

# **Probability and Statistics for Data Science**

## Solutions to Exercises

Carlos Fernandez-Granda



This document contains all solutions to the exercises in the book [Probability and Statistics for Data Science](#) by Carlos Fernandez-Granda, available for purchase [here](#). A free preprint, videos, code, slides and solutions to exercises are available at <https://www.ps4ds.net>

---

## Contents

1	Probability	5
2	Discrete Variables	15
3	Continuous Variables	24
4	Multiple Discrete Variables	36
5	Multiple Continuous Variables	54
6	Discrete and Continuous Variables	69
7	Averaging	80
8	Correlation	94
9	Estimation Of Population Parameters	106
10	Hypothesis Testing	115
11	Principal Component Analysis And Low-Rank Models	124
12	Regression And Classification	134

---

# Probability

## Exercises

1.1 (Conditional probability space) We check that  $\mathcal{C}_A$  satisfies the conditions:

- If  $B \in \mathcal{C}_A$ , then  $B^c \in \mathcal{C}_A$ . If the sample space is  $A$  then  $B^c = A - B$ . If  $B \in \mathcal{C}_A$ , there is some set  $S \in \mathcal{C}$  such that  $B = A \cap S$ . This implies that  $S^c \in \mathcal{C}$  because  $\mathcal{C}$  is a valid collection. As a result,  $S^c \cap A \in \mathcal{C}_A$ . We end the proof proving  $A - B = S^c \cap A$  by showing that they contain each other. (1) If  $\omega \in A - B$ , then  $\omega$  belongs to  $A$  and not to  $B$ . This means that it cannot belong to  $S$  because otherwise it would belong to  $B = A \cap S$ . This implies  $A - B \subseteq S^c \cap A$ . (2) If  $\omega \in S^c \cap A$ ,  $\omega$  belongs to  $A$  and not to  $S$ . It cannot belong to  $B$  because then it would be in  $S$ . This implies  $S^c \cap A \subseteq A - B$ .
- If  $B_1, B_2 \in \mathcal{C}_A$ , then  $B_1 \cup B_2 \in \mathcal{C}_A$ . If  $B_1, B_2 \in \mathcal{C}_A$ , then there exist  $S_1, S_2 \in \mathcal{C}$  such that  $B_1 = A \cap S_1$  and  $B_2 = A \cap S_2$ .  $S_1 \cup S_2$  is in  $\mathcal{C}$ , so  $A \cap (S_1 \cup S_2)$  is in  $\mathcal{C}_A$ . This completes the proof because  $A \cap (S_1 \cup S_2) = (A \cap S_1) \cup (A \cap S_2) = B_1 \cup B_2$ .
- If  $B_1, B_2, \dots \in \mathcal{C}_A$  then  $\bigcup_{i=1}^{\infty} B_i \in \mathcal{C}_A$ . This holds by the same argument as the finite case.
- $\mathcal{C}_A$  contains the sample space.  $A = A \cap A$ , so  $A \in \mathcal{C}_A$ .

Note that by the definition for any  $B \in \mathcal{C}_A$

$$P_A(B) := \frac{P(B)}{P(A)}. \quad (1.1)$$

We check that  $\mathcal{P}_A$  satisfies the conditions of a probability measure:

- $P_A(B) \geq 0$  for any event  $B \in \mathcal{C}_A$ . This just follows from  $P(B) \geq 0$ , and  $P(A) > 0$ .
- If  $B_1, B_2, \dots, B_n \in \mathcal{C}_A$  are disjoint then  $P_A(\bigcup_{i=1}^n B_i) = \sum_{i=1}^n P_A(B_i)$ . Since  $B_1, B_2, \dots, B_n$  are also in  $\mathcal{C}$  we have

$$P_A\left(\bigcup_{i=1}^n B_i\right) := \frac{P\left(\bigcup_{i=1}^n B_i\right)}{P(A)} \quad (1.2)$$

$$= \frac{\sum_{i=1}^n P(B_i)}{P(A)} \quad (1.3)$$

$$= \sum_{i=1}^n P_A(B_i). \quad (1.4)$$

- For a countably infinite sequence of disjoint sets  $B_1, B_2, \dots \in \mathcal{C}_A$   
 $P(\lim_{n \rightarrow \infty} \bigcup_{i=1}^n B_i) = \lim_{n \rightarrow \infty} \sum_{i=1}^n P(B_i)$ . This holds by the same argument as the finite case.
- The probability of the sample space equals 1. By the definition of  $P_A$

$$P_A(A) := \frac{P(A)}{P(A)} = 1. \quad (1.5)$$

1.2 (Empirical probability measure) We check that  $\mathcal{P}$  satisfies the conditions of a probability measure:

- $P(B) \geq 0$  for any event  $B \in \mathcal{C}$ . The numerator and denominator are both nonnegative by definition.
- If  $S_1, S_2, \dots, S_n \in \mathcal{C}$  are disjoint then  $P(\cup_{i=1}^n S_i) = \sum_{i=1}^n P(S_i)$ .

$$P(\cup_{i=1}^n S_i) = \frac{\text{number of data points with value in } \cup_{i=1}^n S_i}{N} \quad (1.6)$$

$$= \frac{\text{number in } S_1 + \text{number in } S_2 + \dots + \text{number in } S_n}{N} \\ = \sum_{i=1}^n P(S_i). \quad (1.7)$$

- For a countably infinite sequence of disjoint sets  $S_1, S_2, \dots \in \mathcal{C}$   
 $P(\lim_{n \rightarrow \infty} \cup_{i=1}^n S_i) = \lim_{n \rightarrow \infty} \sum_{i=1}^n P(S_i)$ . This follows by the same argument as in the finite case.
- The probability of the sample space equals 1.

$$P(\Omega) := \frac{\text{number of data points with value in } \Omega}{N} = 1. \quad (1.8)$$

1.3 (Independence and complements) Yes. If  $A$  and  $B$  are independent then

$$P(B) = P(A^c \cap B) + P(A \cap B) = P(A^c \cap B) + P(A)P(B), \quad (1.9)$$

which implies

$$P(A^c \cap B) = P(B)(1 - P(A)) = P(A^c)P(B). \quad (1.10)$$

1.4 (Conditional independence and complements) Not necessarily. For example, imagine that we flip two independent fair coins.  $A$  is the event that heads occurs on the first coin flip,  $C$  is the event that heads occurs on the second coin flip, and  $B = C \cup A^c$  is the event that heads occurs in the second coin flip, but not in the first. Then  $P(C) = 0.5$ ,  $P(A|C) = 0.5$ ,  $P(A|C^c) = 0.5$ ,  $P(B|A, C) = P(B|A^c, C) = 1$ ,  $P(B|A, C^c) = 0$ ,  $P(B|A^c, C^c) = 1$ . We have

$$P(A, B|C) = P(A|C)P(B|A, C) = 0.5, \quad (1.11)$$

$$P(B|C) = P(A|C)P(B|A, C) + P(A^c|C)P(B|A^c, C) = 1 \quad (1.12)$$

so  $P(A, B|C) = P(A|C)P(B|C) = 0.5$ , which means that  $A$  and  $B$  are conditionally independent given  $C$ . However,

$$P(A, B|C^c) = P(A|C^c)P(B|A, C^c) = 0, \quad (1.13)$$

$$P(B|C^c) = P(A|C^c)P(B|A, C^c) + P(A^c|C^c)P(B|A^c, C^c) = 0.5 \quad (1.14)$$

so  $P(A, B|C^c) \neq P(A|C^c)P(B|C^c) = 0.25$ , which means that  $A$  and  $B$  are not conditionally independent given  $C^c$ .

1.5 (Partition and independence) Any two events  $A$  and  $B$  in a partition are disjoint, which means that  $P(A \cap B) = 0$ . If they are independent then  $P(A)P(B) = P(A \cap B) = 0$  so that either  $P(A) = 0$  or  $P(B) = 0$ .

- 1.6 (Conditional probability and complements) Yes.  $P(A|B) = 1$  implies that  $P(A \cap B) = P(B)$ , so that

$$P(B^c|A^c) = \frac{P(A^c \cap B^c)}{P(A^c)} \quad (1.15)$$

$$= \frac{P((A \cup B)^c)}{P(A^c)} \quad \text{by DeMorgan's law} \quad (1.16)$$

$$= \frac{1 - P(A \cup B)}{1 - P(A)} \quad (1.17)$$

$$= \frac{1 - P(A) - P(B) + P(A \cap B)}{1 - P(A)} \quad (1.18)$$

$$= \frac{1 - P(A)}{1 - P(A)} \quad \text{because } P(A \cap B) = P(B) \quad (1.19)$$

$$= 1. \quad (1.20)$$

- 1.7 (The Linda problem) Let  $D$  be the event that the provided description applies to a person,  $B$  the event that someone is a bank teller and  $F$  the event that someone is active in the feminist movement. Since  $D \cap B \cap F$  and  $D \cap B \cap F^c$  form a partition of  $D \cap B$ , by the definition of conditional probability and the law of total probability,

$$\text{Conditional probability that Linda is a bank teller given the description} \quad (1.21)$$

$$= P(B|D) \quad (1.22)$$

$$= \frac{P(B, D)}{P(D)} \quad (1.23)$$

$$= \frac{P(B, D, F) + P(B, D, F^c)}{P(D)} \quad (1.24)$$

$$\geq \frac{P(B, D, F)}{P(D)} \quad (1.25)$$

$$= P(B, F|D) \quad (1.26)$$

$$= \text{Conditional probability that Linda is a bank teller and active in the feminist movement given the description.} \quad (1.27)$$

- 1.8 (Quiz)

$$P(\text{2nd question right} | \text{1st question right}) \quad (1.28)$$

$$\approx \frac{\text{Number of times 1st and 2nd question are both right}}{\text{Number of times 1st question is right}} \quad (1.29)$$

$$= \frac{4}{8} = \frac{1}{2}. \quad (1.30)$$

$$P(\text{2nd question right} | \text{1st question wrong}) \quad (1.31)$$

$$\approx \frac{\text{Number of times 1st question is wrong and 2nd question is right}}{\text{Number of times 1st question is wrong}} \quad (1.32)$$

$$= \frac{1}{2}. \quad (1.33)$$

The conditional probabilities are the same, suggesting the events could be independent (we can't be sure because of limited data).

- 1.9 (Baby name)

a) We define the events  $G$ ,  $B$  and  $M$  to represent that the baby is a girl, that Bob thinks

the baby is a girl, and that Margaret thinks that the baby is a girl. By the law of total probability and the chain rule

$$P(B \cap M) = P(G \cap B \cap M) + P(G^c \cap B \cap M) \quad (1.34)$$

$$\begin{aligned} &= P(G)P(M|G)P(B|M) + P(G^c)P(M|G^c)P(B|M) \\ &= 0.5 \cdot 0.8 \cdot 0.9 + 0.5 \cdot 0.2 \cdot 0.9 \end{aligned} \quad (1.35)$$

$$= 0.45. \quad (1.36)$$

b) By Bayes' rule, the law of total probability and the chain rule

$$P(G|B) \quad (1.37)$$

$$= \frac{P(G \cap B)}{P(B)} \quad (1.38)$$

$$\begin{aligned} &= \frac{P(G \cap M^c \cap B) + P(G \cap M \cap B)}{P(G^c \cap M^c \cap B) + P(G^c \cap M \cap B) + P(G \cap M^c \cap B) + P(G \cap M \cap B)} \\ &= \frac{0.5 \cdot 0.2 \cdot 0.1 + 0.5 \cdot 0.8 \cdot 0.9}{0.5 \cdot 0.8 \cdot 0.1 + 0.5 \cdot 0.2 \cdot 0.9 + 0.5 \cdot 0.2 \cdot 0.1 + 0.5 \cdot 0.8 \cdot 0.9} \end{aligned} \quad (1.39)$$

$$= 0.74. \quad (1.40)$$

1.10 (Cake)

a) By the definition of conditional probability,

$$P(\text{Scott helps} | \text{On time}) = \frac{P(\text{Scott helps, On time})}{P(\text{On time})}. \quad (1.41)$$

Under the independence assumptions, by the law of total probability and the chain rule

$$P(\text{Scott helps, On time}) \quad (1.42)$$

$$= P(\text{Scott helps, Antonis helps, On time}) + P(\text{Scott helps, Antonis doesn't, On time})$$

$$= P(\text{Scott helps})P(\text{Antonis helps})P(\text{On time} | \text{Scott helps, Antonis helps})$$

$$+ P(\text{Scott helps})P(\text{Antonis doesn't})P(\text{On time} | \text{Scott helps, Antonis doesn't})$$

$$= 0.4 \cdot 0.8 \cdot 1 + 0.4 \cdot 0.2 \cdot 0.5 \quad (1.43)$$

$$= 0.36. \quad (1.44)$$

Similarly,

$$P(\text{Scott doesn't, On time}) \quad (1.45)$$

$$= P(\text{Scott doesn't, Antonis helps, On time}) + P(\text{Scott doesn't, Antonis doesn't, On time})$$

$$= P(\text{Scott doesn't})P(\text{Antonis helps})P(\text{On time} | \text{Scott doesn't, Antonis helps})$$

$$+ P(\text{Scott doesn't})P(\text{Antonis doesn't})P(\text{On time} | \text{Scott doesn't, Antonis doesn't})$$

$$= 0.6 \cdot 0.8 \cdot 0.5 + 0.6 \cdot 0.2 \cdot 0 \quad (1.46)$$

$$= 0.24. \quad (1.47)$$

Since  $P(\text{On time}) = P(\text{Scott helps, On time}) + P(\text{Scott doesn't, On time}) = 0.6$ ,

$$P(\text{Scott helps} | \text{On time}) = \frac{0.36}{0.6} = 0.6. \quad (1.48)$$

b) By the definition of conditional probability,

$$P(\text{Scott helps} | \text{On time, Antonis helps}) = \frac{P(\text{Scott helps, Antonis helps, On time})}{P(\text{On time, Antonis helps})}.$$



Under the independence assumptions, by the chain rule

$$P(\text{Scott helps, Antonis helps, On time}) \quad (1.49)$$

$$= P(\text{Scott helps}) P(\text{Antonis helps}) P(\text{On time} \mid \text{Scott helps, Antonis helps}) \quad (1.50)$$

$$= 0.4 \cdot 0.8 \cdot 1 \quad (1.51)$$

$$= 0.32. \quad (1.52)$$

By the law of total probability

$$P(\text{Antonis helps, On time}) \quad (1.53)$$

$$= P(\text{Antonis helps, Scott helps, On time}) + P(\text{Antonis helps, Scott doesn't, On time})$$

$$= P(\text{Antonis helps}) P(\text{Scott helps}) P(\text{On time} \mid \text{Antonis helps, Scott helps})$$

$$+ P(\text{Antonis helps}) P(\text{Scott doesn't}) P(\text{On time} \mid \text{Antonis helps, Scott doesn't})$$

$$= 0.8 \cdot 0.4 \cdot 1 + 0.8 \cdot 0.6 \cdot 0.5 \quad (1.54)$$

$$= 0.56. \quad (1.55)$$

Consequently,

$$P(\text{Scott helps} \mid \text{On time, Antonis helps}) = \frac{0.32}{0.56} \quad (1.56)$$

$$= 0.57. \quad (1.57)$$

The event *Scott helps* is not conditionally independent from the event *Antonis helps* given the event *Milena finishes on time*, because

$$P(\text{Scott helps} \mid \text{On time, Antonis helps}) \neq P(\text{Scott helps} \mid \text{On time}). \quad (1.58)$$

- 1.11 (Baby sleep) We define the events  $W$ ,  $G$  and  $B$  to represent that the baby wakes up, is a good sleeper and eats bad food respectively.

a) By the chain rule and the law of total probability

$$P(W) = P(B^c, W) + P(B, W) \quad (1.59)$$

$$= P(B^c, G^c, W) + P(B^c, G, W) + P(B) P(W \mid B) \quad (1.60)$$

$$= P(B^c) P(G^c \mid B^c) P(W \mid B^c, G^c) \quad (1.61)$$

$$+ P(B^c) P(G \mid B^c) P(W \mid B^c, G) + P(B) P(W \mid B). \quad (1.62)$$

Here we need independence assumptions to make progress. It seems reasonable that the food being good or not (events  $B$  and  $B^c$ ) be independent from the baby being a good sleeper or not (events  $G$  and  $G^c$ ). Under that assumption, and noticing that  $P(W \mid B) = 1$ ,

$$P(W) = P(B^c) (P(G^c) P(W \mid B^c, G^c) + P(G) P(W \mid B^c, G)) + P(B) \quad (1.63)$$

$$= 0.9 \cdot (0.6 \cdot 0.1 + 0.4 \cdot 0.8) + 0.1 \quad (1.64)$$

$$= 0.442. \quad (1.64)$$

b) By Bayes' rule

$$P(B \mid W) = \frac{P(B) P(W \mid B)}{P(W)} \quad (1.65)$$

$$= \frac{0.1 \cdot 1}{0.442} \quad (1.66)$$

$$= 0.226. \quad (1.67)$$

- c) By the definition of conditional probability, the law of total probability and the chain rule

$$\begin{aligned} P(B|W, G) &= \frac{P(B, G, W)}{P(G, W)} \\ &= \frac{P(B)P(G)P(W|B, G)}{P(B)P(G)P(W|B, G) + P(B^c)P(G)P(W|B^c, G)}. \end{aligned} \quad (1.68)$$

If the food is bad, the baby wakes up no matter what, so  $P(W|B, G) = 1$ . Consequently

$$P(B|W, G) = \frac{P(B)}{P(B) + P(B^c)P(W|B^c, G)} \quad (1.69)$$

$$= \frac{0.1}{0.1 + 0.9 \cdot 0.1} \quad (1.70)$$

$$= 0.526 \quad (1.71)$$

- d) The events  $B$  and  $G$  are not conditionally independent given  $W$  because if they were  $P(B|W, G)$  would equal  $P(B|W)$ . This makes sense because if we know that the baby has woken up, whether the food is bad or not provides information about whether the baby is a good sleeper (and vice versa). In particular, conditioned on  $W$ , if the baby is good, then the food is more likely to be bad.

#### 1.12 (COVID-19 tests)

- a) Yes, it is reasonable to assume that the test only depends on whether that particular employee is ill, and not the others, and the events *Employee  $i$  is ill*, for  $1 \leq i \leq 10$ , are all independent.
- b) We define the events  $I_1, \dots, I_{10}$  to represent each employee being ill, and  $T_1, \dots, T_{10}$  to represent that the corresponding test is positive. The event that at least one test is positive is  $\cup_{i=1}^{10} T_i$ . By De Morgan's laws,

$$P(\cup_{i=1}^{10} T_i) = 1 - P((\cup_{i=1}^{10} T_i)^c) \quad (1.72)$$

$$= 1 - P(\cap_{i=1}^{10} T_i^c). \quad (1.73)$$

By the independence assumption

$$P(\cap_{i=1}^{10} T_i^c) = \prod_{i=1}^{10} P(T_i^c) \quad (1.74)$$

$$= \prod_{i=1}^{10} P(I_i)P(T_i^c | I_i) + P(I_i^c)P(T_i^c | I_i^c) \quad (1.75)$$

$$= (0.01 \cdot 0.02 + 0.99 \cdot 0.95)^{10} \quad (1.76)$$

$$= 0.543. \quad (1.77)$$

We conclude that  $P(\cup_{i=1}^{10} T_i) = 0.457$ .

- c) By the definition of conditional probability,

$$P(\cap_{i=1}^{10} I_i^c | \cup_{j=1}^{10} T_j) = \frac{P(\cap_{i=1}^{10} I_i^c, \cup_{j=1}^{10} T_j)}{P(\cup_{j=1}^{10} T_j)}. \quad (1.78)$$

The denominator was computed above. By the chain rule, the independence assumptions and DeMorgan's laws, the numerator equals

$$P(\cap_{i=1}^{10} I_i^c, \cup_{j=1}^{10} T_j) = P(\cap_{i=1}^{10} I_i^c)P(\cup_{j=1}^{10} T_j | \cap_{k=1}^{10} I_k^c) \quad (1.79)$$

$$= \prod_{i=1}^{10} P(I_i^c) \left(1 - P(\cap_{j=1}^{10} T_j^c | \cap_{k=1}^{10} I_k^c)\right) \quad (1.80)$$

We assume that  $T_j^c$  is conditionally independent of  $T_l^c$ ,  $j \neq l$ , conditioned on  $\cap_k I_k^c$ : if we

know that the other people do not have COVID-19, it is plausible that the outcome of the other tests should not provide information about the  $j$ th test. Under this assumption

$$P(\cap_{j=1}^{10} T_j^c \mid \cap_{k=1}^{10} I_k^c) = \prod_{j=1}^{10} P(T_j^c \mid \cap_{k=1}^{10} I_k^c). \quad (1.81)$$

It is reasonable to assume that  $T_i$  is conditionally independent of  $\cap_{j \neq i} I_j^c$  given  $I_i$  because one would expect that  $T_i$  only depends on  $I_i$ . Under this assumption

$$\prod_{j=1}^{10} P(T_j^c \mid \cap_{k=1}^{10} I_k^c) = \prod_{j=1}^{10} P(T_j^c \mid I_j^c). \quad (1.82)$$

Putting everything together

$$P(\cap_{i=1}^{10} I_i^c \mid \cup_{j=1}^{10} T_j) = \frac{\prod_{i=1}^{10} P(I_i^c)(1 - \prod_{j=1}^{10} P(T_j^c \mid I_j^c))}{P(\cup_{j=1}^{10} T_j)} \quad (1.83)$$

$$= \frac{0.99^{10}(1 - 0.95^{10})}{0.457} \quad (1.84)$$

$$= 0.793. \quad (1.85)$$

### 1.13 (Boxing championship)

a) From the assumptions,

$$P(\text{Manny is champion}) = P(\text{Manny beats Saul}) \cdot P(\text{Manny beats Floyd}) \quad (1.86)$$

$$= 0.4 \cdot 0.25 = 0.1, \quad (1.87)$$

$$P(\text{Saul is champion}) = P(\text{Saul beats Manny}) \cdot P(\text{Saul beats Floyd}) \quad (1.88)$$

$$= 0.6 \cdot 0.1 = 0.06, \quad (1.89)$$

$$P(\text{Floyd loses}) = P(\text{Manny is champion} \cup \text{Saul is champion}) \quad (1.90)$$

$$= P(\text{Manny is champion}) + P(\text{Saul is champion}) \quad (1.91)$$

$$= 0.16. \quad (1.92)$$

The conditional probability of interest is  $P(\text{Manny is champion} \mid \text{Floyd loses})$ . The event *Manny is champion* is included in the event *Floyd loses*, so their intersection is *Manny is champion*. Consequently,

$$P(\text{Manny is champion} \mid \text{Floyd loses}) = \frac{P(\text{Manny is champion})}{P(\text{Floyd loses})} \quad (1.93)$$

$$= \frac{0.1}{0.16} = \frac{5}{8}. \quad (1.94)$$

b) There are 4 simulations in which Floyd loses, and Manny wins in 3 of them, so the estimated probability is  $3/4$ . It is not surprising that the answer is different because we are approximating the probability using only 4 simulations.

