T}^n v)(s_1) = \mathbb{E} \left[-(\rho^* - \min_{s,a} r_a(s))(n \wedge \tau) + v(S_{\tau \wedge n}) \right].$$

Taking the limit as n tends to infinity shows that $v(s_1) \geq v(s_2) - (\rho^* - \min_{s,a} r_a(s))\mathbb{E}[\tau]$, which completes the result.

38.14

- (a) It is clear that Algorithm 25 returns TRUE if and only if (ρ, v) is feasible for Eq. (38.6). Note that feasibility can be written in the compact form $\rho\mathbf{1}v \geq Tv$. It remains to show that when (ρ, v) is not feasible then $u = (1, e_s - P_{a_s^*}(s))$ is such that for any (ρ', v') feasible, $\langle (\rho', v'), u \rangle > \langle (\rho, v), u \rangle$. For this, we have $\langle (\rho', v'), u \rangle = \rho' + v'(s) - \langle P_{a_s^*}(s), v' \rangle \geq r_{a_s^*}(s)$, where the inequality used that (ρ', v') is feasible. Further, $\langle (\rho, v), u \rangle = \rho + v(s) - \langle P_{a_s^*}(s), v \rangle < r_{a_s^*}(s)$, by the construction of u . Putting these together gives the result.
- (b) Relax the constraint that $v(\tilde{s}) = 0$ to $-\varepsilon \leq v(\tilde{s}) \leq \varepsilon$. Then add the ε of slack to the first constraint of Eq. (38.7) and add the additional constraints used in Eq. (38.9). Now the ellipsoid method can be applied as for Eq. (38.9).

38.15

- (a) Let $\pi = (\pi_1, \pi_2, \dots)$ be an arbitrary Markov policy where π_t is the policy followed in time step t . Using the notation and techniques from the proof of Theorem 38.1,

$$\begin{aligned} P_\pi^{(t-1)} r_{\pi_t} &= P_\pi^{(t-1)} (r_{\pi_t} + P_{\pi_t} v - P_{\pi_t} v) \leq P_\pi^{(t-1)} (r_{\tilde{\pi}} + P_{\tilde{\pi}} v - P_{\pi_t} v) \\ &\leq P_\pi^{(t-1)} ((\rho + \varepsilon)\mathbf{1} + v - P_{\pi_t} v) = (\rho + \varepsilon)\mathbf{1} + P_\pi^{(t-1)} v - P_\pi^{(t)} v. \end{aligned}$$

Taking the average and then the limit shows that $\bar{\rho}^\pi(s) \leq \rho + \varepsilon$ for all $s \in \mathcal{S}$. By the claim in Exercise 38.2, $\rho^* \leq \rho + \varepsilon$.

- (b) We have $r_{\tilde{\pi}}(s) + \langle P_{\tilde{\pi}(s)}(s), v \rangle \geq \max_a r_a(s) + \langle P_a(s), v \rangle - \varepsilon' \geq \rho + v(s) - (\varepsilon + \varepsilon')$. Therefore,

$$P_{\tilde{\pi}}^{t-1} r_{\tilde{\pi}} = P_{\tilde{\pi}}^{t-1} (r_{\tilde{\pi}} + P_{\tilde{\pi}} v - P_{\tilde{\pi}} v) \geq P_{\tilde{\pi}}^{t-1} ((\rho - (\varepsilon + \varepsilon'))\mathbf{1} + v - P_{\tilde{\pi}} v).$$

Taking the the average and the limit again shows that $\rho^{\tilde{\pi}}(s) \geq \rho - (\varepsilon + \varepsilon')$. The claim follows from combining this with the previous result.

38.16 Define operator $T : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$ by $(Tu)(s) = \max_{a \in \mathcal{A}} r_a(s) + \langle P_a(s), u \rangle$, which is chosen so that $v_n^* = T^n \mathbf{0}$. Let v be a solution to the Bellman optimality equation with $\min_s v(s) = 0$. Then $Tv = \rho^* \mathbf{1} + v$ and

$$v_n^* = T^n \mathbf{0} \leq T^n v = n\rho^* \mathbf{1} + v \leq n\rho^* \mathbf{1} + D\mathbf{1},$$

where the last inequality follows from the previous exercise and the assumption that $\min_s v(s) = 0$.

38.18

- (a) There are four memoryless policies in this MDP. All are optimal except the policy π that always chooses the dashed action.