### 1.14 (Videogame)

a) By the independence assumption,

$$P(\text{Wins game}) \quad (1.95)$$

$$= P(\text{Beats Honda}) P(\text{Beats Zangief}) P(\text{Beats Blanka}) \quad (1.96)$$

$$= (1 - P(\text{Loses to Honda})) (1 - P(\text{Loses to Zangief})) (1 - P(\text{Loses to Blanka}))$$

$$= (1 - (1 - 0.8)^2)(1 - (1 - 0.5)^2)(1 - (1 - 0.4)^2) \quad (1.97)$$

$$= 0.4608. \quad (1.98)$$

- b) Game 1 (win): Beats Honda in 2 fights, beats Zangief, beats Blanka in 2 fights.  
 Game 2 (loss): Beats Honda in 2 fights, loses to Zangief.  
 Game 3 (loss): Beats Honda, beats Zangief in 2 fights, loses to Blanka.  
 Game 4 (loss): Beats Honda, beats Zangief in 2 fights, loses to Blanka.  
 Game 5 (win): Beats Honda, beats Zangief in 2 fights, beats Blanka.  
 Game 6 (win): Beats Honda, beats Zangief, beats Blanka in 2 fights.  
 Game 7 (win): Beats Honda in two fights, beats Zangief in 2 fights, beats Blanka.  
 Game 8 (loss): Beats Honda, loses to Zangief.  
 Game 9 (loss): Beats Honda, beats Zangief, loses to Blanka.  
 Game 10 (loss): Loses to Honda.  
 Game 11 (win): Beats Honda, beats Zangief in 2 fights, beats Blanka in 2 fights.  
 Game 12 (loss): Beats Honda, beats Zangief in 2 fights, loses to Blanka.

The estimated probability of winning the game is  $5/12 = 0.42$ , which is quite close to the true value. We can improve the accuracy by performing more simulations.

- 1.15 (Rare event) The probability of not observing the event in  $n$  independent simulations is  $(1 - P(A))^n$ , consequently we require

$$1 - (1 - P(A))^n \geq 0.99, \quad (1.99)$$

which is equivalent to

$$0.01 \geq (1 - P(A))^n. \quad (1.100)$$

Taking logarithms on both sides,

$$\log 0.01 \geq n \log (1 - P(A)), \quad (1.101)$$

which implies

$$n \geq \frac{\log 0.01}{\log (1 - P(A))}, \quad (1.102)$$

because  $\log (1 - P(A))$  is negative. For  $P(A) := 0.01$ ,  $n \geq 458.2$ , so we need at least 459 simulations.

- 1.16 (Streak of heads)

- a) We represent heads with 1 and tails with 0. To compute the probabilities, we consider the  $2^5 = 32$  possible sequences:

00000 00001 00010 00011 00100 00101 00110 00111 01000 01001 01010 01011 01100  
 01101 01110 01111 10000 10001 10010 10011 10100 10101 10110 10111 11000 11001  
 11010 11011 11100 11101 11110 11111

We have

$$P(\text{sequence equals } 00000) \quad (1.103)$$

$$= P(\text{1st roll equals 0, 2nd roll equals 0, } \dots, \text{5th roll equals 0}) \quad (1.104)$$

$$= P(\text{1st roll equals 0})P(\text{2nd roll equals 0}) \cdots P(\text{5th roll equals 0}) \quad (1.105)$$

$$= \frac{1}{32} \quad (1.106)$$

and by the same argument, all the other sequences also have probability  $1/32$ .

Since the probability of the union of disjoint events is the sum of the individual probabilities,

$$P(\text{at most 0 heads in a row}) = \frac{1}{32}, \quad (1.107)$$

$$P(\text{at most 1 heads in a row}) = \frac{12}{32}, \quad (1.108)$$

$$P(\text{at most 2 heads in a row}) = \frac{11}{32}, \quad (1.109)$$

$$P(\text{at most 3 heads in a row}) = \frac{5}{32}, \quad (1.110)$$

$$P(\text{at most 4 heads in a row}) = \frac{2}{32}, \quad (1.111)$$

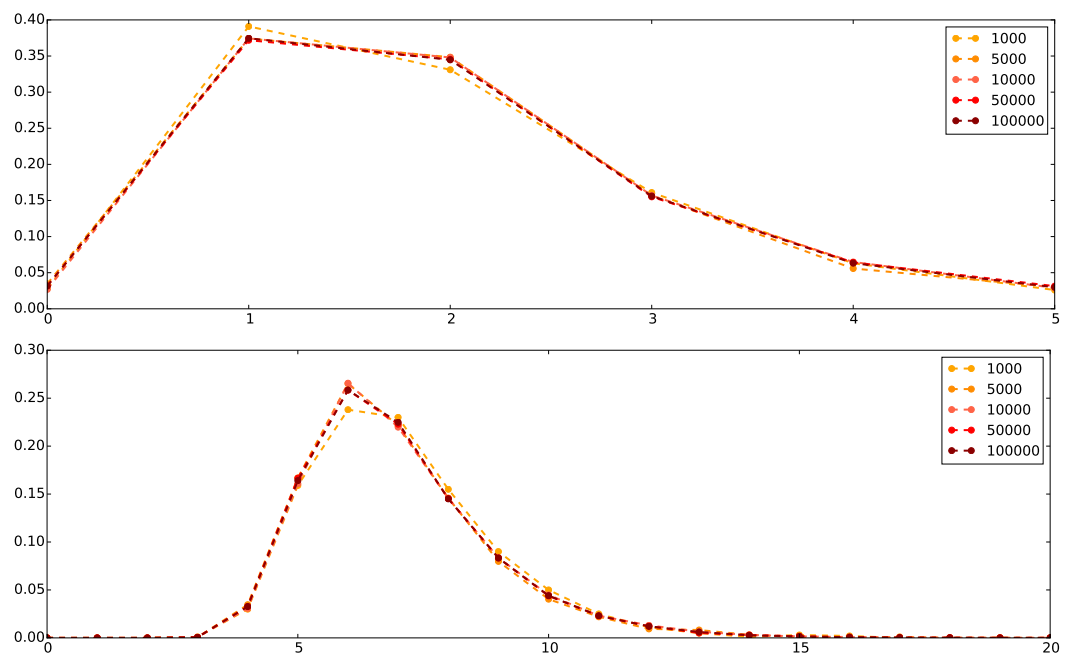
$$P(\text{at most 5 heads in a row}) = \frac{1}{32}. \quad (1.112)$$

b) A possible implementation in Python is:

```
def p_longest_streak(n, tries):
    p_longest = np.zeros(n+1)
    for i in range(tries):
        current_streak = 0
        longest_streak = 0
        for j in range(n):
            if np.random.rand() > 0.5:
                current_streak = current_streak + 1
            else:
                if current_streak > longest_streak:
                    longest_streak = current_streak
                current_streak = 0
        if current_streak > longest_streak:
            longest_streak = current_streak
        p_longest[longest_streak] = p_longest[longest_streak] + 1./tries
    return p_longest
```

The images are shown in Figure 1.1.

c) The probability is 0.319. It is therefore not unlikely to find a streak of 8 or more heads in a sequence of 200 fair coin flips, so it is very possible that the random generator is fine.



**Figure 1.1** Probability of streaks of heads for sequences of length 5 (above) and 200 (below) estimated using different number of Monte Carlo runs (indicated in the legend).

## 2

---

# Discrete Variables

### Exercises

- 2.1 (Flipping until heads) By Definition 1.28, if the coin flips are independent, the probability of first obtaining  $a - 1$  tails and then heads equals

$$p_{\tilde{a}}(a) := P(\tilde{a} = a) \quad (2.1)$$

$$= P(\text{flip } 1 = t, \dots, \text{flip } a - 1 = t, \text{flip } a = h) \quad (2.2)$$

$$= P(\text{flip } 1 = t) \cdots P(\text{flip } a - 1 = t) P(\text{flip } a = h) \quad (2.3)$$

$$= (1 - \alpha)^{a-1} \alpha, \quad (2.4)$$

where  $h$  represents heads and  $t$  tails. Here  $a$  is positive integer.

- 2.2 (Geometric pmf) For any  $\alpha \in (0, 1)$  the pmf is positive, so we only need to prove that it sums to one. By the geometric series identity  $\sum_{k=1}^{\infty} r^k = r(1 - r)^{-1}$ , setting  $r := 1 - \alpha$ ,

$$\sum_{a=1}^{\infty} p_{\tilde{a}}(a) = \frac{\alpha}{1 - \alpha} \sum_{a=1}^{\infty} (1 - \alpha)^a \quad (2.5)$$

$$= \frac{\alpha}{1 - \alpha} \frac{1 - \alpha}{1 - (1 - \alpha)} = 1. \quad (2.6)$$

- 2.3 (Poisson pmf) For any  $\lambda > 0$  the pmf is positive, so we only need to prove that it sums to one. This follows directly from the Taylor series expansion of the exponential function:

$$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{\lambda}. \quad (2.7)$$

Indeed,

$$\sum_{a=0}^{\infty} p_{\tilde{a}}(a) = e^{-\lambda} \sum_{a=0}^{\infty} \frac{\lambda^a}{(a)!} \quad (2.8)$$

$$= e^{-\lambda} e^{\lambda} = 1. \quad (2.9)$$

- 2.4 (Memoryless property of the geometric distribution) By definition of conditional probability, we have

$$P(\tilde{a} = a \mid \tilde{a} > 5) = \frac{P(\tilde{a} = a, \tilde{a} > 5)}{P(\tilde{a} > 5)}. \quad (2.10)$$

The numerator is zero unless  $a > 5$ , in that case

$$P(\tilde{a} = a \mid \tilde{a} > 5) = \frac{(1 - \alpha)^{a-1} \alpha}{\sum_{b=6}^{\infty} (1 - \alpha)^{b-1} \alpha}. \quad (2.11)$$

The denominator contains the geometric sum

$$\sum_{b=6}^{\infty} (1-\alpha)^{b-1} = \frac{(1-\alpha)^6}{\alpha}. \quad (2.12)$$

We conclude

$$P(\tilde{a} = a \mid \tilde{a} > 5) = (1-\alpha)^{a-6} \alpha. \quad (2.13)$$

This is the same as the geometric pmf if we plug in  $a - 5$ . In terms of the coin example, if we have obtained 5 tails, now the probability that we have to flip  $b = a - 5$  more times is distributed like a geometric random variable. This makes sense because the flips are independent, so there is no difference between this and just considering another sequence of flips starting from the beginning.

## 2.5 (Old car)

- a) A binomial distribution with parameters  $n := 50$  and  $\theta := 1/10$ . They need to assume that the car breaks down independently each time.
- b) We denote the binomial by  $\tilde{x}$

$$P(\text{At most once}) = p_{\tilde{x}}(0) + p_{\tilde{x}}(1) \quad (2.14)$$

$$= \binom{50}{0} \left(1 - \frac{1}{10}\right)^{50} + \binom{50}{1} \left(1 - \frac{1}{10}\right)^{49} \frac{1}{10} \quad (2.15)$$

$$= 0.9^{50} + 50 \cdot 0.9^{49} \cdot 0.1 \quad (2.16)$$

$$= 3.38 \cdot 10^{-2}. \quad (2.17)$$

As explained in the section on the Poisson distribution, the pmf of a binomial with parameters  $n$  and  $\lambda/n$  converges to the pmf of a Poisson with parameter  $\lambda$  (in distribution), we set  $\lambda = n/\theta = 5$  for a Poisson random variable  $\tilde{y}$  and compute

$$P(\text{At most once}) \approx p_{\tilde{y}}(0) + p_{\tilde{y}}(1) \quad (2.18)$$

$$= e^{-\lambda} + \lambda e^{-\lambda} \quad (2.19)$$

$$= 4.04 \cdot 10^{-2}. \quad (2.20)$$

- c) The probability that the car breaks down  $k - 1$  times in  $n - 1$  drives is equal to the probability that a binomial random variable with parameters  $n - 1$  and 0.1 equals  $k - 1$ . Consequently, under the assumption that the car breaks down independently each time,

$$P(\text{breaks down } n\text{th time in } k\text{th drive}) \quad (2.21)$$

$$= P(\text{breaks down } k - 1 \text{ times in } n - 1 \text{ drives} \cap \{\text{breaks down in } k\text{th drive}\}) \quad (2.22)$$

$$= P(\text{breaks down } k - 1 \text{ times in } n - 1 \text{ drives}) P(\{\text{breaks down in } k\text{th drive}\}) \quad (2.23)$$

$$= \binom{n-1}{k-1} \left(\frac{9}{10}\right)^{n-k} \left(\frac{1}{10}\right)^k. \quad (2.24)$$

- 2.6 (Oil prospector) Let  $\tilde{x}$  denote the random variable representing the number of times the prospector has to drill to find oil. By the law of total probability and the conditional



independence assumptions:

$$p_{\tilde{x}}(x) \quad (2.25)$$

$$= P(\text{no oil in } x-1 \text{ attempts, oil in } x\text{th}) \quad (2.26)$$

$$= \sum_{r \in \{\text{poor, standard, rich}\}} P(\text{region is } r) P(\text{no oil in } x-1 \text{ attempts, oil in } x\text{th} \mid \text{region is } r)$$

$$= \frac{1}{10} \frac{1}{2^x} + \frac{8}{10} \frac{9^{x-1}}{10^x} + \frac{1}{10} \frac{99^{x-1}}{100^x} \quad (2.27)$$

$$= \frac{1}{10} \left( \frac{1}{2^x} + \frac{8 \cdot 9^{x-1}}{10^x} + \frac{99^{x-1}}{100^x} \right). \quad (2.28)$$

- 2.7 (Darts) Let  $\tilde{a}$  denote the random variable. Note that the last attempt must always be a success. We can therefore decompose the event *a attempts required* into the intersection of *k-1 successes over first a-1 attempts* and *ath attempt is a success*. Since the attempts are all independent,

$$p_{\tilde{a}}(a) = P(k-1 \text{ successes over first } a-1 \text{ attempts}) P(\text{ath attempt is a success}). \quad (2.29)$$

By exactly the same reasoning we use to derive the binomial distribution,

$$P(k-1 \text{ successes over first } a-1 \text{ attempts}) = \binom{a-1}{k-1} \theta^{k-1} (1-\theta)^{a-1-(k-1)}, \quad (2.30)$$

as long as  $a \geq k$ . We conclude that

$$p_{\tilde{a}}(a) = \binom{a-1}{k-1} \theta^k (1-\theta)^{a-k}, \quad \text{for } a \geq k, \quad (2.31)$$

and zero otherwise.

- 2.8 (Binomial random variable)

$$P(\tilde{a} = 1 \mid \tilde{a} \leq 1) = \frac{P(\tilde{a} = 1, \tilde{a} \leq 1)}{P(\tilde{a} \leq 1)} \quad (2.32)$$

$$= \frac{P(\tilde{a} = 1)}{P(\tilde{a} = 0) + P(\tilde{a} = 1)} \quad (2.33)$$

$$= \frac{n\theta(1-\theta)^{n-1}}{(1-\theta)^n + n\theta(1-\theta)^{n-1}} \quad (2.34)$$

$$= \frac{n\theta}{1 + (n-1)\theta}. \quad (2.35)$$

- 2.9 (Intersection) The parametric model is  $p_{\tilde{c}}(-1) = p_{\tilde{c}}(1) = \theta$  and  $p_{\tilde{c}}(0) = 1 - 2\theta$ . The log likelihood is

$$\log \mathcal{L}_X(\theta) = \sum_{i=1}^6 \log p_{\tilde{c}}(x_i) \quad (2.36)$$

$$= 4 \log \theta + 2 \log(1 - 2\theta). \quad (2.37)$$

The derivatives of the log likelihood are

$$\frac{d \log \mathcal{L}_X(\theta)}{d\theta} = \frac{4}{\theta} - \frac{4}{1-2\theta}, \quad (2.38)$$

$$\frac{d^2 \log \mathcal{L}_X(\theta)}{d\theta^2} = -\frac{4}{\theta^2} - \frac{8}{(1-2\theta)^2} < 0 \quad \text{for all } \theta \in [0, 1]. \quad (2.39)$$

The function is concave so we can set the first derivative to zero, obtaining  $\theta_{\text{ML}} = 1/3$ .

2.10 (Bad apples)

a) The probability of picking a bad apple is  $\theta$ , so

$$p_{\bar{t}}(1) = (\text{Test is positive}) \quad (2.40)$$

$$= P(\text{Test is positive, Bad apple}) + P(\text{Test is positive, Good apple}) \quad (2.41)$$

$$= P(\text{Bad apple}) P(\text{Test is positive} | \text{Bad apple}) \quad (2.42)$$

$$+ P(\text{Good apple}) P(\text{Test is positive} | \text{Good apple}) \quad (2.43)$$

$$= 0.9\theta + 0.4(1 - \theta) \quad (2.44)$$

$$= 0.5\theta + 0.4. \quad (2.45)$$

Since the pmf must add to one,

$$p_{\bar{t}}(0) = 1 - p_{\bar{t}}(1) \quad (2.46)$$

$$= 0.6 - 0.5\theta. \quad (2.47)$$

b) The log-likelihood equals

$$\log \mathcal{L}(\theta) = \sum_{i=1}^n \log p_{\bar{t}}(x_i) \quad (2.48)$$

$$= 5 \log(0.5\theta + 0.4) + 5 \log(0.6 - 0.5\theta). \quad (2.49)$$

The first derivative is

$$\frac{d \log \mathcal{L}(\theta)}{d\theta} = \frac{5 \cdot 0.5}{0.5\theta + 0.4} + \frac{5 \cdot (-0.5)}{0.6 - 0.5\theta} \quad (2.50)$$

$$= \frac{2.5}{0.5\theta + 0.4} - \frac{2.5}{0.6 - 0.5\theta} \quad (2.51)$$

and the second

$$\frac{d^2 \log \mathcal{L}(\theta)}{d\theta^2} = -\frac{2.5 \cdot 0.5}{(0.5\theta + 0.4)^2} - \frac{2.5 \cdot 0.5}{(0.6 - 0.5\theta)^2} < 0, \quad (2.52)$$

so the function is concave. Setting the first derivative to zero yields,

$$0.5\theta_{\text{ML}} + 0.4 = 0.6 - 0.5\theta_{\text{ML}}, \quad (2.53)$$

$$\theta_{\text{ML}} = 0.2. \quad (2.54)$$

2.11 (False positives)

a) A test is negative only if the person does not have the disease and there is no false positive so

$$p_{\bar{t}}(0) = (\text{Test is negative}) \quad (2.55)$$

$$= P(\text{No disease, No false positive}) \quad (2.56)$$

$$= P(\text{No disease}) P(\text{No false positive} | \text{No disease}) \quad (2.57)$$

$$= 0.9(1 - \theta). \quad (2.58)$$

Since the pmf must add to one,

$$p_{\bar{t}}(1) = 1 - p_{\bar{t}}(0) \quad (2.59)$$

$$= 0.1 + 0.9\theta. \quad (2.60)$$

b) Let  $n_{\text{neg}} := n - n_{\text{pos}}$ . The log likelihood is

$$\log \mathcal{L}(\theta) = \sum_{i=1}^n \log p_i(x_i) \quad (2.61)$$

$$= n_{\text{neg}} \log 0.9(1 - \theta) + n_{\text{pos}} \log (0.1 + 0.9\theta) \quad (2.62)$$

$$= n_{\text{neg}} \log 0.9 + n_{\text{neg}} \log (1 - \theta) + n_{\text{pos}} \log (0.1 + 0.9\theta). \quad (2.63)$$

The first derivative is

$$\frac{d \log \mathcal{L}(\theta)}{d\theta} = -\frac{n_{\text{neg}}}{1 - \theta} + \frac{0.9n_{\text{pos}}}{0.1 + 0.9\theta} \quad (2.64)$$

and the second

$$\frac{d^2 \log \mathcal{L}(\theta)}{d\theta^2} = -\frac{n_{\text{neg}}}{(1 - \theta)^2} - \frac{0.9^2 n_{\text{pos}}}{(0.1 + 0.9\theta)^2} < 0, \quad (2.65)$$

so the function is concave. Setting the first derivative to zero yields,

$$0.1n_{\text{neg}} + 0.9\theta_{\text{max}}n_{\text{neg}} = 0.9n_{\text{pos}} - 0.9n_{\text{pos}}\theta_{\text{max}}, \quad (2.66)$$

which implies

$$\theta_{\text{max}} = \frac{0.9n_{\text{pos}} - 0.1n_{\text{neg}}}{0.9(n_{\text{pos}} + n_{\text{neg}})} \quad (2.67)$$

$$= \frac{n_{\text{pos}}}{n} - \frac{0.1n_{\text{neg}}}{0.9n} \quad (2.68)$$

$$= \phi - \frac{1}{9}(1 - \phi) \quad (2.69)$$

$$= \frac{10\phi - 1}{9}. \quad (2.70)$$

Notice that for  $\phi < 0.1$  the maximum is negative ( $\theta_{\text{max}} < 0$ ), and consequently out of the interval  $[0, 1]$ . Since we know that the log-likelihood is monotone decreasing for values of  $\theta$  larger than the maximum (as it is concave), in such cases the maximum within  $[0, 1]$  is at zero. We conclude that

$$\theta_{\text{ML}} = \begin{cases} 0 & \text{if } \phi < 0.1, \\ \frac{10\phi - 1}{9} & \text{if } \phi \geq 0.1. \end{cases} \quad (2.71)$$

## 2.12 (Chess games)

a) Under the independence assumption, the likelihood and log-likelihood equal

$$\mathcal{L}_X(\theta) = \theta^4 \alpha^2 (1 - \theta - \alpha)^4, \quad (2.72)$$

$$\log \mathcal{L}_X(\theta) = 4 \log \theta + 2 \log \alpha + 4 \log (1 - \theta - \alpha). \quad (2.73)$$

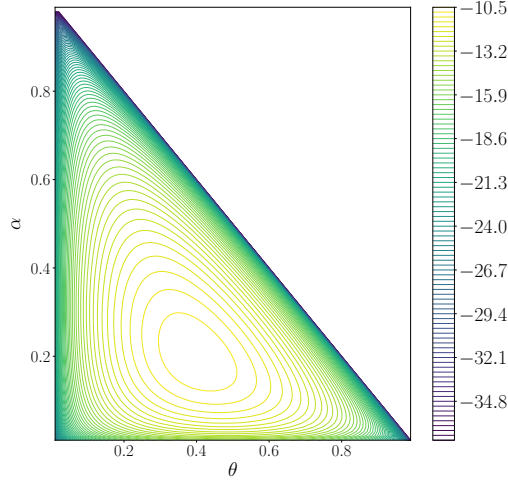
The plot is shown in Figure 2.1.

b) From the plot we can see that the function has a single maximum. To find it, we set the partial derivatives to zero:

$$\frac{d \log \mathcal{L}_X(\theta)}{d\theta} = \frac{4}{\theta} - \frac{4}{1 - \theta - \alpha}, \quad (2.74)$$

$$\frac{d \log \mathcal{L}_X(\alpha)}{d\alpha} = \frac{2}{\alpha} - \frac{4}{1 - \theta - \alpha}. \quad (2.75)$$

$$(2.76)$$



**Figure 2.1** Log-likelihood of the parametric model for the *Chess games* problem.

Setting the first expression equal to zero yields  $\alpha_{\text{ML}} = 1 - 2\theta_{\text{ML}}$ . Plugging into the second and solving the equation, we conclude  $\theta_{\text{ML}} = 0.4$  and  $\alpha_{\text{ML}} = 0.2$ .

- c) The empirical pmf would assign  $4/10 = 0.4$  to the probability of Garry winning,  $2/10 = 0.2$  to the probability of Anish winning, and  $4/10 = 0.4$  to the probability of a draw. This is exactly equivalent to the parametric model.

2.13 (Maximum-likelihood estimator for the geometric distribution) The log-likelihood equals

$$\log \mathcal{L}_{\{x_1, \dots, x_n\}}(\alpha) = \sum_{i=1}^n \log p_{\alpha}(x_i) \quad (2.77)$$

$$= \sum_{i=1}^n \log \left( (1 - \alpha)^{x_i - 1} \alpha \right) \quad (2.78)$$

$$= \sum_{i=1}^n (\log \alpha + (x_i - 1) \log (1 - \alpha)) \quad (2.79)$$

$$= n \log \alpha + \left( \sum_{i=1}^n x_i - n \right) \log (1 - \alpha). \quad (2.80)$$

The derivative and second derivative of the log-likelihood are

$$\frac{d \log \mathcal{L}_{\{x_1, \dots, x_n\}}(\alpha)}{d\alpha} = \frac{n}{\alpha} - \frac{\sum_{i=1}^n x_i - n}{1 - \alpha}, \quad (2.81)$$

$$\frac{d^2 \log \mathcal{L}_{\{x_1, \dots, x_n\}}(\alpha)}{d\alpha^2} = -\frac{n}{\alpha^2} - \frac{\sum_{i=1}^n x_i - n}{(1 - \alpha)^2} \quad (2.82)$$

$$= -n \cdot \frac{1 - 2\alpha + \alpha^2 \frac{1}{n} \sum_{i=1}^n x_i}{\alpha^2 (1 - \alpha)^2} \quad (2.83)$$

$$< -n \cdot \frac{1 - 2\alpha + \alpha^2}{\alpha^2 (1 - \alpha)^2} = -\frac{n}{\alpha^2} < 0. \quad (2.84)$$

Where we have used the assumption that  $x_i \geq 1$  for all  $i$ , so that  $\frac{1}{n} \sum_{i=1}^n x_i \geq 1$ . The function is concave, as the second derivative is negative. The maximum-likelihood estimator of the parameter (see Definition 2.25) is consequently at the point where the first derivative equals zero, namely

$$\alpha_{\text{ML}} := \arg \max_{\alpha} \log \mathcal{L}_{\{x_1, \dots, x_n\}}(\alpha) \quad (2.85)$$

$$= \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^{-1}. \quad (2.86)$$

#### 2.14 (Multinoulli distribution)

- a) The parameters  $\theta_1, \dots, \theta_m$  must be between 0 and 1, as they represent entries of a pmf. They also need to sum to one, so that the pmf is valid.
- b) We are interested in deriving the maximum likelihood estimate of  $\theta_k$ , for an arbitrary  $1 \leq k \leq m$ . The likelihood and log-likelihood equal

$$\mathcal{L}_X(\theta) = \prod_{j=1}^m \theta_j^{n_j} \quad (2.87)$$

$$= \left( 1 - \sum_{l \neq k} \theta_l \right)^{n_k} \prod_{j \neq k} \theta_j^{n_j}, \quad (2.88)$$

$$\log \mathcal{L}_X(\theta) = n_k \log \left( 1 - \sum_{l \neq k} \theta_l \right) + \sum_{j \neq k} n_j \log \theta_j, \quad (2.89)$$

where  $n_j := |\{i : x_i = v_j\}|$  is the number of data points equal to  $v_j$  for  $1 \leq j \leq m$ . We have used the fact that  $\theta_k = 1 - \sum_{l \neq k} \theta_l$  because the parameters need to sum up to one. Assuming the log-likelihood is concave, we can obtain the maximum-likelihood estimator setting the partial derivatives (and hence the gradient) to zero:

$$\frac{d \log \mathcal{L}_X(\theta)}{d \theta_j} = -\frac{n_k}{1 - \sum_{l \neq k} \theta_l} + \frac{n_j}{\theta_j} = 0, \quad 1 \leq j \leq m. \quad (2.90)$$

This yields

$$n_k \theta_j = n_j \theta_k, \quad 1 \leq j \leq m, \quad (2.91)$$

and consequently

$$n_k = \sum_{j=1}^m n_k \theta_j = \sum_{j=1}^m n_j \theta_k \quad (2.92)$$

$$= n \theta_k, \quad (2.93)$$

because  $\sum_{j=1}^m n_j = n$ . We conclude that the maximum-likelihood estimator of  $\theta_k$  is  $\frac{n_k}{n}$ .

- c) The parametric model is exactly equivalent to the nonparametric empirical-pmf estimator.

#### 2.15 (Adaptive test)

a) Let  $\tilde{c}$  be the number of correct answers.

$$p_{\tilde{c}}(0) = P(1\text{st wrong}) P(2\text{nd wrong} \mid 1\text{st wrong}) \quad (2.94)$$

$$= (1 - \theta^2) \quad (2.95)$$

$$= 1 - 2\theta + \theta^2, \quad (2.96)$$

$$p_{\tilde{c}}(1) = P(1\text{st correct}) P(2\text{nd wrong} \mid 1\text{st correct}) \quad (2.97)$$

$$+ P(1\text{st wrong}) P(2\text{nd correct} \mid 1\text{st wrong}) \quad (2.98)$$

$$= \theta \left(1 - \frac{\theta}{2}\right) + (1 - \theta)\theta \quad (2.99)$$

$$= 2\theta - \frac{3\theta^2}{2}, \quad (2.100)$$

$$p_{\tilde{c}}(2) = P(1\text{st correct}) P(2\text{nd correct} \mid 1\text{st correct}) \quad (2.101)$$

$$= \frac{\theta^2}{2}. \quad (2.102)$$

b)

$$\log \mathcal{L}(\theta) = \sum_{i=1}^n \log p_{\tilde{c}}(x_i) \quad (2.103)$$

$$= \log p_{\tilde{c}}(2) + 2 \log p_{\tilde{c}}(1) \quad (2.104)$$

$$= \log \left(\frac{\theta^2}{2}\right) + 2 \log \left(2\theta - \frac{3\theta^2}{2}\right) \quad (2.105)$$

$$= 2 \log \theta - \log 2 + 2 \log \left(2\theta - \frac{3\theta^2}{2}\right). \quad (2.106)$$

$$\frac{d \log \mathcal{L}(\theta)}{d\theta} = \frac{2}{\theta} + \frac{2(2 - 3\theta)}{2\theta - \frac{3\theta^2}{2}}. \quad (2.107)$$

Setting the derivative to zero yields

$$\frac{1}{\theta_{\text{ML}}} = \frac{-2 + 3\theta_{\text{ML}}}{2\theta_{\text{ML}} - \frac{3\theta_{\text{ML}}^2}{2}}, \quad (2.108)$$

$$2\theta_{\text{ML}} - \frac{3\theta_{\text{ML}}^2}{2} = -2\theta_{\text{ML}} + 3\theta_{\text{ML}}^2, \quad (2.109)$$

$$4\theta_{\text{ML}} = \frac{9\theta_{\text{ML}}^2}{2}. \quad (2.110)$$

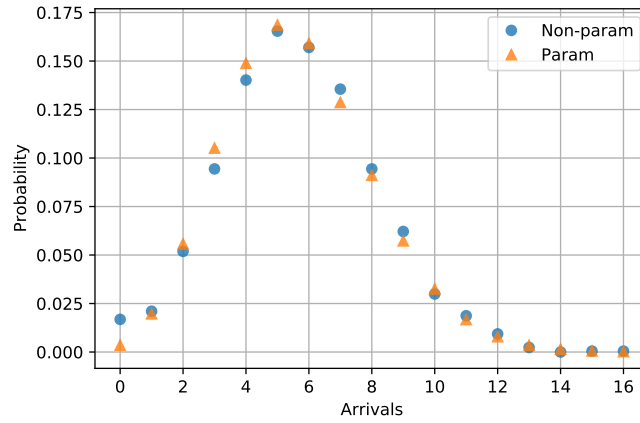
This expression is zero if we set  $\theta_{\text{ML}}$  equal to zero, but then the derivative would be infinite! The only value for which the derivative is zero is

$$\theta_{\text{ML}} = \frac{8}{9}. \quad (2.111)$$

By the assumption that the log-likelihood is concave, this is the maximum-likelihood estimate.

## 2.16 (Air traffic)

a) To design the parametric model, we assume that:



**Figure 2.2** Parametric and nonparametric pmfs for the *Air traffic* problem.

- a) For any period of time of length  $t$ , if  $t$  is small enough, the probability of a flight arriving in that period is equal to  $\lambda t$  and the probability of more than one flight is negligible.  $\lambda$  is a fixed parameter representing the total flights that we expect to arrive every 10 minutes.
- b) Each flight arrives independently from other flights..

These are the assumptions we made to derive the Poisson distribution, so we model the number of arrivals using a Poisson parametric model. The code to fit the nonparametric and parametric models is the following: The parametric and nonparametric pmfs are shown in Figure 2.2.

- b) The test RMSD is 0.016 for the nonparametric model and 0.021 for the parametric model. In this case, the nonparametric model is more accurate.

### 3

## Continuous Variables

### Exercises

- 3.1 (Probability of individual points) Let  $S := \{s \in \mathbb{R} : P(\tilde{a} = s) \neq 0\}$  be the set of values that are assigned nonzero probability. Since the events  $\tilde{a} = s$  are all disjoint (the function cannot map the same outcome to two values), we must have

$$\sum_{s \in S} P(\tilde{a} = s) = 1. \quad (3.1)$$

A technical detail is that if  $S$  is uncountable the sum may not be well defined. We define it as the supremum of the sums of any countable subset of  $S$ . Now consider the following partition of  $S$

$$A_n := \left\{ a \in S : \frac{1}{n-1} \geq P(\tilde{a} = a) > \frac{1}{n} \right\}, \quad n = 2, 3, \dots \quad (3.2)$$

None of these sets can contain more than  $n$  elements (otherwise the sum of the probabilities would be larger than one, contradicting (3.1)). The union of these subsets is equal to  $S$ , and there are a countable number of them, each containing a finite number of elements. We therefore conclude that  $S$  is countable.

- 3.2 (Continuous cdf) For any fixed  $a$ , leveraging the properties of probability measures, we have

$$F_{\tilde{a}}(a) = P(\tilde{a} \leq a) \quad (3.3)$$

$$= P\left(\lim_{n \rightarrow \infty} \left\{ \tilde{a} \leq a - \frac{1}{n} \right\} \cup \bigcup_{i=2}^n \left\{ a - \frac{1}{n-1} < \tilde{a} \leq a - \frac{1}{n} \right\} \cup \{\tilde{a} = a\}\right) \quad (3.4)$$

$$= P(\tilde{a} \leq a - 1) + P\left(\lim_{n \rightarrow \infty} \bigcup_{i=2}^n \left\{ a - \frac{1}{n-1} < \tilde{a} \leq a - \frac{1}{n} \right\}\right) + P(\tilde{a} = a) \quad (3.5)$$

$$= P(\tilde{a} \leq a - 1) + \lim_{n \rightarrow \infty} \sum_{i=2}^n P\left(a - \frac{1}{n-1} < \tilde{a} \leq a - \frac{1}{n}\right) + P(\tilde{a} = a) \quad (3.6)$$

$$= \lim_{n \rightarrow \infty} P\left(\tilde{a} \leq a - \frac{1}{n}\right) + P(\tilde{a} = a) \quad (3.7)$$

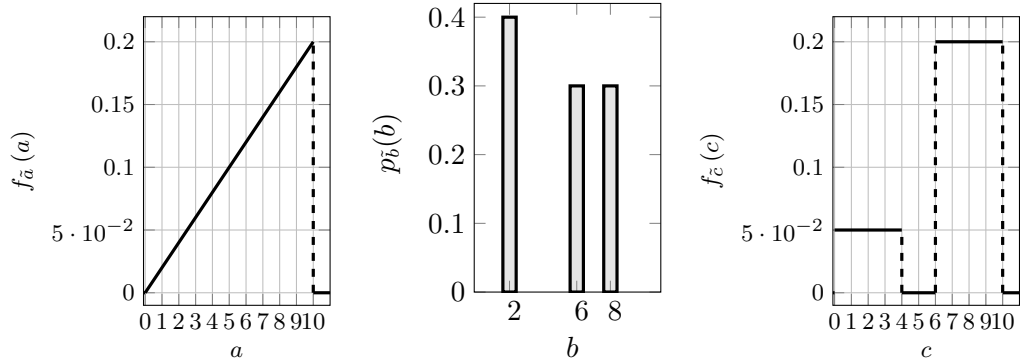
$$= \lim_{n \rightarrow \infty} F_{\tilde{a}}\left(a - \frac{1}{n}\right) + P(\tilde{a} = a). \quad (3.8)$$

Consequently,  $P(\tilde{a} = a)$  is zero, if and only if

$$\lim_{n \rightarrow \infty} F_{\tilde{a}}\left(a - \frac{1}{n}\right) = F_{\tilde{a}}(a), \quad (3.9)$$

i.e. if and only if  $F_{\tilde{a}}$  is continuous at  $a$ .





**Figure 3.1** Pdfs of  $\tilde{a}$  and  $\tilde{c}$  and pmf of  $\tilde{b}$  in Exercise 3.4.

3.3 (Median of affine transformation) If the median of  $\tilde{a}$  is  $m$ , then

$$P(\tilde{a} \leq m) = \frac{1}{2}, \quad (3.10)$$

which implies

$$P(\tilde{b} \leq \alpha m + \beta) = P(\alpha \tilde{a} + \beta \leq \alpha m + \beta) \quad (3.11)$$

$$= P(\tilde{a} \leq m), \quad (3.12)$$

so the median of  $\tilde{b}$  is  $\alpha m + \beta$ .

3.4 (Three cdfs)

a) The median is the value for which the cdf equals 0.5. The median of  $\tilde{a}$  is 7, the median of  $\tilde{b}$  is 6 and the median of  $\tilde{c}$  is 7.5.

b) Figure 3.1 shows the pdfs of  $\tilde{a}$  and  $\tilde{c}$ , and the pmf of  $\tilde{b}$ .

3.5 (Cumulative distribution function)

a) The median is 0.2, because the probability that  $\tilde{x} \leq 0.2$  is 0.5.

b) The slope of the cdf is clearly larger at 0.2, so the density at 0.2 is higher than at 0.8.

c)

$$P(\tilde{x} \leq 0.4 | \tilde{x} < 0.2) = \frac{P(0.2 \leq \tilde{x} \leq 0.4 | \tilde{x} \geq 0.2)}{P(\tilde{x} \geq 0.2)} \quad (3.13)$$

$$= \frac{F_{\tilde{x}}(0.4) - F_{\tilde{x}}(0.2)}{1 - F_{\tilde{x}}(0.2)} \quad (3.14)$$

$$= \frac{0.8 - 0.5}{0.5} \quad (3.15)$$

$$= 0.6. \quad (3.16)$$

d) We apply inverse-transform sampling. Since  $F_{\tilde{x}}(0.08) = 0.2$ ,  $F_{\tilde{x}}(0.15) = 0.4$ , and  $F_{\tilde{x}}(0.25) = 0.55$ , the simulated samples are 0.08, 0.15, and 0.25.

3.6 (Step pdf)

a) The pdf must integrate to 1, the value of the integral is  $0.8 \cdot b/2 + 0.2 \cdot b/2 = 0.5b$  so  $b$  must equal 2.

- b) For  $a < 0$ , the cdf is equal to zero, because the probability that  $\tilde{a} \leq a$  is zero. For  $0 \leq a < 1$ ,

$$F_{\tilde{a}}(a) = P(\tilde{a} \leq a) \quad (3.17)$$

$$= \int_{t=0}^a 0.8 \, dt \quad (3.18)$$

$$= 0.8t. \quad (3.19)$$

For  $1 \leq a < 2$ ,

$$F_{\tilde{a}}(a) = P(\tilde{a} \leq a) \quad (3.20)$$

$$= \int_{t=0}^1 0.8 \, dt + \int_{t=1}^a 0.2 \, dt \quad (3.21)$$

$$= 0.8 + 0.2(a - 1) \quad (3.22)$$

$$= 0.6 + 0.2a. \quad (3.23)$$

For  $a \geq 2$ , the cdf is equal to one, because the probability that  $\tilde{a} \leq a$  is one.

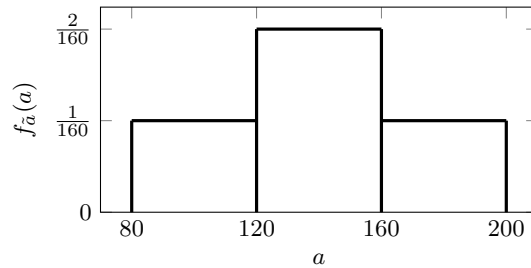
The median of  $\tilde{a}$ , which we denote by  $m$  satisfies  $F_{\tilde{a}}(m) = 0.5$ . We notice that the probability that  $\tilde{a} < 1$  is more than 0.5 (it's 0.8), so the median must be below, which implies

$$F_{\tilde{a}}(m) = 0.8m = 0.5, \quad (3.24)$$

so  $m = 5/8$ .

### 3.7 (Pigs)

- a) The histogram is shown below.



- b) The maximum likelihood estimate for the  $\mu$  parameter of the Gaussian is the average of the data, which equals 145 kg. This is also the median because the Gaussian pdf is centered at  $\mu$  and is symmetric, so the probability that the corresponding random variable is smaller than  $\mu$  is 0.5.
- 3.8 (Scaled exponential distribution) The cdf of  $\tilde{a}$  equals  $F_{\tilde{a}}(a) = 1 - \exp(-\lambda a)$  for  $a \geq 0$  and

0 for  $a < 0$ . The cdf of  $\tilde{b}$  therefore equals

$$F_{\tilde{b}}(b) := P(\tilde{b} \leq b) \quad (3.25)$$

$$= P(\alpha \tilde{a} \leq b) \quad (3.26)$$

$$= P\left(\tilde{a} \leq \frac{b}{\alpha}\right) \quad (3.27)$$

$$= F_{\tilde{a}}\left(\frac{b}{\alpha}\right) \quad (3.28)$$

$$= 1 - \exp\left(-\frac{\lambda}{\alpha}b\right), \quad (3.29)$$

so  $\tilde{b}$  is an exponential random variable with parameter  $\lambda/\alpha$ .