(b) The optimistic rewards are given by

$$\begin{aligned}\tilde{r}_{k,\text{STAY}}(s) &= \frac{1}{2} + \sqrt{\frac{L}{2(1 \vee T_{k-1}(s, \text{STAY}))}} \\ \tilde{r}_{k,\text{GO}}(s) &= \sqrt{\frac{L}{2(1 \vee T_{k-1}(s, \text{GO}))}}.\end{aligned}$$

Whenever $T_{k-1}(s, a) \geq 1$ the transition estimates $\hat{P}_{k-1,a}(s) = P_a(s)$. Let S'_t be the state with $S_t \neq S'_t$. Suppose that $T_{k-1}(S_t, \text{STAY}) > T_{k-1}(S'_t, \text{STAY})$ and $\tilde{r}_{k,\text{STAY}}(S_t) < 1$. Then $\tilde{r}_{k,\text{STAY}}(S'_t) > \tilde{r}_{k,\text{STAY}}(S_t)$. Once this occurs the optimal policy in the optimistic MDP is to choose action GO. It follows easily that once t is sufficiently large the algorithm will alternate between choosing actions GO and STAY and subsequently suffer linear regret. Note that the uncertainty in the transitions does not play a big role here. In the optimistic MDP they will always be chosen to maximize the probability of transitioning to the state with the largest optimistic reward.

38.20 Here we abuse notation by letting $\hat{P}_{u,a}(s)$ be the empirical next-state transitions after u visits to state/action pair (s, a) . By a union bound and the result in Exercise 5.19,

$$\begin{aligned}\mathbb{P}(F) &\leq \mathbb{P}\left(\exists u \in \mathbb{N}, s, a \in \mathcal{S} \times \mathcal{A} : \|P - \hat{P}_{u,a}(s)\|_1 \geq \sqrt{\frac{2S \log(4SAu(u+1)/\delta)}{u}}\right) \\ &\leq \sum_{s,a \in \mathcal{S} \times \mathcal{A}} \sum_{u=1}^{\infty} \mathbb{P}\left(\|P - \hat{P}_{u,a}(s)\|_1 \geq \sqrt{\frac{2S \log(4SAu(u+1)/\delta)}{u}}\right) \\ &\leq \sum_{s,a \in \mathcal{S} \times \mathcal{A}} \sum_{u=1}^{\infty} \frac{\delta}{2SAu(u+1)} = \frac{\delta}{2}.\end{aligned}$$

The statement in Exercise 5.19 makes an independence assumption that is not exactly satisfied here. We are saved by the Markov property, which provides the conditional independence required.

38.21 We follow the hint. If $\sum_{k=1}^{m-1} a_k \leq 1$, then $A_0 = A_1 = \dots = A_{m-1} = 1$. Hence, $a_m \leq 1$ also holds and thus, using $A_m \geq 1$,

$$\sum_{k=1}^m \frac{a_k}{A_{k-1}} = \sum_{k=1}^{m-1} a_k + a_m \leq 1 + 1 \leq (\sqrt{2} + 1)A_m.$$

Let us now use induction on m . As long as m is so that $\sum_{k=1}^{m-1} a_k \leq 1$, the previous argument covers us. Thus, consider any $m > 1$ such that $\sum_{k=1}^{m-1} a_k > 1$ and assume that the statement holds for $m-1$ (note that $m=1$ implies $\sum_{k=1}^{m-1} a_k \leq 1$). Let

$c = \sqrt{2} + 1$. Then,

$$\begin{aligned}
\sum_{k=1}^m \frac{a_k}{A_{k-1}} &\leq c\sqrt{A_{m-1}} + \frac{a_m}{\sqrt{A_{m-1}}} && \text{(split sum, induction hypothesis)} \\
&= \sqrt{c^2 A_{m-1} + 2ca_m + \frac{a_m^2}{A_{m-1}}} \\
&= \sqrt{c^2 A_{m-1} + (2c+1)a_m} && (a_m \leq A_{m-1}) \\
&= c\sqrt{A_{m-1}} + a_m && \text{(choice of } c\text{)} \\
&= c\sqrt{A_m}. && (A_{m-1} \geq 1 \text{ and definition of } A_m)
\end{aligned}$$

38.22

- (a) The bound on (A) is DK where K is the number of phases. If there are no phases, $K = n$, and this term becomes linear in n .
- (b) Yes, this variant enjoys a similar regret bound. One obvious change to the proof is that now K is bounded by \sqrt{n} . The other change is that now

$$\sum_{k=1}^K \frac{T_{(k)}(s, a)}{\sqrt{1 \vee T_{\tau_k-1}(s, a)}}$$

needs to be bounded differently. For this, we can use Eq. (38.21).