- 3.9 (Gaussian pdf) Let  $f_{\tilde{a}}$  denote the pdf of a Gaussian random variable with mean  $\mu$  and standard deviation  $\sigma$ . By the change of variables  $t = (a - \mu)/\sqrt{2}\sigma$ ,

$$\int_{-\infty}^{\infty} f_{\tilde{a}}(a) da = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} \exp(-t^2) dt. \quad (3.30)$$

Let us define

$$I = \int_{-\infty}^{\infty} \exp(-t^2) dt. \quad (3.31)$$

We perform a change of variable to polar coordinates to the square of  $I$ :

$$I^2 = \int_{-\infty}^{\infty} \exp(-x^2) dx \int_{-\infty}^{\infty} \exp(-y^2) dy \quad (3.32)$$

$$= \int_{x=-\infty}^{\infty} \int_{y=-\infty}^{\infty} \exp(-(x^2 + y^2)) dx dy \quad (3.33)$$

$$= \int_{\theta=0}^{2\pi} \int_{r=-\infty}^{\infty} r \exp(-r^2) d\theta dr \quad (3.34)$$

$$= \pi \exp(-r^2) \Big|_0^{\infty} = \pi. \quad (3.35)$$

Consequently,

$$\int_{-\infty}^{\infty} f_{\tilde{a}}(a) da = \frac{I}{\sqrt{\pi}} = 1. \quad (3.36)$$

- 3.10 (Uniform distribution) Recall that  $F_{\tilde{u}}(u) = u$  for  $0 \leq u \leq 1$ . For  $0 \leq w \leq 1$ , the cdf of  $\tilde{w}$  equals

$$F_{\tilde{w}}(w) = P(1 - \tilde{u} \leq w) \quad (3.37)$$

$$= P(\tilde{u} \geq 1 - w) \quad (3.38)$$

$$= 1 - P(\tilde{u} \leq 1 - w) \quad (3.39)$$

$$= 1 - F_{\tilde{u}}(1 - w) = 1 - (1 - w) = w. \quad (3.40)$$

For  $w > 1$ ,  $1 - \tilde{u}$  is always smaller than  $w$ , so  $F_{\tilde{w}}(w) = 1$ . For  $w < 0$ ,  $1 - \tilde{u}$  is always greater than  $w$ , so  $F_{\tilde{w}}(w) = 0$ . We conclude that  $\tilde{w}$  indeed follows a uniform distribution in  $[0, 1]$ .

- 3.11 (Nuclear power plant)

a) The pdf should integrate to one. We have

$$\int_{-\infty}^{\infty} f_{\tilde{t}}(t) dt = \int_{-1}^0 \alpha dt + \int_0^{\infty} \alpha \exp(t) dt \quad (3.41)$$

$$= \alpha(0 - (-1)) + \alpha(\exp(0) - \exp(-\infty)) \quad (3.42)$$

$$= 2\alpha, \quad (3.43)$$

so  $\alpha = 1/2$ .

b) We compute the cdf of  $\tilde{t}$  conditioned on the event  $\{\tilde{t} < 0\}$ :

$$F_{\tilde{t}|\tilde{t}<0}(t) := P(\tilde{t} \leq t | \tilde{t} < 0) \quad (3.44)$$

$$= \frac{P(\tilde{t} \leq \min\{0, t\})}{P(\tilde{t} < 0)} \quad (3.45)$$

$$= \begin{cases} 0 & \text{if } t < -1, \\ \frac{\int_{a=-1}^t f_{\tilde{t}}(a) da}{\int_{a=-1}^0 f_{\tilde{t}}(a) da} & \text{if } -1 \leq t < 0, \\ 1 & \text{if } t \geq 0. \end{cases} \quad (3.46)$$

For  $-1 \leq t < 0$ ,

$$\int_{a=-1}^t f_{\tilde{t}}(a) da = \frac{t+1}{2}, \quad (3.47)$$

so

$$F_{\tilde{t}|\tilde{t}<0}(t) = \begin{cases} 0 & \text{if } t < -1, \\ 1+t & \text{if } -1 \leq t < 0, \\ 1 & \text{if } t \geq 0. \end{cases} \quad (3.48)$$

Differentiating, we obtain

$$f_{\tilde{t}|\tilde{t}<0}(t) = \begin{cases} 0 & \text{if } t < -1, \\ 1 & \text{if } -1 \leq t < 0, \\ 0 & \text{if } t \geq 0. \end{cases} \quad (3.49)$$

The conditional distribution is uniform between -1 and 0.

### 3.12 (Evaluating survival analysis)

a) Since  $S_{\text{true}}(t) := P(\tilde{t} > t) = 1 - F_{\tilde{t}}(t)$ ,

$$\tilde{w} := S_{\text{true}}(\tilde{t}) \quad (3.50)$$

$$= 1 - F_{\tilde{t}}(\tilde{t}). \quad (3.51)$$

By the probability integral transform  $F_{\tilde{t}}(\tilde{t})$  is a uniform random variable in  $[0, 1]$ . Consequently, the distribution of  $\tilde{w}$  is also uniform in  $[0, 1]$  by Exercise 3.10.

b) If  $S_{\text{est}}$  is a good approximation to  $S_{\text{true}}$ , then  $S_{\text{est}}(\tilde{t})$  should be approximately uniformly distributed. We can evaluate whether this is the case by computing the transformed samples

$$w_i := S_{\text{est}}(t_i), \quad 1 \leq i \leq n, \quad (3.52)$$

and checking whether their empirical cdf resembles the linear cdf of a uniform random variable in  $[0, 1]$ .

3.13 (Feeding data into its own empirical cdf) Since  $x_1, x_2, \dots, x_n$  are sorted,

$$y_i = F_X(x_i) \quad (3.53)$$

$$= \frac{1}{n} \sum_{j=1}^n 1(x_j \leq x_i) = \frac{i}{n}. \quad (3.54)$$

The answer does not depend on the dataset.

3.14 (Uniform temperature)

- a) We denote the temperature by  $\tilde{t}$ . Its pdf is uniform between 60 and 80, so it equals  $1/20$ . For  $t \leq 60$  the cdf equals zero, because the probability that  $\tilde{t} \leq 60$  is zero. For  $t \geq 80$  the cdf equals one, because the probability that  $\tilde{t} \leq 80$  is one. For  $60 < t < 80$ ,

$$F_{\tilde{t}}(t) = P(\tilde{t} \leq t) \quad (3.55)$$

$$= \int_{t=60}^t \frac{1}{20} dt \quad (3.56)$$

$$= \frac{t - 60}{20}. \quad (3.57)$$

The first quartile  $q_1$  satisfies

$$F_{\tilde{t}}(q_1) = \frac{q_1 - 60}{20} = \frac{1}{4}, \quad (3.58)$$

so  $q_1 = 65$ .

- b) We apply the inverse transform method. The inverse of the cdf satisfies

$$F_{\tilde{t}}(F_{\tilde{t}}^{-1}(u)) = \frac{F_{\tilde{t}}^{-1}(u) - 60}{20} \quad (3.59)$$

$$= u, \quad (3.60)$$

so  $F_{\tilde{t}}^{-1}(u) = 60 + 20u$ . We apply this inverse cdf to generate the samples. For  $u=0.1$ , the sample is  $t = 62$ . For  $u = 0.8$  the sample is  $t = 76$ .

3.15 (Rounded-up measurements)

- a) Let  $\tilde{d}$  be the time the particle takes to decay, so  $\tilde{r} = \lceil \tilde{d} \rceil$ . Its pmf equals

$$P(\tilde{r} = r) = P(r - 1 \leq \tilde{d} < r) = \int_{r-1}^r \lambda e^{-\lambda x} dx = e^{-\lambda(r-1)} - e^{-\lambda r} \quad (3.61)$$

$$= (e^{-\lambda})^{r-1} (1 - e^{-\lambda}) \quad \text{for } r = 1, 2, 3, \dots \quad (3.62)$$

The reading is a geometric random variable with parameter  $1 - e^{-\lambda}$ .

- b) Let  $\tilde{x}$  be the error. For  $x \leq 0$ ,  $F_{\tilde{x}}(x) = 0$  since the error cannot be smaller than 0. It

also cannot be larger than 1, so for  $x \geq 1$   $F_{\tilde{x}}(x) = 1$ . For  $0 \leq x \leq 1$ ,

$$F_{\tilde{x}}(x) = P(\tilde{x} \leq x) \quad (3.63)$$

$$= P(\lceil \tilde{d} \rceil - \tilde{d} \leq x) \quad (3.64)$$

$$= P\left(\bigcup_{i=1}^{\infty} \{i - x \leq \tilde{d} \leq i\}\right) \quad \text{union of disjoint events} \quad (3.65)$$

$$= \sum_{i=1}^{\infty} P(i - x \leq \tilde{d} \leq i) \quad (3.66)$$

$$= \sum_{i=1}^{\infty} \lambda \int_{i-x}^i e^{-\lambda x} dx \quad (3.67)$$

$$= \sum_{i=1}^{\infty} e^{-\lambda(i-x)} - e^{-\lambda i} \quad (3.68)$$

$$= (e^{\lambda x} - 1) \sum_{i=1}^{\infty} e^{-\lambda i} \quad (3.69)$$

$$= \frac{e^{-\lambda} (e^{\lambda x} - 1)}{1 - e^{-\lambda}} = \frac{e^{\lambda x} - 1}{e^{\lambda} - 1}. \quad (3.70)$$

Differentiating we obtain

$$f_{\tilde{x}}(x) = \begin{cases} \frac{\lambda e^{\lambda x}}{e^{\lambda} - 1} & 0 \leq x \leq 1 \\ 0 & \text{otherwise.} \end{cases} \quad (3.71)$$

### 3.16 (Half life)

a) We have

$$P(\tilde{t} > t_{1/2}) = \int_{t_{1/2}}^{\infty} \lambda \exp(-\lambda x) dx \quad (3.72)$$

$$= \exp(-\lambda t_{1/2}). \quad (3.73)$$

Setting equal to 1/2 yields  $t_{1/2} = \frac{\ln 2}{\lambda}$ . This is the time that it takes for the particle to decay with probability 1/2. Consequently, from a large group of particles with this distribution, about half of them will decay after that time, so this is a reasonable definition of half life.

b) We have

$$P(t_{1/2} < \tilde{t} < t) = \int_{t_{1/2}}^t \lambda \exp(-\lambda x) dx \quad (3.74)$$

$$= \exp(-\lambda t_{1/2}) - \exp(-\lambda t) \quad (3.75)$$

$$= \frac{1}{2} - \exp(-\lambda t). \quad (3.76)$$

Setting equal to 1/4 yields  $t_{1/4} = \frac{\ln 4}{\lambda} = \frac{2 \ln 2}{\lambda} = 2t_{1/2}$ . Intuitively, this means that after two half-lives, one fourth of the particles remain. This makes sense: the number of particles is halved after the first half-time, and then halved again after the second.

c)

$$P(\tilde{t} > kt_{1/2}) = \int_{kt_{1/2}}^{\infty} \lambda \exp(-\lambda x) dx \quad (3.77)$$

$$= \exp(-\lambda kt_{1/2}) \quad (3.78)$$

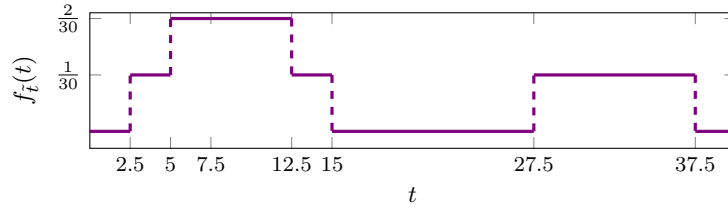
$$= \exp(-\ln 2^k) \quad (3.79)$$

$$= \frac{1}{2^k}. \quad (3.80)$$

This says that if we wait for  $k$  half times, the remaining particles are halved  $k$  times, which is again consistent with the intuitive definition of half time.

## 3.17 (Earthquake)

a) The KDE estimate looks like this:



The probability is equal to

$$P(\tilde{t} > 10) = \frac{2 \cdot 2.5}{30} + \frac{2.5}{30} + \frac{10}{30} \quad (3.81)$$

$$= \frac{17.5}{30} \quad (3.82)$$

$$= 0.583, \quad (3.83)$$

which is the area under the pdf to the right of 10.

b) From the notes, the ML estimate for the parameter of the exponential is

$$\lambda_{\text{ML}} = \frac{3}{7.5 + 10 + 32.5} \quad (3.84)$$

$$= \frac{3}{50}. \quad (3.85)$$

The probability equals

$$P(\tilde{t} > 10) = \int_{10}^{\infty} f_{\tilde{t}}(t) dt \quad (3.86)$$

$$= \int_{10}^{\infty} \frac{3 \exp(-3t/50)}{50} dt \quad (3.87)$$

$$= -\exp(-3t/50)]_{10}^{\infty} \quad (3.88)$$

$$= 0.549. \quad (3.89)$$

c) The parametric method requires less data but makes a stronger assumption about the distribution. The nonparametric method is more flexible but requires more data.

## 3.18 (Uniform distribution with a bump)

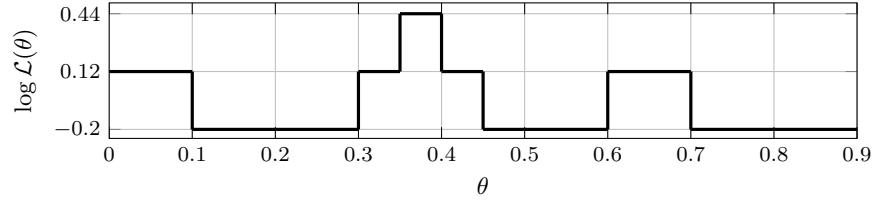
a)

$$F_{\tilde{a}}(a) = \begin{cases} 0 & \text{if } a \leq 0, \\ 0.9a & \text{if } 0 \leq a \leq \theta, \\ 0.9\theta + 1.9(a - \theta) = 1.9a - \theta & \text{if } \theta \leq a \leq \theta + 0.1, \\ 0.9\theta + 0.19 + 0.9(a - \theta - 0.1) = 0.9a + 0.1 & \text{if } \theta \leq a \leq \theta + 0.1, \\ 1 & \text{if } a \geq 1. \end{cases} \quad (3.90)$$

b)

$$\log \mathcal{L}(\theta) = \sum_{i=1}^4 \log f_{\tilde{a}}(x_i) = \begin{cases} 3 \log 0.9 + \log 1.9 = 0.12 & \text{if } 0 \leq \theta \leq 0.1, \\ 4 \log 0.9 = -0.2 & \text{if } 0.1 < \theta < 0.3, \\ 0.12 & \text{if } 0.3 \leq \theta < 0.35, \\ 2 \log 0.9 + 2 \log 1.9 = 0.44 & \text{if } 0.35 \leq \theta \leq 0.4, \\ 0.12 & \text{if } 0.4 < \theta \leq 0.45, \\ -0.2 & \text{if } 0.45 < \theta < 0.6, \\ 0.12 & \text{if } 0.6 \leq \theta \leq 0.7, \\ -0.2 & \text{if } 0.7 < \theta \leq 0.9. \end{cases} \quad (3.91)$$

The log-likelihood function looks like this:



### 3.19 (Planet)

a) Since

$$F_{\tilde{t}}(t) = \int_{-\infty}^t f_{\tilde{t}}(u) \, du, \quad (3.92)$$

for  $t < 0$ ,

$$F_{\tilde{t}}(t) = \int_{-\infty}^t \frac{\lambda \exp(\lambda u)}{2} \, du \quad (3.93)$$

$$= \frac{\exp(\lambda t)}{2}, \quad (3.94)$$

and for  $t \geq 0$ ,

$$F_{\tilde{t}}(t) = \int_{-\infty}^0 \frac{\lambda \exp(\lambda u)}{2} \, du + \int_0^t \frac{\lambda \exp(-\lambda u)}{2} \, du \quad (3.95)$$

$$= \frac{1}{2} + \frac{1 - \exp(-\lambda t)}{2} \quad (3.96)$$

$$= 1 - \frac{\exp(-\lambda t)}{2}. \quad (3.97)$$



b) The log-likelihood is

$$\log \mathcal{L}_X(\lambda) = \sum_{i=1}^n \log f_\lambda(x_i) \quad (3.98)$$

$$= \sum_{i=1}^n \log \frac{\lambda \exp(-\lambda |x_i|)}{2} \quad (3.99)$$

$$= n \log \lambda - n \log 2 - \lambda \sum_{i=1}^n |x_i|. \quad (3.100)$$

The derivative and second derivative of the log-likelihood function are

$$\frac{d \log \mathcal{L}_{x_1, \dots, x_n}(\lambda)}{d\lambda} = \frac{n}{\lambda} - \sum_{i=1}^n |x_i|, \quad (3.101)$$

$$\frac{d^2 \log \mathcal{L}_{x_1, \dots, x_n}(\lambda)}{d\lambda^2} = -\frac{n}{\lambda^2} < 0 \quad \text{for all } \lambda > 0. \quad (3.102)$$

The function is concave, as the second derivative is negative, so there cannot be different local maxima. The maximum is obtained by setting the first derivative to zero, which yields

$$\lambda_{\text{ML}} = \frac{1}{\frac{1}{n} \sum_{i=1}^n |x_i|} \quad (3.103)$$

$$= \frac{1}{39}. \quad (3.104)$$

c) The conditional cdf of  $\tilde{t}$  given  $\tilde{t} \geq 0$  evaluated at  $t > 0$  is

$$F_{\tilde{t}|\tilde{t}>0}(t) = P(\tilde{t} \leq t | \tilde{t} > 0) \quad (3.105)$$

$$= \frac{P(0 < \tilde{t} \leq t)}{P(\tilde{t} > 0)} \quad (3.106)$$

$$= \frac{F_{\tilde{t}}(t) - F_{\tilde{t}}(0)}{1 - F_{\tilde{t}}(0)} \quad (3.107)$$

$$= \frac{1 - \frac{\exp(-\lambda t)}{2} - \frac{1}{2}}{\frac{1}{2}} \quad (3.108)$$

$$= 1 - \exp(-\lambda t). \quad (3.109)$$

Differentiating with respect to  $t$  yields an exponential pdf  $f_{\tilde{t}|\tilde{t}>0}(t) = \lambda e^{-\lambda t}$  with parameter  $\lambda$ .

3.20 (Triangular pdf)

a) The possible values of  $w$  are  $w \geq \max(x_1, \dots, x_n) = 1.5$ .

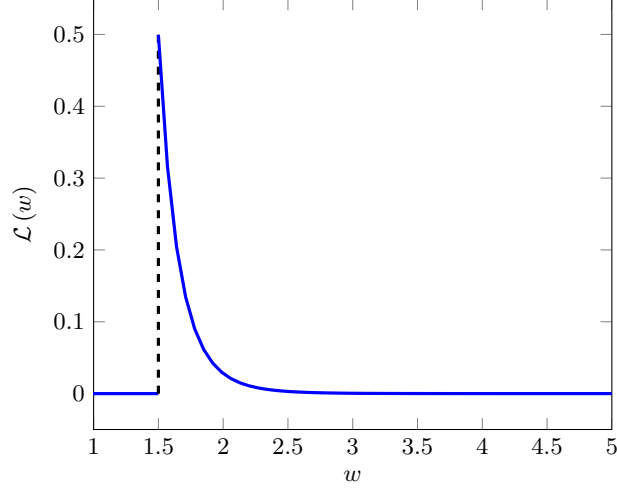
b) The likelihood equals

$$\mathcal{L}(w) = \prod_i^n f_w(x_i) \quad (3.110)$$

$$= \frac{2^5}{w^{10}} \prod_i^n x_i \quad (3.111)$$

$$= \frac{28.8}{w^{10}} \quad (3.112)$$

if  $w \geq \max(x_1, \dots, x_n)$  and 0 otherwise because in that case  $f_w(\max(x_1, \dots, x_n)) = 0$ .



- c) The maximum likelihood estimate is 1.5 because the likelihood is decreasing over  $[1.5, \infty)$ .
- d) The ML estimate systematically underestimates the true parameter. It is equal to the maximum value in the data, which has to be smaller than  $w_{\text{true}}$ .
- e) We apply inverse-transform sampling. The cdf equals

$$F_w(x) = \begin{cases} 0, & x \leq 0, \\ \frac{x^2}{w^2}, & \text{for } 0 \leq x \leq w, \\ 1, & \text{for } x \geq w. \end{cases} \quad (3.113)$$

The inverse of the cdf in the interval of interest is  $F_w^{-1}(y) = w\sqrt{y}$ . The sample is therefore  $2\sqrt{0.64} = 1.6$ .

### 3.21 (Rat)

- a) The pdf needs to be nonnegative, which requires  $0 \leq \alpha \leq 1$ . We also need the pdf to integrate to one, which it does if  $\alpha$  is in that range:

$$\int_{a=0}^1 f_{\tilde{a}}(a) da = \int_{a=0}^{0.5} 2\alpha da + \int_{a=0.5}^1 2(1-\alpha) da \quad (3.114)$$

$$= 2\alpha \cdot 0.5 + 2(1-\alpha) \cdot 0.5 \quad (3.115)$$

$$= 1. \quad (3.116)$$

- b) Expressing the pdf as a function of  $\alpha$  the likelihood of each data point is equal to  $2\alpha$  if the point is between 0 and 0.5, and to  $2(1-\alpha)$  if it is between 0.5 and 1. Let  $n$  be the number of data,  $n_{[0,0.5]}$  the number of data between 0 and 0.5, and  $n_{(0.5,1]}$  the number of points between 0.5 and 1. We have,

$$\log \mathcal{L}(X) = \sum_{i=1}^n \log f_{\alpha}(x_i) \quad (3.117)$$

$$= n_{[0,0.5]} \log 2\alpha + n_{(0.5,1]} \log(2(1-\alpha)). \quad (3.118)$$

The first and second derivatives of the log likelihood equal

$$(\log \mathcal{L}(X))' = \frac{n_{[0,0.5]}}{\alpha} - \frac{n_{(0.5,1]}}{1-\alpha}, \quad (3.119)$$

$$(\log \mathcal{L}(X))'' = -\frac{n_{[0,0.5]}}{\alpha^2} - \frac{n_{(0.5,1]}}{(1-\alpha)^2}. \quad (3.120)$$

The function is concave, so we can set the first derivative to zero to find the ML estimate, it equals

$$\alpha_{\text{ML}} = \frac{n_{[0,0.5]}}{n} \quad (3.121)$$

$$= \frac{2}{5}. \quad (3.122)$$

c) The probability of the rat being in the first half is

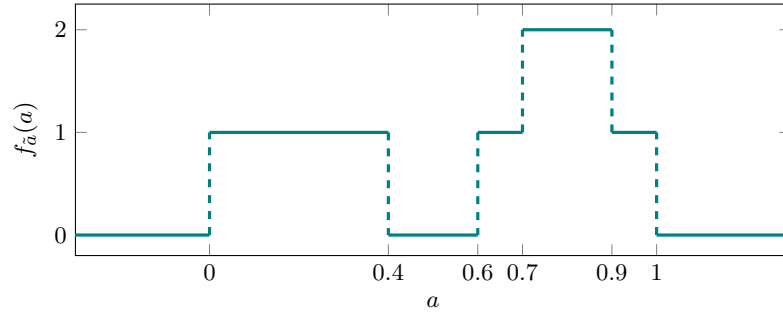
$$P(0 \leq \tilde{a} \leq 0.5) = \int_{a=0}^{0.5} 2\alpha \, da \quad (3.123)$$

$$= \alpha. \quad (3.124)$$

Estimating this probability using empirical probabilities yields exactly the same estimate for  $\alpha$  as the ML estimate,

$$\alpha_{\hat{\text{ML}}} = \frac{n_{[0,0.5]}}{n}. \quad (3.125)$$

d) The estimated pdf is:



What is problematic is that the density is zero between 0.4 and 0.6, which seems an artifact of the limited number of data. This can be alleviated by increasing the width of the rectangular kernel.

---

## Multiple Discrete Variables

### Exercises

#### 4.1 (Halloween)

- a) Lou grabs  $\tilde{l}$  chocolate bars, where  $\tilde{l}$  is a random variable with values between 0 and 2. Since all three values have the same probability,  $p_{\tilde{l}}(l) = \frac{1}{3}$ , for  $l \in \{0, 1, 2\}$  and zero otherwise. Ellie grabs  $\tilde{e}$  chocolate bars, where  $\tilde{e}$  is a random variable with values between 0 and  $2 - \tilde{l}$ . Since they have the same probability, the conditional pmf equals

$$p_{\tilde{e}|\tilde{l}}(a|l) = \begin{cases} \frac{1}{3-l} & \text{if } 0 \leq a \leq 2-l, \\ 0 & \text{otherwise.} \end{cases} \quad (4.1)$$

By the chain rule, the joint pmf is:

$$p_{\tilde{e},\tilde{l}}(a,l) = \frac{1}{3(3-l)} \text{ for } l = 0, \dots, 2 \text{ and } a = 0, \dots, 2-l. \quad (4.2)$$

In the  $(a, l)$  plane, the points on which the joint pmf is defined form a triangle.

- b) To find the marginal pmf, we sum over  $l$ . For each value of  $a$  the sum has a different number of terms:

$$p_{\tilde{e}}(0) = \frac{1}{3} \left( \frac{1}{3} + \frac{1}{2} + \frac{1}{1} \right) = \frac{11}{18}, \quad (4.3)$$

$$p_{\tilde{e}}(1) = \frac{1}{3} \left( \frac{1}{3} + \frac{1}{2} \right) = \frac{5}{18}, \quad (4.4)$$

$$p_{\tilde{e}}(2) = \frac{1}{3} \frac{1}{3} = \frac{1}{9}. \quad (4.5)$$

- c) When  $\tilde{e} = 1$ ,  $\tilde{l}$  must be either 0 or 1. By the definition of conditional pmf,

$$p_{\tilde{l}|\tilde{e}}(0|1) = \frac{p_{\tilde{e},\tilde{l}}(1,0)}{p_{\tilde{e}}(1)} \quad (4.6)$$

$$= \frac{\frac{1}{3} \frac{1}{3}}{\frac{5}{18}} = \frac{2}{5}, \quad (4.7)$$

$$p_{\tilde{l}|\tilde{e}}(1|1) = \frac{p_{\tilde{e},\tilde{l}}(1,1)}{p_{\tilde{e}}(1)} \quad (4.8)$$

$$= \frac{\frac{1}{3} \frac{1}{2}}{\frac{5}{18}} = \frac{3}{5}. \quad (4.9)$$

4.2 (Empirical conditional pmf) The number of data points in  $X$  equal to  $a$  is  $n_a$ , so

$$p_X(a) := \frac{1}{n} \sum_{i=1}^n 1(x_i = a) \quad (4.10)$$

$$= \frac{n_a}{n}. \quad (4.11)$$

We can write the number of data point pairs such that the first entry is  $a$  and the second entry is  $b$  in the two following ways,

$$\sum_{\{y \in Y_a\}} 1(y = b) = \sum_{i=1}^n 1(x_i = a, y_i = b). \quad (4.12)$$

Consequently,

$$p_{Y|X}(b|a) := \frac{1}{n_a} \sum_{\{y \in Y_a\}} 1(y = b) \quad (4.13)$$

$$= \frac{1}{np_X(a)} \sum_{i=1}^n 1(x_i = a, y_i = b) \quad (4.14)$$

$$= \frac{p_{X,Y}(a, b)}{p_X(a)}. \quad (4.15)$$

4.3 (Noisy data)

a)

$$P(\tilde{z} = 1 | \tilde{x} = 1) \quad (4.16)$$

$$= \frac{P(\tilde{z} = 1, \tilde{x} = 1)}{P(\tilde{x} = 1)} \quad (4.17)$$

$$= \frac{P(\tilde{z} = 1, \tilde{x} = 1, \tilde{y} = 0) + P(\tilde{z} = 1, \tilde{x} = 1, \tilde{y} = 1)}{P(\tilde{x} = 1)} \quad (4.18)$$

$$= \frac{P(\tilde{x} = 1)P(\tilde{y} = 0 | \tilde{x} = 1)P(\tilde{z} = 1 | \tilde{x} = 1, \tilde{y} = 0) + P(\tilde{x} = 1)P(\tilde{y} = 1 | \tilde{x} = 1)P(\tilde{z} = 1 | \tilde{x} = 1, \tilde{y} = 1)}{P(\tilde{x} = 1)} \quad (4.19)$$

$$= P(\tilde{y} = 0 | \tilde{x} = 1)P(\tilde{z} = 1 | \tilde{y} = 0) + P(\tilde{y} = 1 | \tilde{x} = 1)P(\tilde{z} = 1 | \tilde{x} = 1, \tilde{y} = 1) \quad (4.20)$$

$$= 0.82. \quad (4.21)$$

b)  $P(\tilde{z} = 1 | \tilde{x} = 1, \tilde{y} = 1) = P(\tilde{z} = 1 | \tilde{y} = 1) = 0.9 \neq P(\tilde{z} = 1 | \tilde{x} = 1)$ , so  $\tilde{y}$  and  $\tilde{z}$  are not conditionally independent given  $\tilde{x}$ .

c) By definition of the problem  $P(\tilde{z} = z | \tilde{x} = x, \tilde{y} = y) = P(\tilde{z} = z | \tilde{y} = y)$  for any values of  $x$ ,  $y$  and  $z$ , so  $\tilde{z}$  and  $\tilde{x}$  are conditionally independent given  $\tilde{y}$ .

4.4 (Fire alarm)

a)

$$p_{\bar{a}}(1) = \sum_{x \in \{0,1\}} \sum_{b \in \{0,1\}} p_{\bar{a},\bar{x},\bar{b}}(1, x, b) \quad (4.22)$$

$$= p_{\bar{a},\bar{x},\bar{b}}(1, 0, 0) + p_{\bar{a},\bar{x},\bar{b}}(1, 0, 1) + p_{\bar{a},\bar{x},\bar{b}}(1, 1, 0) + p_{\bar{a},\bar{x},\bar{b}}(1, 1, 1) \quad (4.23)$$

$$= p_{\bar{x}}(0)p_{\bar{b}}(0)p_{\bar{a}|\bar{x},\bar{b}}(1|0,0) + p_{\bar{x}}(0)p_{\bar{b}}(1)p_{\bar{a}|\bar{x},\bar{b}}(1|0,1) \quad (4.24)$$

$$+ p_{\bar{x}}(1)p_{\bar{b}}(0)p_{\bar{a}|\bar{x},\bar{b}}(1|1,0) + p_{\bar{x}}(1)p_{\bar{b}}(1)p_{\bar{a}|\bar{x},\bar{b}}(1|1,1) \quad (4.25)$$

$$= 0 + 0.8 \cdot 0.5 \cdot 0.2 + 0.2 \cdot 0.5 \cdot 0.8 + 0.2 \cdot 0.5 \cdot 0.4 \quad (4.26)$$

$$= 0.2. \quad (4.27)$$

b)

$$p_{\bar{x},\bar{b}|\bar{a}}(0, 0|1) = \frac{p_{\bar{a},\bar{x},\bar{b}}(1, 0, 0)}{p_{\bar{a}}(1)} \quad (4.28)$$

$$= \frac{p_{\bar{x}}(0)p_{\bar{b}}(0)p_{\bar{a}|\bar{x},\bar{b}}(1|0,0)}{p_{\bar{a}}(1)} \quad (4.29)$$

$$= 0, \quad (4.30)$$

$$p_{\bar{x},\bar{b}|\bar{a}}(0, 1|1) = \frac{p_{\bar{a},\bar{x},\bar{b}}(1, 0, 1)}{p_{\bar{a}}(1)} \quad (4.31)$$

$$= \frac{p_{\bar{x}}(0)p_{\bar{b}}(1)p_{\bar{a}|\bar{x},\bar{b}}(1|0,1)}{p_{\bar{a}}(1)} \quad (4.32)$$

$$= \frac{0.8 \cdot 0.5 \cdot 0.2}{0.2} \quad (4.33)$$

$$= 0.4, \quad (4.34)$$

$$p_{\bar{x},\bar{b}|\bar{a}}(1, 0|1) = \frac{p_{\bar{a},\bar{x},\bar{b}}(1, 1, 0)}{p_{\bar{a}}(1)} \quad (4.35)$$

$$= \frac{p_{\bar{x}}(1)p_{\bar{b}}(0)p_{\bar{a}|\bar{x},\bar{b}}(1|1,0)}{p_{\bar{a}}(1)} \quad (4.36)$$

$$= \frac{0.2 \cdot 0.5 \cdot 0.8}{0.2} \quad (4.37)$$

$$= 0.4, \quad (4.38)$$

$$p_{\bar{x},\bar{b}|\bar{a}}(1, 1|1) = \frac{p_{\bar{a},\bar{x},\bar{b}}(1, 1, 1)}{p_{\bar{a}}(1)} \quad (4.39)$$

$$= \frac{p_{\bar{x}}(1)p_{\bar{b}}(1)p_{\bar{a}|\bar{x},\bar{b}}(1|1,1)}{p_{\bar{a}}(1)} \quad (4.40)$$

$$= \frac{0.2 \cdot 0.5 \cdot 0.4}{0.2} \quad (4.41)$$

$$= 0.2. \quad (4.42)$$

c) We compare  $p_{\tilde{x}|\tilde{a},\tilde{b}}(0|1,0)$  and  $p_{\tilde{x}|\tilde{a}}(0|1)$ ,

$$p_{\tilde{x}|\tilde{a},\tilde{b}}(0|1,0) = \frac{p_{\tilde{a},\tilde{x},\tilde{b}}(1,0,0)}{p_{\tilde{a},\tilde{b}}(1,0)} \quad (4.43)$$

$$= \frac{p_{\tilde{x}}(0)p_{\tilde{b}}(0)p_{\tilde{a}|\tilde{x},\tilde{b}}(1|0,0)}{\sum_{x=0}^1 p_{\tilde{x}}(x)p_{\tilde{b}}(0)p_{\tilde{a}|\tilde{x},\tilde{b}}(1|x,0)} \quad (4.44)$$

$$= 0, \quad (4.45)$$

$$p_{\tilde{x}|\tilde{a}}(0|1) = \frac{p_{\tilde{a},\tilde{x}}(1,0)}{p_{\tilde{a}}(1)} \quad (4.46)$$

$$= \frac{\sum_{b=0}^1 p_{\tilde{a},\tilde{x},\tilde{b}}(1,0,b)}{p_{\tilde{a}}(1)} \quad (4.47)$$

$$= \frac{\sum_{b=0}^1 p_{\tilde{x}}(0)p_{\tilde{b}}(b)p_{\tilde{a}|\tilde{x},\tilde{b}}(1|0,b)}{p_{\tilde{a}}(1)} \quad (4.48)$$

$$= \frac{0.8 \cdot 0.5 \cdot 0.2}{0.2} = 0.4 \neq 0, \quad (4.49)$$

which implies that  $\tilde{x}$  and  $\tilde{b}$  are not conditionally independent given  $\tilde{a}$ . Intuitively, if we know that the alarm is ringing, then the state of the battery does provide information about fire. For example, if the battery is fine, then the alarm is not ringing for that reason, so it is more likely that there is a fire.

#### 4.5 (Volcano eruption)

a)

$$p_{\tilde{v}|\tilde{s}_1,\tilde{s}_2,\tilde{s}_3}(1|1,1,1) \quad (4.50)$$

$$= \frac{p_{\tilde{v},\tilde{s}_1,\tilde{s}_2,\tilde{s}_3}(1,1,1,1)}{p_{\tilde{s}_1,\tilde{s}_2,\tilde{s}_3}(1,1,1,1)} \quad (4.51)$$

$$= \frac{p_{\tilde{v},\tilde{s}_1,\tilde{s}_2,\tilde{s}_3}(1,1,1,1)}{p_{\tilde{v},\tilde{s}_1,\tilde{s}_2,\tilde{s}_3}(1,1,1,1) + p_{\tilde{v},\tilde{s}_1,\tilde{s}_2,\tilde{s}_3}(0,1,1,1)} \quad (4.52)$$

$$= \frac{p_{\tilde{v}}(1)p_{\tilde{s}_1|\tilde{v}}(1|1)p_{\tilde{s}_2|\tilde{v}}(1|1)p_{\tilde{s}_3|\tilde{v}}(1|1)}{p_{\tilde{v}}(1)p_{\tilde{s}_1|\tilde{v}}(1|1)p_{\tilde{s}_2|\tilde{v}}(1|1)p_{\tilde{s}_3|\tilde{v}}(1|1) + p_{\tilde{v}}(0)p_{\tilde{s}_1|\tilde{v}}(1|0)p_{\tilde{s}_2|\tilde{v}}(1|0)p_{\tilde{s}_3|\tilde{v}}(1|0)} \quad (4.53)$$

$$= \frac{0.1 \cdot 0.5 \cdot 0.8 \cdot 0.9}{0.1 \cdot 0.5 \cdot 0.8 \cdot 0.9 + 0.9 \cdot 0.2 \cdot 0.2 \cdot 0.2} \quad (4.54)$$

$$= \frac{5}{5+1} = \frac{5}{6}. \quad (4.54)$$

b) We compare the marginal probability that the first sensor is activated,

$$p_{\tilde{s}_1}(1) = p_{\tilde{v},\tilde{s}_1}(1,1) + p_{\tilde{v},\tilde{s}_1}(0,1) \quad (4.55)$$

$$= p_{\tilde{v}}(1)p_{\tilde{s}_1|\tilde{v}}(1|1) + p_{\tilde{v}}(0)p_{\tilde{s}_1|\tilde{v}}(1|0) \quad (4.56)$$

$$= 0.9 \cdot 0.2 + 0.1 \cdot 0.5 = 0.23, \quad (4.57)$$

and the conditional probability that the first sensor is activated given that the second

is activated,

$$p_{\tilde{s}_1 | \tilde{s}_2}(1 | 1) = \frac{p_{\tilde{s}_1, \tilde{s}_2}(1, 1)}{p_{\tilde{s}_2}(1)} \quad (4.58)$$

$$= \frac{p_{\tilde{v}}(1) p_{\tilde{s}_1 | \tilde{v}}(1 | 1) p_{\tilde{s}_2 | \tilde{v}}(1 | 1) + p_{\tilde{v}}(0) p_{\tilde{s}_1 | \tilde{v}}(1 | 0) p_{\tilde{s}_2 | \tilde{v}}(1 | 0)}{p_{\tilde{v}}(1) p_{\tilde{s}_2 | \tilde{v}}(1 | 1) + p_{\tilde{v}}(0) p_{\tilde{s}_2 | \tilde{v}}(1 | 0)} \quad (4.59)$$

$$= \frac{0.9 \cdot 0.2 \cdot 0.2 + 0.1 \cdot 0.5 \cdot 0.8}{0.9 \cdot 0.2 + 0.1 \cdot 0.8} \quad (4.60)$$

$$= 0.29. \quad (4.61)$$

Since they are different, the two random variables are not independent. Intuitively, if the second sensor is activated then an eruption is more likely, which makes it more likely for the first sensor to be activated. Therefore the state of the second sensor provides information about the state of the first sensor.

#### 4.6 (Flu and COVID-19)

a)

$$p_{\tilde{t}}(1) = \sum_{c \in \{0,1\}} \sum_{x \in \{0,1\}} \sum_{s \in \{0,1\}} p_{\tilde{s}, \tilde{t}, \tilde{c}, \tilde{x}}(s, 1, c, x) \quad (4.62)$$

$$= \sum_{c \in \{0,1\}} \sum_{x \in \{0,1\}} \sum_{s \in \{0,1\}} p_{\tilde{c}}(c) p_{\tilde{x}}(x) p_{\tilde{s}, \tilde{t} | \tilde{c}, \tilde{x}}(s, 1 | c, x) \quad (4.63)$$

$$= 0.2 \cdot 0.9(0.2 + 0.4) + 0.8 \cdot 0.1(0.4 + 0.4) + 0.2 \cdot 0.1 \quad (4.64)$$

$$= 0.192. \quad (4.65)$$

b)

$$p_{\tilde{c} | \tilde{t}}(1 | 1) = \frac{p_{\tilde{c}, \tilde{t}}(1, 1)}{p_{\tilde{t}}(1)} \quad (4.66)$$

$$= \frac{\sum_{x \in \{0,1\}} \sum_{s \in \{0,1\}} p_{\tilde{s}, \tilde{t}, \tilde{c}, \tilde{x}}(s, 1, 1, x)}{p_{\tilde{t}}(1)} \quad (4.67)$$

$$= \frac{\sum_{x \in \{0,1\}} \sum_{s \in \{0,1\}} p_{\tilde{c}}(1) p_{\tilde{x}}(x) p_{\tilde{s}, \tilde{t} | \tilde{c}, \tilde{x}}(s, 1 | 1, x)}{p_{\tilde{t}}(1)} \quad (4.68)$$

$$= \frac{0.2 \cdot 0.9(0.2 + 0.4) + 0.2 \cdot 0.1}{0.192} \quad (4.69)$$

$$= 0.67. \quad (4.70)$$

#### 4.7 (Vibrations)

a) By conditional independence of  $\tilde{s}$  and  $\tilde{e}$  given  $\tilde{v}$

$$p_{\tilde{s}}(1) = \sum_{e=0}^1 \sum_{v=0}^2 p_{\tilde{e}, \tilde{v}}(e, v) p_{\tilde{s} | \tilde{e}, \tilde{v}}(1 | e, v) \quad (4.71)$$

$$= \sum_{e=0}^1 \sum_{v=0}^2 p_{\tilde{e}, \tilde{v}}(e, v) p_{\tilde{s} | \tilde{v}}(1 | v) \quad (4.72)$$

$$= p_{\tilde{s} | \tilde{v}}(1 | 1)(p_{\tilde{e}, \tilde{v}}(0, 1) + p_{\tilde{e}, \tilde{v}}(1, 1)) + p_{\tilde{s} | \tilde{v}}(1 | 2)(p_{\tilde{e}, \tilde{v}}(0, 2) + p_{\tilde{e}, \tilde{v}}(1, 2)) \quad (4.73)$$

$$= 0.5(0.05 + 0.05) + 0.1 \quad (4.74)$$

$$= 0.15 \quad (4.75)$$



and consequently  $p_{\bar{s}}(0) = 1 - p_{\bar{s}}(1) = 0.85$ .

b)

$$p_{\bar{e}|\bar{s}}(1|1) = \frac{p_{\bar{e},\bar{s}}(1,1)}{p_{\bar{s}}(1)} \quad (4.76)$$

$$= \frac{\sum_{v=0}^2 p_{\bar{e},\bar{v},\bar{s}}(1,v,1)}{p_{\bar{s}}(1)} \quad (4.77)$$

$$= \frac{\sum_{v=0}^2 p_{\bar{e},\bar{v},\bar{s}}(1,v)p_{\bar{s}|\bar{v}}(1|v)}{p_{\bar{s}}(1)} \quad (4.78)$$

$$= \frac{p_{\bar{s}|\bar{v}}(1|1)p_{\bar{e},\bar{v}}(1,1) + p_{\bar{s}|\bar{v}}(1|2)p_{\bar{e},\bar{v}}(1,2)}{p_{\bar{s}}(1)} \quad (4.79)$$

$$= \frac{0.5 \cdot 0.05 + 0.1}{0.15} \quad (4.80)$$

$$= 0.833. \quad (4.81)$$

c) We have

$$p_{\bar{e}}(1) = \sum_{v=0}^2 p_{\bar{e},\bar{v}}(1,v) \quad (4.82)$$

$$= 0.05 + 0.1 \quad (4.83)$$

$$= 0.15 \neq p_{\bar{e}|\bar{s}}(1|1), \quad (4.84)$$

so they are not independent. This makes sense, because the sensor reading is more likely to be one if there are vibrations, which is more likely if there is an earthquake.