38.23 We only outline the necessary changes to the proof of Theorem 38.4. The first step is to augment the failure event to include the event that there exists a phase k and state-action pair s, a such that

$$|\tilde{r}_{k,a}(s) - r_a(s)| \geq \sqrt{\frac{2L}{T_t(s, a)}}.$$

The likelihood of this event is at most $\delta/2$ by Hoeffding's bound combined with a union bound. Like in the proof of Theorem 38.4 we now restrict our attention to the regret on the event that the failure does not occur. The first change is that $r_{\pi_k}(s)$ in Eq. (38.18) must be replaced with $\tilde{r}_{k,\pi_k}(s)$. Then the reward terms no longer cancel in Eq. (38.20), which means that now

$$\begin{aligned}
\tilde{R}_k &= \sum_{t \in E_k} (-v_k(S_t) + \langle P_{k,A_t}(S_t), v_k \rangle + \tilde{r}_{k,A_t}(S_t) - r_{A_t}(S_t)) \\
&\leq \sum_{t \in E_k} (\langle P_{A_t}, v_k \rangle - v_k(S_t)) + \frac{D}{2} \sum_{t \in E_k} \|P_{k,A_t}(S_t) - P_{A_t}(S_t)\|_1 \\
&\quad + \sum_{t \in E_k} \sqrt{\frac{2L}{1 \vee T_{\tau_k-1}(S_t, A_t)}}.
\end{aligned}$$

The first two terms are the same as the proof of Theorem 38.4 and are bounded in the same way, which result in the same contribution to the regret. Only the last

term is new. Summing over all phases and applying the result from Exercise 38.21 and Cauchy-Schwarz,

$$\begin{aligned} \sum_{k=1}^K \sum_{t \in E_k} \sqrt{\frac{2L}{1 \vee T_{\tau_k-1}(S_t, A_t)}} &= \sqrt{2L} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{k=1}^K \frac{T_{(k)}(s, a)}{\sqrt{1 \vee T_{\tau_k-1}(s, a)}} \\ &\leq (\sqrt{2} + 1) \sqrt{2LSAn}. \end{aligned}$$

This term is small relative to the contribution due to the uncertainty in the transitions. Hence there exists a universal constant C such that with probability at least $1 - 3\delta/2$ the regret of this modified algorithm is at most

$$\hat{R}_n \leq CD(M)S \sqrt{nA \log \left(\frac{nSA}{\delta} \right)}.$$

38.24

- (a) An easy calculation shows that the depth d of the tree is bounded by $2 + \log_A S$, which by the conditions in the statement of the lower bound implies that $d + 1 \leq D/2$. The diameter is the maximum over all distinct pairs of states of the expected travel time between those two states. It is not hard to see that this is maximized by the pair s_g and s_b , so we restrict our attention to bounding the expected travel time between these two states under some policy. Let $\tau = \min\{t : S_t = s_b\}$ and let π be a policy that traverses the tree to a decision state with $\varepsilon(s, a) = 0$. We will show that for this policy

$$\mathbb{E}[\tau \mid S_1 = s_g] \leq D.$$

Let X_1, X_2, \dots be a sequence of random variables where $X_i \in \mathbb{N}^+$ is the number of rounds until the policy leaves state s_g on the i th series of visits to s_g . Then let M be the number of visits to state s_g before s_b is reached. All of these random variables are independent and geometrically distributed. An easy calculation shows that

$$\mathbb{E}[X_i] = 1/\delta \qquad \mathbb{E}[M] = 2.$$

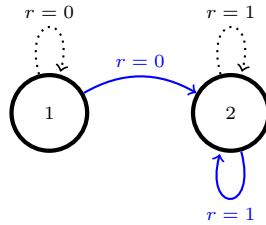
Then $\tau = \sum_{i=1}^M (X_i + d + 1)$, which has expectation $\mathbb{E}[\tau] = 2(1/\delta + d + 1) \leq 2/\delta + D/2 \leq D$.

- (b) The definition of stopping time τ ensures that $T_\sigma \leq n/D + 1 \leq 2n/D$ almost surely and hence $D\mathbb{E}[T_\sigma]/n \leq 2$ is immediate. For the second part note that

$$\mathbb{P} \left(\sum_{t=1}^{n/(2D)} D_t \geq \frac{n}{D} \right)$$

- (c)

38.30 Consider the deterministic MDP below with two states and two actions, $\mathcal{A} = \{\text{SOLID}, \text{DASHED}\}$.



Clearly the optimal policy is to choose $\pi(1) = \text{SOLID}$ and $\pi(2)$ arbitrarily which leads to a gain of 1. On the other hand, choosing $\rho = 1$ and $v = (2, 1)$ satisfies the linear program in Eq. (38.6) and the greedy policy with respect to this value function chooses $\pi(1) = \text{DASHED}$ and $\pi(2)$ arbitrary.

38.31 Let $\phi_k(x) = \text{TRUE}$ if $\langle a_k, x \rangle \geq b_k$ and $\phi_k(x) = a_k$ otherwise. Then the new separation oracle returns true if $\phi(x)$ and $\phi_k(x)$ are true for all k . Otherwise return the separating hyperplane provided by some ϕ or ϕ_k that did not return true.