4.8 (Footprints)

a) The joint pmf  $p_{\tilde{a},\tilde{s},\tilde{h}}(a,s,h)$  is obtained by dividing the counts by the total number of data (200):

		s					s		
h	a = Coyote	Small	Medium	Large	h	a = Wolf	Small	Medium	Large
	Oval	0.4	0.15	0.025		Oval	0	0.025	0.05
	Circular	0.15	0.025	0		Circular	0.025	0.075	0.075

b)

$$p_{\tilde{a}|\tilde{s},\tilde{h}}(\text{Wolf}|\text{Med},\text{Circular}) = \frac{p_{\tilde{a},\tilde{s},\tilde{h}}(\text{Wolf},\text{Med},\text{Circular})}{p_{\tilde{s},\tilde{h}}(\text{Med},\text{Circular})} \quad (4.85)$$

$$= \frac{p_{\tilde{a},\tilde{s},\tilde{h}}(\text{Wolf},\text{Med},\text{Circular})}{p_{\tilde{a},\tilde{s},\tilde{h}}(\text{Wolf},\text{Med},\text{Circular}) + p_{\tilde{a},\tilde{s},\tilde{h}}(\text{Coyote},\text{Med},\text{Circular})} \quad (4.86)$$

$$= \frac{0.075}{0.075 + 0.025} = 0.75. \quad (4.87)$$

4.9 (Election)

- a) Candidate A wins if they win 2 or 3 states. Consequently, under the independence assumptions,

$$p_{\bar{o}}(1) = P(\text{A wins all states}) + \sum_{i=1}^3 P(\text{A wins all states except } i) \quad (4.88)$$

$$= 0.6^3 + 3 \cdot 0.4 \cdot 0.6^2 \quad (4.89)$$

$$= 0.648. \quad (4.90)$$

- b) By the definition of conditional pmf,

$$p_{\bar{s}_2 | \bar{o}}(0 | 0) = \frac{p_{\bar{s}_2, \bar{o}}(0, 0)}{p_{\bar{o}}(0)} \quad (4.91)$$

$$= \frac{\sum_{s_1=0}^1 \sum_{s_3=0}^1 p_{\bar{s}_1, \bar{s}_2, \bar{s}_3, \bar{o}}(s_1, 0, s_3, 0)}{p_{\bar{o}}(0)} \quad (4.92)$$

$$= \frac{p_{\bar{s}_1, \bar{s}_2, \bar{s}_3, \bar{o}}(0, 0, 0, 0) + p_{\bar{s}_1, \bar{s}_2, \bar{s}_3, \bar{o}}(1, 0, 0, 0) + p_{\bar{s}_1, \bar{s}_2, \bar{s}_3, \bar{o}}(0, 0, 1, 0)}{p_{\bar{o}}(0)} \quad (4.93)$$

$$= \frac{p_{\bar{s}_1}(0)p_{\bar{s}_2}(0)p_{\bar{s}_3}(0) + p_{\bar{s}_1}(1)p_{\bar{s}_2}(0)p_{\bar{s}_3}(0) + p_{\bar{s}_1}(0)p_{\bar{s}_2}(0)p_{\bar{s}_3}(1)}{p_{\bar{o}}(0)} \quad (4.94)$$

$$= \frac{0.4^3 + 2 \cdot 0.4^2 \cdot 0.6}{1 - 0.648} \quad (4.95)$$

$$= 0.727. \quad (4.96)$$

- c) We have

$$p_{\bar{s}_1 | \bar{o}}(1 | 1) \quad (4.97)$$

$$= \frac{p_{\bar{s}_1, \bar{o}}(1, 1)}{p_{\bar{o}}(1)} \quad (4.98)$$

$$= \frac{P(\text{A wins all states}) + P(\text{A wins 1 and 2 but not 3}) + P(\text{A wins 1 and 3 but not 2})}{p_{\bar{o}}(1)}$$

$$= \frac{0.6^3 + 2 \cdot 0.4 \cdot 0.6^2}{0.648} \quad (4.99)$$

$$= 0.778. \quad (4.100)$$

However,  $p_{\bar{s}_1 | \bar{o}, \bar{s}_2}(1 | 1, 0) = 1$  because if A wins the election then they cannot lose state 1 and state 2. Intuitively, even if the state results are independent, they both determine the result of the election, so revealing who wins the election *connects* them. For example, if candidate A has won the election but lost state 2, then this completely determines the result of state 1.

a) By the independence assumptions

$$p_{\tilde{d}_1, \tilde{d}_2 | \tilde{q}}(1, 1 | 1) = \frac{P(\tilde{e}_1 \tilde{q} = 1, \tilde{e}_2 \tilde{q} = 1, \tilde{q} = 1)}{P(\tilde{q} = 1)} \quad (4.101)$$

$$= \frac{P(\tilde{q} = 1)P(\tilde{e}_1 \tilde{q} = 1, \tilde{e}_2 \tilde{q} = 1 | \tilde{q} = 1)}{P(\tilde{q} = 1)} \quad (4.102)$$

$$= P(\tilde{e}_1 \tilde{q} = 1, \tilde{e}_2 \tilde{q} = 1 | \tilde{q} = 1) \quad (4.103)$$

$$= P(\tilde{e}_1 = 1, \tilde{e}_2 = 1 | \tilde{q} = 1) \quad (4.104)$$

$$= P(\tilde{e}_1 = 1, \tilde{e}_2 = 1) \quad (4.105)$$

$$= P(\tilde{e}_1 = 1)P(\tilde{e}_2 = 1) \quad (4.106)$$

$$= 0.64. \quad (4.107)$$

b)

$$p_{\tilde{q} | \tilde{d}_1, \tilde{d}_2}(1 | 1, 1) = \frac{p_{\tilde{q}, \tilde{d}_1, \tilde{d}_2}(1, 1, 1)}{p_{\tilde{d}_1, \tilde{d}_2}(1, 1)}. \quad (4.108)$$

The numerator equals

$$p_{\tilde{q}, \tilde{d}_1, \tilde{d}_2}(1, 1, 1) = P(\tilde{e}_1 \tilde{q} = 1, \tilde{e}_2 \tilde{q} = 1 | \tilde{q} = 1) \quad (4.109)$$

$$= P(\tilde{q} = 1)P(\tilde{e}_1 = 1, \tilde{e}_2 = 1 | \tilde{q} = 1) \quad (4.110)$$

$$= P(\tilde{q} = 1)P(\tilde{e}_1 = 1, \tilde{e}_2 = 1) \quad (4.111)$$

$$= P(\tilde{q} = 1)P(\tilde{e}_1 = 1)P(\tilde{e}_2 = 1) \quad (4.112)$$

$$= 0.16. \quad (4.113)$$

The denominator equals

$$p_{\tilde{d}_1, \tilde{d}_2}(1, 1) = P(\tilde{e}_1 \tilde{q} = 1, \tilde{e}_2 \tilde{q} = 1) \quad (4.114)$$

$$= P(\tilde{e}_1 = 1, \tilde{e}_2 = 1, \tilde{q} = 1) + P(\tilde{e}_1 = -1, \tilde{e}_2 = -1, \tilde{q} = -1) \quad (4.115)$$

$$= P(\tilde{e}_1 = 1)P(\tilde{e}_2 = 1)P(\tilde{q} = 1) + P(\tilde{e}_1 = -1)P(\tilde{e}_2 = -1)P(\tilde{q} = -1) \quad (4.116)$$

$$= 0.19. \quad (4.116)$$

Consequently,

$$p_{\tilde{q} | \tilde{d}_1, \tilde{d}_2}(1 | 1, 1) = 0.84. \quad (4.117)$$

c) As shown above  $p_{\tilde{d}_1, \tilde{d}_2}(1, 1) = 0.19$ , but

$$p_{\tilde{d}_1}(1) = P(\tilde{e}_1 \tilde{q} = 1) \quad (4.118)$$

$$= P(\tilde{e}_1 = 1, \tilde{q} = 1) + P(\tilde{e}_1 = -1, \tilde{q} = -1) \quad (4.119)$$

$$= P(\tilde{e}_1 = 1)P(\tilde{q} = 1) + P(\tilde{e}_1 = -1)P(\tilde{q} = -1) \quad (4.120)$$

$$= 0.35, \quad (4.121)$$

and by the same reasoning,

$$p_{\tilde{d}_2}(1) = 0.35. \quad (4.122)$$

Since  $0.19 \neq 0.35^2$ , the random variables are not independent.

d) For any  $x_1, x_2$  and  $q$  in  $\{-1, 1\}$ ,

$$p_{\tilde{d}_1, \tilde{d}_2 | \tilde{q}}(x_1, x_2 | q) = P(\tilde{e}_1 \tilde{q} = x_1, \tilde{e}_2 \tilde{q} = x_2 | \tilde{q} = q) \quad (4.123)$$

$$= P(\tilde{e}_1 = x_1/q, \tilde{e}_2 = x_2/q | \tilde{q} = q) \quad (4.124)$$

$$= P(\tilde{e}_1 = x_1/q)P(\tilde{e}_2 = x_2/q) \quad \text{by independence} \quad (4.125)$$

$$= P(\tilde{e}_1 = x_1/q | \tilde{q} = q)P(\tilde{e}_2 = x_2/q | \tilde{q} = q) \quad \text{by independence}$$

$$= P(\tilde{e}_1 \tilde{q} = x_1 | \tilde{q} = q)P(\tilde{e}_2 \tilde{q} = x_2 | \tilde{q} = q) \quad (4.126)$$

$$= p_{\tilde{d}_1 | \tilde{q}}(x_1 | q)p_{\tilde{d}_2 | \tilde{q}}(x_2 | q), \quad (4.127)$$

so  $\tilde{d}_1$  and  $\tilde{d}_2$  are conditionally independent given  $\tilde{q}$ .

#### 4.11 (Surgery)

a) Let us define the potential outcomes  $\widetilde{\text{po}}_A$  and  $\widetilde{\text{po}}_B$ .  $\widetilde{\text{po}}_A = 1$  and  $\widetilde{\text{po}}_B = 1$  indicate recovery after procedure A and B respectively.  $\widetilde{\text{po}}_A = 0$  and  $\widetilde{\text{po}}_B = 0$  indicate non-recovery. The random variable  $\tilde{y}$  represents the observed outcome and  $\tilde{t}$  the treatment, so that

$$\tilde{y} := \begin{cases} \widetilde{\text{po}}_A & \text{if } \tilde{t} = A, \\ \widetilde{\text{po}}_B & \text{if } \tilde{t} = B. \end{cases} \quad (4.128)$$

We also define a random variable  $\tilde{d}$  to indicate whether each case is mild ( $\tilde{d} = m$ ) or serious ( $\tilde{d} = s$ ). We denote the probability that a case is mild if the treatment is A or B by  $\alpha_A$  and  $\alpha_B$  respectively. The observed probability of recovery for patients undergoing procedure A is

$$P(\tilde{y} = 1 | \tilde{t} = A) = P(\widetilde{\text{po}}_A = 1 | \tilde{t} = A) \quad (4.129)$$

$$= P(\widetilde{\text{po}}_A = 1, \tilde{d} = m | \tilde{t} = A) + P(\widetilde{\text{po}}_A = 1, \tilde{d} = s | \tilde{t} = A) \quad (4.130)$$

$$= P(\tilde{d} = m | \tilde{t} = A) P(\widetilde{\text{po}}_A = 1 | \tilde{t} = A, \tilde{d} = m) \quad (4.131)$$

$$+ P(\tilde{d} = s | \tilde{t} = A) P(\widetilde{\text{po}}_A = 1 | \tilde{t} = A, \tilde{d} = s) \quad (4.132)$$

$$= 0.9\alpha_A + 0.5(1 - \alpha_A). \quad (4.133)$$

Setting this equal to 0.58 yields  $\alpha_A = 0.2$ .

Similarly,

$$P(\tilde{y} = 1 | \tilde{t} = B) = P(\widetilde{\text{po}}_B = 1 | \tilde{t} = B) \quad (4.134)$$

$$= P(\widetilde{\text{po}}_B = 1, \tilde{d} = m | \tilde{t} = B) + P(\widetilde{\text{po}}_B = 1, \tilde{d} = s | \tilde{t} = B) \quad (4.135)$$

$$= P(\tilde{d} = m | \tilde{t} = B) P(\widetilde{\text{po}}_B = 1 | \tilde{t} = B, \tilde{d} = m) \quad (4.136)$$

$$+ P(\tilde{d} = s | \tilde{t} = B) P(\widetilde{\text{po}}_B = 1 | \tilde{t} = B, \tilde{d} = s) \quad (4.137)$$

$$= 0.8\alpha_B + 0.2(1 - \alpha_B). \quad (4.138)$$

Setting this equal to 0.68 yields  $\alpha_B = 0.8$ .

What is happening is that the fraction of patients with mild cases is much lower for procedure A (20%) than for procedure B (80%). Since those patients are more likely to recover, irrespective of the procedure, this inflates the recovery rate for procedure B.

b) We can correct for the confounder by taking into account whether each case is mild or serious. Under the assumption that the treatment  $\tilde{t}$  and the potential outcomes  $\widetilde{\text{po}}_A$

and  $\widetilde{\text{po}}_B$  are conditionally independent given the degree of severity  $\tilde{d}$ , then

$$P(\tilde{y} | \tilde{t} = A, \tilde{d} = m) = P(\widetilde{\text{po}}_A | \tilde{t} = A, \tilde{d} = m) \quad (4.139)$$

$$= P(\widetilde{\text{po}}_A | \tilde{d} = m). \quad (4.140)$$

By the same argument,

$$P(\tilde{y} = 1 | \tilde{t} = A, \tilde{d} = s) = P(\widetilde{\text{po}}_A = 1 | \tilde{d} = s), \quad (4.141)$$

$$P(\tilde{y} = 1 | \tilde{t} = B, \tilde{d} = m) = P(\widetilde{\text{po}}_B = 1 | \tilde{d} = m), \quad (4.142)$$

$$P(\tilde{y} = 1 | \tilde{t} = B, \tilde{d} = s) = P(\widetilde{\text{po}}_B = 1 | \tilde{d} = s). \quad (4.143)$$

Consequently, we can compute the *true* efficacy as follows,

$$P(\widetilde{\text{po}}_A) = \sum_{d \in \{m, s\}} P(\widetilde{\text{po}}_A = 1 | \tilde{d} = d) P(\tilde{d} = d) \quad (4.144)$$

$$= \sum_{d \in \{m, s\}} P(\tilde{y} = 1 | \tilde{t} = A, \tilde{d} = d) P(\tilde{d} = d), \quad (4.145)$$

$$P(\widetilde{\text{po}}_B) = \sum_{d \in \{m, s\}} P(\widetilde{\text{po}}_B = 1 | \tilde{d} = d) P(\tilde{d} = d) \quad (4.146)$$

$$= \sum_{d \in \{m, s\}} P(\tilde{y} = 1 | \tilde{t} = B, \tilde{d} = d) P(\tilde{d} = d). \quad (4.147)$$

This only works if the conditional independence assumption holds. Intuitively, we are assuming that once we control for the severity, there are no other systematic differences between the patients that undergo each of the procedures.

- c) Randomizing what patient undergoes each procedure would enable us to neutralize all confounding factors, even if we don't know what they are.

#### 4.12 (Admissions)

- a) The empirical conditional probabilities are

$$P(\text{admitted} | \text{men}) = \frac{40 + 25}{500 + 500} \quad (4.148)$$

$$= 0.065, \quad (4.149)$$

$$P(\text{admitted} | \text{women}) = \frac{50 + 25}{800 + 200} \quad (4.150)$$

$$= 0.075. \quad (4.151)$$

b) By the assumptions,

$$P(\widetilde{\text{po}}_{\text{man}} = 1) = P(\widetilde{\text{po}}_{\text{man}} = 1, \tilde{x} = \text{Med}) \quad (4.152)$$

$$+ P(\widetilde{\text{po}}_{\text{man}} = 1, \tilde{x} = \text{Art}) \quad (4.153)$$

$$= P(\tilde{x} = \text{Med}) P(\widetilde{\text{po}}_{\text{man}} = 1 \mid \tilde{x} = \text{Med}) \quad (4.154)$$

$$+ P(\tilde{x} = \text{Art}) P(\widetilde{\text{po}}_{\text{man}} = 1 \mid \tilde{x} = \text{Art}) \quad (4.155)$$

$$= P(\tilde{x} = \text{Med}) P(\widetilde{\text{po}}_{\text{man}} = 1 \mid \tilde{x} = \text{Med}, \tilde{s} = \text{man}) \quad (4.156)$$

$$+ P(\tilde{x} = \text{Art}) P(\widetilde{\text{po}}_{\text{man}} = 1 \mid \tilde{x} = \text{Art}, \tilde{s} = \text{man}) \quad (4.157)$$

$$= P(\tilde{x} = \text{Med}) P(\tilde{y} = 1 \mid \tilde{x} = \text{Med}, \tilde{s} = \text{man}) \quad (4.158)$$

$$+ P(\tilde{x} = \text{Art}) P(\tilde{y} = 1 \mid \tilde{x} = \text{Art}, \tilde{s} = \text{man}) \quad (4.159)$$

$$= \frac{1300}{2000} \frac{40}{800} + \frac{700}{2000} \frac{25}{200} \quad (4.160)$$

$$= 0.076, \quad (4.161)$$

$$P(\widetilde{\text{po}}_{\text{woman}} = 1) = P(\widetilde{\text{po}}_{\text{man}} = 1, \tilde{x} = \text{Med}) \quad (4.162)$$

$$+ P(\widetilde{\text{po}}_{\text{woman}} = 1, \tilde{x} = \text{Art}) \quad (4.163)$$

$$= P(\tilde{x} = \text{Med}) P(\widetilde{\text{po}}_{\text{woman}} = 1 \mid \tilde{x} = \text{Med}) \quad (4.164)$$

$$+ P(\tilde{x} = \text{Art}) P(\widetilde{\text{po}}_{\text{woman}} = 1 \mid \tilde{x} = \text{Art}) \quad (4.165)$$

$$= P(\tilde{x} = \text{Med}) P(\widetilde{\text{po}}_{\text{woman}} = 1 \mid \tilde{x} = \text{Med}, \tilde{s} = \text{woman}) \quad (4.166)$$

$$+ P(\tilde{x} = \text{Art}) P(\widetilde{\text{po}}_{\text{woman}} = 1 \mid \tilde{x} = \text{Art}, \tilde{s} = \text{woman}) \quad (4.167)$$

$$= P(\tilde{x} = \text{Med}) P(\tilde{y} = 1 \mid \tilde{x} = \text{Med}, \tilde{s} = \text{woman}) \quad (4.168)$$

$$+ P(\tilde{x} = \text{Art}) P(\tilde{y} = 1 \mid \tilde{x} = \text{Art}, \tilde{s} = \text{woman}) \quad (4.169)$$

$$= \frac{1300}{2000} \frac{25}{500} + \frac{700}{2000} \frac{50}{500} \quad (4.170)$$

$$= 0.068. \quad (4.171)$$

Once we adjust by the school they apply to, men are more likely to be admitted than women, so the data do not support their hypothesis.

#### 4.13 (Shop)

a) From the data the probability of a purchase if there is music is  $120/160=0.75$ . If there is no music it is  $109/190 = 0.57$ .

b) We define random variables  $\tilde{a}$  and  $\tilde{m}$  to represent age and music respectively:

$$p_{\widetilde{\text{po}}_0}(1) = \sum_{a \in \{\text{young}, \text{mid}, \text{old}\}} p_{\tilde{a}}(a) p_{\widetilde{\text{po}}_0 \mid \tilde{a}}(1 \mid a) \quad (4.172)$$

$$= \sum_{a \in \{\text{young}, \text{mid}, \text{old}\}} p_{\tilde{a}}(a) p_{\widetilde{\text{po}}_0 \mid \tilde{a}, \tilde{m}}(1 \mid a, 0) \quad (4.173)$$

$$= \sum_{a \in \{\text{young}, \text{mid}, \text{old}\}} p_{\tilde{a}}(a) p_{\tilde{y} \mid \tilde{a}, \tilde{m}}(1 \mid a, 0) \quad (4.174)$$

$$= \frac{100}{350} \cdot \frac{10}{20} + \frac{110}{350} \cdot \frac{90}{100} + \frac{140}{350} \cdot \frac{20}{40} \quad (4.175)$$

$$= 0.626. \quad (4.176)$$

$$p_{\widetilde{\text{po}}_1}(1) = \sum_{a \in \{\text{young}, \text{mid}, \text{old}\}} p_{\tilde{a}}(a) p_{\widetilde{\text{po}}_1 | \tilde{a}}(1 | a) \quad (4.177)$$

$$= \sum_{a \in \{\text{young}, \text{mid}, \text{old}\}} p_{\tilde{a}}(a) p_{\widetilde{\text{po}}_1 | \tilde{a}, \tilde{m}}(1 | a, 1) \quad (4.178)$$

$$= \sum_{a \in \{\text{young}, \text{mid}, \text{old}\}} p_{\tilde{a}}(a) p_{\tilde{y} | \tilde{a}, \tilde{m}}(1 | a, 1) \quad (4.179)$$

$$= \frac{100}{350} \cdot \frac{40}{80} + \frac{110}{350} \cdot \frac{9}{10} + \frac{140}{350} \cdot \frac{60}{100} \quad (4.180)$$

$$= 0.666. \quad (4.181)$$

Playing music decreases purchases under the conditional-independence assumption.

#### 4.14 (Rackets)

- a) True. Let  $\tilde{s}$  be a random variable that represents the surface. If  $\alpha := P(\tilde{s} = \text{clay} | \tilde{r} = A)$  equals  $P(\tilde{s} = \text{clay})$ , then  $1 - \alpha := P(\tilde{s} = \text{grass} | \tilde{r} = A)$  equals  $P(\tilde{s} = \text{grass})$ , so that

$$P(\tilde{y} = 1 | \tilde{r} = A) = \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{y} = 1, \tilde{s} = s | \tilde{r} = A) \quad (4.182)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\widetilde{\text{po}}_A = 1, \tilde{s} = s | \tilde{r} = A) \quad (4.183)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{s} = s | \tilde{r} = A) P(\widetilde{\text{po}}_A = 1 | \tilde{r} = A, \tilde{s} = s) \quad (4.184)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{s} = s | \tilde{r} = A) P(\widetilde{\text{po}}_A = 1 | \tilde{s} = s) \quad (4.185)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{s} = s) P(\widetilde{\text{po}}_A = 1 | \tilde{s} = s) \quad (4.186)$$

$$= P(\widetilde{\text{po}}_A = 1). \quad (4.187)$$

b)

$$P(\tilde{y} = 1 | \tilde{r} = A) = \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{y} = 1, \tilde{s} = s | \tilde{r} = A) \quad (4.188)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\widetilde{\text{po}}_A = 1, \tilde{s} = s | \tilde{r} = A) \quad (4.189)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{s} = s | \tilde{r} = A) P(\widetilde{\text{po}}_A = 1 | \tilde{r} = A, \tilde{s} = s) \quad (4.190)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{s} = s | \tilde{r} = A) P(\widetilde{\text{po}}_A = 1 | \tilde{s} = s) \quad (4.191)$$

$$= 0.8\alpha + 0.2(1 - \alpha) = 0.2 + 0.6\alpha. \quad (4.192)$$

Similarly,

$$P(\tilde{y} = 1 \mid \tilde{r} = B) = \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{y} = 1, \tilde{s} = s \mid \tilde{r} = B) \quad (4.193)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\widetilde{\text{po}}_A = 1, \tilde{s} = s \mid \tilde{r} = B) \quad (4.194)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{s} = s \mid \tilde{r} = A) P(\widetilde{\text{po}}_B = 1 \mid \tilde{r} = B, \tilde{s} = s) \quad (4.195)$$

$$= \sum_{s \in \{\text{clay}, \text{grass}\}} P(\tilde{s} = s \mid \tilde{r} = A) P(\widetilde{\text{po}}_B = 1 \mid \tilde{s} = s) \quad (4.196)$$

$$= 0.7 \cdot 0.5 + 0.1 \cdot 0.5 = 0.4. \quad (4.197)$$

We need  $0.2 + 0.6\alpha < 0.4$ , which occurs for  $\alpha < 1/3$ .

#### 4.15 (Missing data)

- No, otherwise there wouldn't be such a large discrepancy between the number of missing data for men and women.
- 70 out of 140 observed cases have side effects, so the probability is 0.5.
- Let  $\tilde{\text{se}}$ ,  $\tilde{o}$  and  $\tilde{s}$  be Bernoulli random variables representing the side effects, whether the data are observed ( $\tilde{o} = 1$ ) or not ( $\tilde{o} = 0$ ) and the sex respectively. We have

$$p_{\tilde{\text{se}}}(1) = \sum_{s \in \{\text{man}, \text{woman}\}} p_{\tilde{s}}(s) p_{\tilde{\text{se}} \mid \tilde{s}}(1 \mid s). \quad (4.198)$$

If the side effects are conditionally independent from the data being observed or missing given the patient sex, then

$$p_{\tilde{\text{se}} \mid \tilde{s}}(1 \mid \text{man}) = p_{\tilde{\text{se}} \mid \tilde{o}, \tilde{s}}(1 \mid 1, \text{man}) \quad (4.199)$$

$$= \frac{40}{50} = 0.8, \quad (4.200)$$

$$p_{\tilde{\text{se}} \mid \tilde{s}}(1 \mid \text{woman}) = p_{\tilde{\text{se}} \mid \tilde{o}, \tilde{s}}(1 \mid 1, \text{woman}) \quad (4.201)$$

$$= \frac{30}{90} = 0.33. \quad (4.202)$$

Consequently, since  $p_{\tilde{s}}(\text{man}) = p_{\tilde{s}}(\text{woman}) = \frac{1}{2}$ ,

$$p_{\tilde{\text{se}}}(1) = \frac{0.8 + 0.33}{2} = 0.57. \quad (4.203)$$

#### 4.16 (Three players)

- We cannot because there are no games where all three players are absent, so we cannot compute what fraction are wins or losses.
- Let us represent the presence or absence of the players using a random vector  $\tilde{x}$ :  $\tilde{x}[1]$ ,  $\tilde{x}[2]$  and  $\tilde{x}[3]$  are 1 if James, Kevin and Kyrie are present respectively, and 0 if they're absent. We represent the result of the game using a Bernoulli random variable  $\tilde{y}$  (win is  $\tilde{y} = 1$  and loss  $\tilde{y} = 0$ ). Using the naive Bayes assumption and empirical estimates of



$p_{\tilde{y}}$  and  $p_{\tilde{x}[i] | \tilde{y}}$

$$\begin{aligned} p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3]}(1, 0, 0, 0) &= p_{\tilde{y}}(1) p_{\tilde{x}[1] | \tilde{y}}(0 | 1) p_{\tilde{x}[2] | \tilde{y}}(0 | 1) p_{\tilde{x}[3] | \tilde{y}}(0 | 1) \\ &= \frac{6}{10} \frac{2}{6} \frac{1}{6} \frac{2}{6} = \frac{1}{90}, \end{aligned} \quad (4.204)$$

$$\begin{aligned} p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3]}(0, 0, 0, 0) &= p_{\tilde{y}}(0) p_{\tilde{x}[1] | \tilde{y}}(0 | 0) p_{\tilde{x}[2] | \tilde{y}}(0 | 0) p_{\tilde{x}[3] | \tilde{y}}(0 | 0) \\ &= \frac{4}{10} \frac{1}{4} \frac{3}{4} \frac{3}{4} = \frac{9}{160}. \end{aligned} \quad (4.205)$$

As a result,

$$p_{\tilde{y} | \tilde{x}[1], \tilde{x}[2], \tilde{x}[3]}(1 | 0, 0, 0) \quad (4.206)$$

$$= \frac{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3]}(1, 0, 0, 0)}{p_{\tilde{x}[1], \tilde{x}[2], \tilde{x}[3]}(0, 0, 0)} \quad (4.207)$$

$$= \frac{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3]}(1, 0, 0, 0)}{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3]}(1, 0, 0, 0) + p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3]}(0, 0, 0, 0)} \quad (4.208)$$

$$= \frac{\frac{1}{90}}{\frac{1}{90} + \frac{9}{160}} \quad (4.209)$$

$$= 0.165. \quad (4.210)$$

#### 4.17 (Spam detector)

- a) The problem is that there are  $2^4 = 16$  different possible values for the entries of  $\tilde{x}$ , and hence 16 different conditional distributions. However, we only have 10 data points.
- b) The email corresponds to  $\tilde{x}[1] = 0$ ,  $\tilde{x}[2] = 0$ ,  $\tilde{x}[3] = 1$ ,  $\tilde{x}[4] = 1$ . By the naive Bayes assumption combined with empirical estimates of  $p_{\tilde{y}}$  and  $p_{\tilde{x}[j] | \tilde{y}}$

$$\begin{aligned} p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1, 0, 0, 1, 1) &= p_{\tilde{y}}(1) p_{\tilde{x}[1] | \tilde{y}}(0 | 1) p_{\tilde{x}[2] | \tilde{y}}(0 | 1) p_{\tilde{x}[3] | \tilde{y}}(1 | 1) p_{\tilde{x}[4] | \tilde{y}}(1 | 1) \\ &= \frac{5}{10} \frac{1}{5} \frac{3}{5} \frac{2}{5} \frac{1}{5}, \end{aligned} \quad (4.211)$$

$$\begin{aligned} p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(0, 0, 0, 1, 1) &= p_{\tilde{y}}(0) p_{\tilde{x}[1] | \tilde{y}}(0 | 0) p_{\tilde{x}[2] | \tilde{y}}(0 | 0) p_{\tilde{x}[3] | \tilde{y}}(1 | 0) p_{\tilde{x}[4] | \tilde{y}}(1 | 0) \\ &= \frac{5}{10} \frac{4}{5} \frac{2}{5} \frac{2}{5} \frac{4}{5}. \end{aligned} \quad (4.212)$$

As a result,

$$p_{\tilde{y} | \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1 | 0, 0, 1, 1) \quad (4.213)$$

$$= \frac{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1, 0, 0, 1, 1)}{p_{\tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(0, 0, 1, 1)} \quad (4.214)$$

$$= \frac{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1, 0, 0, 1, 1)}{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1, 0, 0, 1, 1) + p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(0, 0, 0, 1, 1)} \quad (4.215)$$

$$= \frac{3}{3 + 4 \cdot 2 \cdot 4} \quad (4.216)$$

$$= \frac{3}{35}, \quad (4.217)$$

so  $p_{\tilde{y} | \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(0 | 0, 0, 1, 1) = \frac{32}{35}$  and we classify the email as not spam.

- c) The email corresponds to  $\tilde{x}[1] = 0$ ,  $\tilde{x}[2] = 1$ ,  $\tilde{x}[3] = 1$ ,  $\tilde{x}[4] = 0$ . By the naive Bayes

assumption combined with empirical estimates of  $p_{\tilde{y}}$  and  $p_{\tilde{x}[j] | \tilde{y}}$

$$\begin{aligned} p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1, 0, 1, 1, 0) &= p_{\tilde{y}}(1) p_{\tilde{x}[1] | \tilde{y}}(0 | 1) p_{\tilde{x}[2] | \tilde{y}}(1 | 1) p_{\tilde{x}[3] | \tilde{y}}(1 | 1) p_{\tilde{x}[4] | \tilde{y}}(0 | 1) \\ &= \frac{5}{10} \frac{1}{5} \frac{2}{5} \frac{2}{5} \frac{4}{5}, \end{aligned} \quad (4.218)$$

$$\begin{aligned} p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(0, 0, 1, 1, 0) &= p_{\tilde{y}}(0) p_{\tilde{x}[1] | \tilde{y}}(0 | 0) p_{\tilde{x}[2] | \tilde{y}}(1 | 0) p_{\tilde{x}[3] | \tilde{y}}(1 | 0) p_{\tilde{x}[4] | \tilde{y}}(0 | 0) \\ &= \frac{5}{10} \frac{4}{5} \frac{3}{5} \frac{2}{5} \frac{1}{5}. \end{aligned} \quad (4.219)$$

As a result,

$$p_{\tilde{y} | \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1 | 0, 1, 1, 0) \quad (4.220)$$

$$= \frac{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1, 0, 1, 1, 0)}{p_{\tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(0, 1, 1, 0)} \quad (4.221)$$

$$= \frac{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1, 0, 1, 1, 0)}{p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(1, 0, 1, 1, 0) + p_{\tilde{y}, \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(0, 0, 1, 1, 0)} \quad (4.222)$$

$$= \frac{2 \cdot 2 \cdot 4}{4 \cdot 3 \cdot 2 + 2 \cdot 2 \cdot 4} \quad (4.223)$$

$$= \frac{2}{5}, \quad (4.224)$$

so  $p_{\tilde{y} | \tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]}(0 | 0, 1, 1, 0) = \frac{3}{5}$  and we classify the email as not spam.

In our data the words *alternative* and *medicine* only appear together in spam emails, but the naive Bayes classifier ignores these dependencies.

#### 4.18 (Symbols)

- a) We use empirical probabilities to estimate the transition matrix.

*	†	□	
0	$\frac{1}{6}$	$\frac{4}{11}$	*
$\frac{3}{7}$	$\frac{2}{3}$	$\frac{3}{11}$	†
$\frac{4}{7}$	$\frac{1}{6}$	$\frac{4}{11}$	□

The Markov assumption is not consistent with the data. Under the assumption, the transition  $\dagger \rightarrow \dagger$  is equally likely no matter what symbols occur before, but this is clearly violated in the data: if the previous three symbols are  $\square \dagger \dagger$  or  $* \dagger \dagger$  then the following symbol is always  $\dagger$ , but if the previous three symbols are  $\dagger \dagger \dagger$  then the following symbol is never  $\dagger$ .

- b) The next symbol is either  $\dagger$  or  $\square$  because the transition  $* \rightarrow *$  has probability zero. If the next symbol is  $\dagger$  the most likely sequence is  $\dagger \dagger$ , which has probability

$$P(\text{next 2 symbols} = \dagger \dagger) = P(* \rightarrow \dagger) P(\dagger \rightarrow \dagger) \quad (4.225)$$

$$= \frac{3}{7} \cdot \frac{2}{3} = \frac{2}{7}. \quad (4.226)$$

If the next symbol is  $\square$  the most likely sequence is  $\square \square$  or  $\square *$ , which have the same probability, equal to

$$P(\text{next 2 symbols} = \square \square) = P(* \rightarrow \square) P(\square \rightarrow \square) \quad (4.227)$$

$$= \frac{4}{7} \cdot \frac{4}{11} = \frac{16}{77} < \frac{22}{77} = \frac{2}{7}. \quad (4.228)$$

We conclude that the most likely next two symbols according to the Markov-chain model are  $\dagger\dagger$ .

4.19 (The Markov property) Following the hint, we first establish that

$$p_{\tilde{a}_{i+1} | \tilde{a}_{i-1}, \tilde{a}_i} (a_{i+1} | a_{i-1}, a_i) = \frac{p_{\tilde{a}_{i-1}, \tilde{a}_i, \tilde{a}_{i+1}} (a_{i-1}, a_i, a_{i+1})}{p_{\tilde{a}_{i-1}, \tilde{a}_i} (a_{i-1}, a_i)} \quad (4.229)$$

$$= \frac{\sum_{a_1, \dots, a_{i-2}} p_{\tilde{a}_1, \dots, \tilde{a}_{i+1}} (a_1, \dots, a_{i+1})}{\sum_{a_1, \dots, a_{i-2}} p_{\tilde{a}_1, \dots, \tilde{a}_i} (a_1, \dots, a_i)} \quad (4.230)$$

$$= \frac{\sum_{a_1, \dots, a_{i-2}} p_{\tilde{a}_1} (a_1) \prod_{j=1}^i p_{\tilde{a}_{j+1} | \tilde{a}_j} (a_{j+1} | a_j)}{\sum_{a_1, \dots, a_{i-2}} p_{\tilde{a}_1} (a_1) \prod_{j=1}^{i-1} p_{\tilde{a}_{j+1} | \tilde{a}_j} (a_{j+1} | a_j)} \quad (4.231)$$

$$= \frac{p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i) \sum_{a_1, \dots, a_{i-2}} p_{\tilde{a}_1} (a_1) \prod_{j=1}^{i-1} p_{\tilde{a}_{j+1} | \tilde{a}_j} (a_{j+1} | a_j)}{\sum_{a_1, \dots, a_{i-2}} p_{\tilde{a}_1} (a_1) \prod_{j=1}^{i-1} p_{\tilde{a}_{j+1} | \tilde{a}_j} (a_{j+1} | a_j)} \\ = p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i). \quad (4.232)$$

We then use the result to conclude

$$p_{\tilde{a}_{i+1}, \tilde{a}_{i-1} | \tilde{a}_i} (a_{i+1}, a_{i-1} | a_i) = \frac{p_{\tilde{a}_{i-1}, \tilde{a}_i, \tilde{a}_{i+1}} (a_{i-1}, a_i, a_{i+1})}{p_{\tilde{a}_i} (a_i)} \quad (4.233)$$

$$= \frac{p_{\tilde{a}_{i-1}, \tilde{a}_i} (a_{i-1}, a_i) p_{\tilde{a}_{i+1} | \tilde{a}_{i-1}, \tilde{a}_i} (a_{i+1} | a_{i-1}, a_i)}{p_{\tilde{a}_i} (a_i)} \quad (4.234)$$

$$= \frac{p_{\tilde{a}_i} (a_i) p_{\tilde{a}_{i-1} | \tilde{a}_i} (a_{i-1} | a_i) p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i)}{p_{\tilde{a}_i} (a_i)} \quad (4.235)$$

$$= p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i) p_{\tilde{a}_{i-1} | \tilde{a}_i} (a_{i-1} | a_i). \quad (4.236)$$

4.20 (Employment dynamics)

a) The pmf of  $\tilde{s}$  equals

$$p_{\tilde{s}}(s) = P(\tilde{x}_1 = \text{student}, \dots, \tilde{x}_s = \text{student}, \tilde{x}_{s+1} = \text{not a student}) \quad (4.237)$$

$$= 0.8^s 0.2. \quad (4.238)$$

b)

$$P(\tilde{x}_1 = \text{student} | \tilde{x}_3 = \text{employed}) \quad (4.239)$$

$$= \frac{P(\tilde{x}_1 = \text{student}, \tilde{x}_3 = \text{employed})}{P(\tilde{x}_3 = \text{employed})} \quad (4.240)$$

$$= \frac{\sum_{x_2 \in \{\text{st.}, \text{em.}, \text{un.}\}} P(\tilde{x}_1 = \text{student}, \tilde{x}_2 = x_2, \tilde{x}_3 = \text{employed})}{\sum_{x_1 \in \{\text{st.}, \text{em.}, \text{un.}\}} \sum_{x_2 \in \{\text{st.}, \text{em.}, \text{un.}\}} P(\tilde{x}_1 = x_1, \tilde{x}_2 = x_2, \tilde{x}_3 = \text{employed})} \quad (4.241)$$

$$= \frac{0.8(0.8 \cdot 0.2 + 0.2 \cdot 0.9)}{0.8(0.8 \cdot 0.2 + 0.2 \cdot 0.9) + 0.2(0.9 \cdot 0.9 + 0.1 \cdot 0.4)} \quad (4.242)$$

$$= 0.615. \quad (4.243)$$

c) A stationary distribution should satisfy  $T\pi_* = \pi_*$ , which implies  $0.8\pi_*[1] = \pi_*[1]$ . This is only possible if  $\pi_*[1] = 0$ , which implies

$$\pi_* = \begin{bmatrix} 0 \\ \alpha \\ 1 - \alpha \end{bmatrix}, \quad (4.244)$$

for some  $0 \leq \alpha \leq 1$ . Setting  $T\pi_* = \pi_*$  yields the equations,

$$0.9\alpha + 0.4(1 - \alpha) = \alpha, \quad (4.245)$$

$$0.1\alpha + 0.6(1 - \alpha) = 1 - \alpha. \quad (4.246)$$

Solving the equations, we obtain  $\alpha = 0.8$ .

#### 4.21 (Cellphones)

a) The initial state vector and the transition matrix of the Markov chain are

$$\pi_1 := \begin{bmatrix} 0.9 \\ 0.1 \\ 0 \end{bmatrix}, \quad T = \begin{bmatrix} 0.7 & 0 & 0 \\ 0.2 & 1 & 0 \\ 0.1 & 0 & 1 \end{bmatrix}. \quad (4.247)$$

The transition matrix  $T$  has three eigenvectors

$$q_1 := \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \quad q_2 := \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad q_3 := \begin{bmatrix} 0.80 \\ -0.53 \\ -0.27 \end{bmatrix}. \quad (4.248)$$

The corresponding eigenvalues are  $\lambda_1 := 1$ ,  $\lambda_2 := 1$  and  $\lambda_3 := 0.7$ . We gather the eigenvectors and eigenvalues into two matrices

$$Q := [q_1 \quad q_2 \quad q_3], \quad \Lambda := \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}, \quad (4.249)$$

so that the eigendecomposition of  $T$  is

$$T := Q\Lambda Q^{-1}. \quad (4.250)$$

We compute

$$\alpha := Q^{-1}\pi_1 = \begin{bmatrix} 0.3 \\ 0.7 \\ 1.122 \end{bmatrix}. \quad (4.251)$$

By Theorem 4.39,

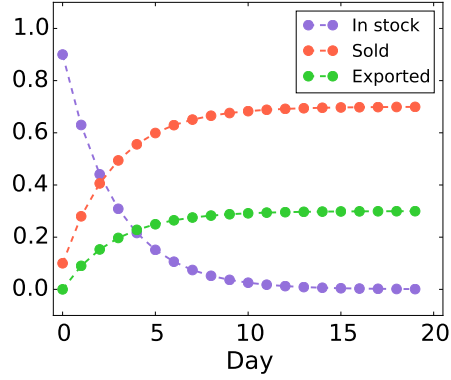
$$\lim_{i \rightarrow \infty} \pi_i = \lim_{i \rightarrow \infty} \sum_{j=1}^3 \lambda^{i-1} \alpha_j q_j \quad (4.252)$$

$$= 0.3 q_1 + 0.7 q_2 \quad (4.253)$$

$$= \begin{bmatrix} 0 \\ 0.7 \\ 0.3 \end{bmatrix}. \quad (4.254)$$

This means that eventually the probability that each phone has been sold locally converges to 0.7 and the probability that it has been exported converges to 0.3.

b) The following graph shows the evolution of the state vector of the simulated Markov chain.



- 4.22 (Empirical state vector and transition matrix) For  $1 \leq j \leq m$ ,  $\pi_X[j]$  is the fraction of data equal to  $s_j$ . Since  $n_j$  is equal to the number of transitions starting at  $s_j$ ,  $\pi_X[j]$  is equal to  $(n_j - 1)/n$  if  $x_n$  is equal to  $s_j$  or to  $n_j/n$  otherwise, so unless  $n_j$  is very small

$$\pi_X[j] \approx \frac{n_j}{n}. \quad (4.255)$$

For any  $l$  between 2 and  $n$ ,

$$\sum_{j=1}^m 1(x_{l-1} = s_j, x_l = s_i) = 1(x_l = s_i), \quad (4.256)$$

because  $x_{l-1}$  is equal to exactly one of the possible values  $s_1, \dots, s_m$ . Consequently,

$$(T_X \pi_X)[i] = \sum_{j=1}^m T_X[i, j] \pi_X[j] \quad (4.257)$$

$$= \sum_{j=1}^m \frac{\pi_X[j]}{n_j} \sum_{l=2}^n 1(x_{l-1} = s_j, x_l = s_i) \quad (4.258)$$

$$\approx \frac{1}{n} \sum_{l=2}^n \sum_{j=1}^m 1(x_{l-1} = s_j, x_l = s_i) \quad (4.259)$$

$$= \frac{1}{n} \sum_{l=2}^n 1(x_l = s_i) \quad (4.260)$$

$$\approx \frac{1}{n} \sum_{l=1}^n 1(x_l = s_i) \quad (4.261)$$

$$= \pi_X[i]. \quad (4.262)$$

---

## Multiple Continuous Variables

### Exercises

- 5.1 (Identities) The only identities that hold are

$$\int_{x=-\infty}^{\infty} f_{\tilde{x}|\tilde{y},\tilde{z}}(x|y,z) \, dx = 1, \quad (5.1)$$

$$\int_{x=-\infty}^{\infty} f_{\tilde{x},\tilde{y}|\tilde{z}}(x,y|z) \, dx = f_{\tilde{y}|\tilde{z}}(y|z). \quad (5.2)$$

The function  $f_{\tilde{x}|\tilde{y},\tilde{z}}(x|y,z)$  is a valid (conditional) probability density of  $\tilde{x}$  for any fixed  $y$  and  $z$  (but not of  $\tilde{y}$  or  $\tilde{z}$ ), so if we integrate it with respect to  $x$ , the integral equals one. Similarly, the function  $f_{\tilde{x},\tilde{y}|\tilde{z}}(x,y|z)$  is a valid (conditional) joint probability density of  $\tilde{x}$  and  $\tilde{y}$  for any fixed  $z$  (but not of  $\tilde{z}$ ), so if we integrate it with respect to  $x$ , we obtain the (conditional) marginal pdf of  $\tilde{y}$ .

- 5.2 (Cross)

- a) The value of the pdf is  $1/12$ , because the area of the cross equals 12. If  $x < -2$  or  $x > 2$ ,  $f_{\tilde{x}}(x) = 0$ . If  $-2 \leq x \leq -1$  or  $1 \leq x \leq 2$ ,

$$f_{\tilde{x}}(x) = \int_{y \in \mathbb{R}} f_{\tilde{x},\tilde{y}}(x,y) \, dy \quad (5.3)$$

$$= \int_{-1}^1 \frac{1}{12} \, dy \quad (5.4)$$

$$= \frac{1}{6}. \quad (5.5)$$

If  $-1 \leq x \leq 1$

$$f_{\tilde{x}}(x) = \int_{y \in \mathbb{R}} f_{\tilde{x},\tilde{y}}(x,y) \, dy \quad (5.6)$$

$$= \int_{-2}^2 \frac{1}{12} \, dy \quad (5.7)$$

$$= \frac{1}{3}. \quad (5.8)$$

- b) If  $x < -2$  or  $x > 2$  the conditional pdf is not defined. If  $-2 \leq x \leq -1$  or  $1 \leq x \leq 2$ , the conditional pdf is uniform between  $-1$  and  $1$ , and equal to

$$f_{\tilde{y}|\tilde{x}}(y|x) = \frac{f_{\tilde{x},\tilde{y}}(x,y)}{f_{\tilde{x}}(x)} \quad (5.9)$$

$$= \frac{\frac{1}{12}}{\frac{1}{6}} = \frac{1}{2}. \quad (5.10)$$

If  $-1 \leq x \leq 1$  the conditional pdf is uniform between -2 and 2, and equal to

$$f_{\tilde{y}|\tilde{x}}(y|x) = \frac{f_{\tilde{x},\tilde{y}}(x,y)}{f_{\tilde{x}}(x)} \quad (5.11)$$

$$= \frac{\frac{1}{12}}{\frac{1}{3}} = \frac{1}{4}. \quad (5.12)$$

Since the conditional distributions are different,  $\tilde{x}$  and  $\tilde{y}$  are not independent.

### 5.3 (Two random variables)

a) The marginal pdf of  $\tilde{b}$  equals

$$f_{\tilde{b}}(b) = \int_{a=0}^2 f_{\tilde{a}}(a) f_{\tilde{b}|\tilde{a}}(b|a) da = \begin{cases} \int_{a=0}^1 \frac{1}{2} da + \int_{a=1}^2 \frac{1}{2} \cdot \frac{1}{2} da = \frac{3}{4}, & \text{if } 0 \leq b \leq 1, \\ \int_{a=1}^2 \frac{1}{2} \cdot \frac{1}{2} da = \frac{1}{4}, & \text{if } 1 \leq b \leq 2, \\ 0, & \text{otherwise.} \end{cases} \quad (5.13)$$

If  $0 \leq b \leq 1$ ,

$$f_{\tilde{a}|\tilde{b}}(a|b) = \frac{f_{\tilde{a}}(a) f_{\tilde{b}|\tilde{a}}(b|a)}{f_{\tilde{b}}(b)} = \begin{cases} \frac{1/2}{3/4} = \frac{2}{3}, & \text{if } 0 \leq a \leq 1, \\ \frac{1/2 \cdot 1/2}{1/4} = \frac{1}{3}, & \text{if } 1 \leq a \leq 2, \\ 0, & \text{otherwise.} \end{cases} \quad (5.14)$$

If  $1 \leq b \leq 2$ ,

$$f_{\tilde{a}|\tilde{b}}(a|b) = \frac{f_{\tilde{a}}(a) f_{\tilde{b}|\tilde{a}}(b|a)}{f_{\tilde{b}}(b)} = \begin{cases} \frac{1/2 \cdot 1/2}{1/4} = 1, & \text{if } 1 \leq a \leq 2, \\ 0, & \text{otherwise.} \end{cases} \quad (5.15)$$

b) We first sample from the marginal pdf of  $\tilde{a}$ . The cdf of  $\tilde{a}$  equals

$$F_{\tilde{a}}(a) = \int_{x=0}^a f_{\tilde{a}}(a) da \quad (5.16)$$

$$= \frac{a}{2}, \quad (5.17)$$

so its inverse is  $F_{\tilde{a}}^{-1}(u) = 2u$ . The sample  $a$  from  $\tilde{a}$  can therefore be obtained by setting:

$$a := 2u_1. \quad (5.18)$$

In our example, that yields  $a = 1.4$ .

Then, we sample from the conditional distribution of  $\tilde{b}$  given  $\tilde{a} = a$ , where  $a$  is the sample from  $\tilde{a}$ . If  $0 \leq 1 \leq a$ , then we just set  $b := u_2$  because the conditional pdf of  $\tilde{b}$  given  $\tilde{a} = a$  is uniform in  $[0, 1]$ . If  $1 \leq a \leq 2$ , as is the case for our example, the conditional pdf of  $\tilde{b}$  given  $\tilde{a} = a$  is uniform in  $[0, 2]$ , just like the marginal distribution of  $\tilde{a}$ , so we set  $b := 2u_2$ , which for our example yields  $b = 0.4$ .

### 5.4 (Piecewise-constant density)

a) If  $0 < a \leq 1$ , the marginal pdf of  $\tilde{a}$  at  $a$  equals

$$f_{\tilde{a}}(a) = \int_{b=0}^2 f_{\tilde{a},\tilde{b}}(a,b) db \quad (5.19)$$

$$= \int_{b=0}^1 \frac{1}{10} db + \int_{b=1}^2 \frac{3}{10} db \quad (5.20)$$

$$= \frac{2}{5}. \quad (5.21)$$

If  $1 < a \leq 2$ ,

$$f_{\tilde{a}}(a) = \int_{b=0}^2 f_{\tilde{a},\tilde{b}}(a,b) \, db \quad (5.22)$$

$$= \int_{b=0}^1 \frac{1}{10} \, db + \int_{b=1}^2 \frac{5}{10} \, db \quad (5.23)$$

$$= \frac{3}{5}. \quad (5.24)$$

For  $a \notin [0, 2]$   $f_{\tilde{a},\tilde{b}}(a,b) = 0$ , so the marginal pdf equals zero.

b) If  $0 < b \leq 1$ , the marginal cdf of  $\tilde{b}$  at  $b$  equals

$$F_{\tilde{b}}(b) = P(\tilde{b} \leq b) \quad (5.25)$$

$$= \int_{a=0}^2 \int_{b=0}^b f_{\tilde{a},\tilde{b}}(a,b) \, da \, db \quad (5.26)$$

$$= \int_{a=0}^1 \int_{b=0}^b \frac{1}{10} \, da \, db + \int_{a=1}^2 \int_{b=0}^b \frac{1}{10} \, da \, db \quad (5.27)$$

$$= \frac{b}{5}. \quad (5.28)$$

If  $1 < b \leq 2$ ,

$$F_{\tilde{b}}(b) = P(\tilde{b} \leq b) \quad (5.29)$$

$$= \int_{a=0}^2 \int_{b=0}^b f_{\tilde{a},\tilde{b}}(a,b) \, da \, db \quad (5.30)$$

$$= \int_{a=0}^1 \int_{b=0}^1 \frac{1}{10} \, da \, db + \int_{a=0}^1 \int_{b=1}^b \frac{1}{10} \, da \, db \quad (5.31)$$

$$+ \int_{a=1}^2 \int_{b=0}^1 \frac{3}{10} \, da \, db + \int_{a=1}^2 \int_{b=1}^b \frac{5}{10} \, da \, db \quad (5.32)$$

$$= \frac{1}{5} + \frac{4(b-1)}{5} = \frac{4b-3}{5}. \quad (5.33)$$

For  $b < 0$   $F_{\tilde{a},\tilde{b}}(a,b) = 0$  and for  $b > 2$   $F_{\tilde{a},\tilde{b}}(a,b) = 1$ .

c) The conditional pdf of  $\tilde{b}$  given  $\tilde{a} = 1.6$  when  $0 < b \leq 1$  equals,

$$f_{\tilde{b}|\tilde{a}}(b|1.6) = \frac{f_{\tilde{a},\tilde{b}}(1.6,b)}{f_{\tilde{a}}(1.6)} \quad (5.34)$$

$$= \frac{\frac{1}{10}}{\frac{3}{5}} = \frac{1}{6}, \quad (5.35)$$

so the desired probability is

$$P(\tilde{b} < 0.5 | \tilde{a} = 1.6) = \int_{b=0}^{0.5} \frac{1}{6} \, db \quad (5.36)$$

$$= \frac{1}{12}. \quad (5.37)$$

d) We first obtain a sample from  $\tilde{b}$  using inverse transform sampling. Since  $u_1 = 0.1$ , this yields  $b = F_{\tilde{b}}^{-1}(u_1) = 0.5$ . Then we obtain a sample from  $\tilde{a}$  applying inverse transform sampling to the conditional cdf of  $\tilde{a}$  given  $\tilde{b} = 0.5$ . Differentiating (5.28)



yields  $f_{\tilde{b}}(0.5) = 1/5$ , so the conditional pdf of  $\tilde{a}$  given  $\tilde{b} = 0.5$  equals

$$f_{\tilde{a}|\tilde{b}}(a|0.5) = \frac{f_{\tilde{a},\tilde{b}}(a,0.5)}{f_{\tilde{b}}(0.5)} \quad (5.38)$$

$$= \frac{1}{2} \quad (5.39)$$

for  $0 < a \leq 2$  (and 0 otherwise). The conditional cdf therefore equals  $a/2$ . We conclude that the second sample is  $a = F_{\tilde{a}}^{-1}(u_2) = 1.4$ .

- 5.5 (Samples and independence) In (1) and (2) the density of the samples in the vertical direction for  $0 \leq a \leq 1$  is completely different to the density for  $1 \leq a \leq 2$ , so the random variables are dependent. In (3) the density of the samples in the vertical direction looks the same for all values of  $a$  between 0 and 1, and is zero otherwise. This suggests that the conditional pdf of  $\tilde{b}$  given  $\tilde{a} = a$  is the same for all values of  $a$ , which would imply that the random variables are independent.

- 5.6 (Independence of continuous random variables)

- a) True. By the assumption we can define a function  $g$  such that  $g(y) := f_{\tilde{y}|\tilde{x}}(y|x)$  for all  $x$  and  $y$ . Marginalizing the joint pdf we obtain

$$f_{\tilde{y}}(y) = \int_{x=-\infty}^{\infty} f_{\tilde{x},\tilde{y}}(x,y) \, dx \quad (5.40)$$

$$= \int_{x=-\infty}^{\infty} f_{\tilde{x}}(x) f_{\tilde{y}|\tilde{x}}(y|x) \, dx \quad (5.41)$$

$$= \int_{x=-\infty}^{\infty} f_{\tilde{x}}(x) g(y) \, dx \quad (5.42)$$

$$= g(y) \int_{x=-\infty}^{\infty} f_{\tilde{x}}(x) \, dx \quad (5.43)$$

$$= g(y) = f_{\tilde{y}|\tilde{x}}(y|x), \quad (5.44)$$

for any  $x$ , so the random variables are independent.

- b) True. At points like  $x = 2$  and  $y = 8$  we have  $f_{\tilde{x},\tilde{y}}(x,y) = 0$  but both  $f_{\tilde{x}}(x)$  and  $f_{\tilde{y}}(y)$  are nonzero (the joint pdf is nonzero on the lines  $x = 2$  and  $y = 8$ , so it will integrate to nonzero values when we marginalize). This implies that  $f_{\tilde{x},\tilde{y}}(x,y) = f_{\tilde{x}}(x)f_{\tilde{y}}(y)$  cannot hold.
- c) False. If the variables are independent then

$$f_{\tilde{x},\tilde{y}}(x_1,y) = f_{\tilde{x}}(x_1)f_{\tilde{y}}(y), \quad (5.45)$$

$$f_{\tilde{x},\tilde{y}}(x_2,y) = f_{\tilde{x}}(x_2)f_{\tilde{y}}(y), \quad (5.46)$$

so  $f_{\tilde{x},\tilde{y}}(x_1,y) \neq f_{\tilde{x},\tilde{y}}(x_2,y)$  can hold for all  $y$  as long as  $f_{\tilde{x}}(x_1) \neq f_{\tilde{x}}(x_2)$ . For a concrete example, take  $\tilde{y}$  uniform between 0 and 1, and  $\tilde{x}$  exponential with parameter  $\lambda$ . If they are independent, for any  $y \in [0,1]$ ,

$$f_{\tilde{x},\tilde{y}}(1,y) = \lambda \exp(-\lambda) \quad (5.47)$$

$$\neq \lambda \exp(-2\lambda) \quad (5.48)$$

$$= f_{\tilde{x},\tilde{y}}(2,y). \quad (5.49)$$

- 5.7 (Regression to the mean)

- a) Let  $q$  denote the  $x$ th percentile of the distribution of  $\tilde{a}$  and  $\tilde{b}$ . By independence between

$\tilde{a}$  and  $\tilde{b}$ , the assumption that they have the same distribution and the definition of percentile,

$$P(\tilde{b} \leq \tilde{a} \mid \tilde{a} = q) = P(\tilde{b} \leq q \mid \tilde{a} = q) \quad (5.50)$$

$$= P(\tilde{b} \leq q) \quad (5.51)$$

$$= x\%. \quad (5.52)$$

In particular if  $x$  is equal to 95, the probability that  $\tilde{b}$  is smaller or equal to  $\tilde{a}$  is 0.95.

- b) The previous answer shows that even if two measurements are completely independent, and are generated by the exact same distribution, the more extreme the first measurement is, the more likely it is for the second to be less extreme.
- c) The value of  $\tilde{a}$  has no effect on the distribution of  $\tilde{b}$ . In fact we use this for our proof! It only affects the probability we are computing, because we are comparing  $\tilde{b}$  to  $\tilde{a}$ .

5.8 (Frog)

- a) From the assumptions we know that the density is zero outside of the ponds, equal to a constant  $c_L$  in the large pond and to another constant  $c_S$  in the small pond. The probability of the frog being in the large pond is  $1/4$  so

$$P(\text{large pond}) = \int_0^{10} \int_0^{10} c_L \, dx \, dy \quad (5.53)$$

$$= 100c_L = \frac{1}{4}, \quad (5.54)$$

which implies  $c_L = \frac{1}{400}$ . Similarly

$$P(\text{small pond}) = \int_{12}^{17} \int_0^5 c_S \, dx \, dy \quad (5.55)$$

$$= 25c_S = \frac{3}{4}, \quad (5.56)$$

which implies  $c_S = \frac{3}{100}$ . Therefore,

$$f_{\tilde{x}, \tilde{y}}(x, y) = \begin{cases} \frac{1}{400} & \text{if } 0 \leq x, y \leq 10, \\ \frac{3}{100} & \text{if } 0 \leq x \leq 5, 12 \leq y \leq 17, \\ 0 & \text{otherwise.} \end{cases} \quad (5.57)$$

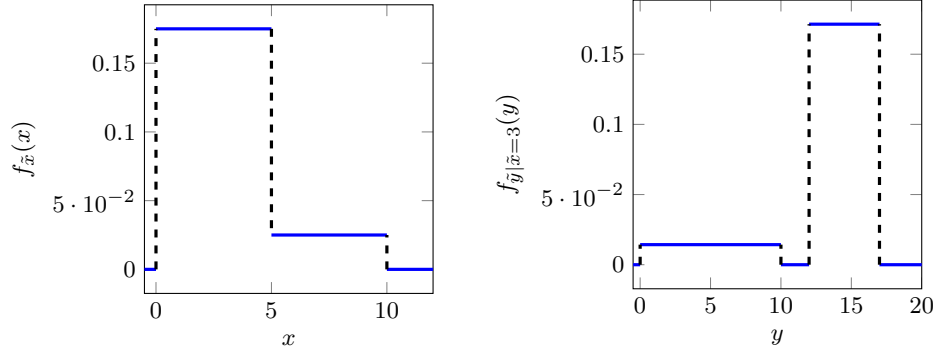
- b) The marginal pdf of the horizontal position of the frog is

$$f_{\tilde{x}}(x) = \int_0^{17} f_{\tilde{x}, \tilde{y}}(x, y) \, dy = \begin{cases} \int_0^{10} \frac{1}{400} \, dy + \int_{12}^{17} \frac{3}{100} \, dy = \frac{7}{40} & \text{if } 0 \leq x \leq 5, \\ \int_0^{10} \frac{1}{400} \, dy = \frac{1}{40} & \text{if } 5 \leq x \leq 10, \\ 0 & \text{otherwise.} \end{cases} \quad (5.58)$$

- c)

$$f_{\tilde{y} \mid \tilde{x}=3}(y) = \frac{f(3, y)}{f_{\tilde{x}}(3)} = \frac{40}{7} \begin{cases} \frac{1}{400} & \text{if } 0 \leq y \leq 10 \\ \frac{3}{100} & \text{if } 12 \leq y \leq 17 \\ 0 & \text{otherwise} \end{cases} \quad (5.59)$$

$$= \begin{cases} \frac{1}{70} & \text{if } 0 \leq y \leq 10 \\ \frac{6}{35} & \text{if } 12 \leq y \leq 17 \\ 0 & \text{otherwise} \end{cases} \quad (5.60)$$



d) The vertical position of the frog is not independent from the horizontal position since

$$f_{\tilde{y}}(13) = \int_0^5 f_{\tilde{x}, \tilde{y}}(x, y) \, dx \quad (5.61)$$

$$= 5 \frac{3}{100} = \frac{3}{20} \neq \frac{6}{35} = f_{\tilde{y}|\tilde{x}=3}(13). \quad (5.62)$$

e) Conditioned on the frog being in the small pond the joint pdf is a constant  $c$  such that

$$P(\text{small pond}) = \int_{12}^{17} \int_0^5 c \, dx \, dy \quad (5.63)$$

$$= 25c = 1, \quad (5.64)$$

so  $c = \frac{1}{25}$ . The joint pdf is consequently

$$f_{\tilde{x}, \tilde{y}|\text{small pond}}(x, y) = \begin{cases} \frac{1}{25} & \text{if } 0 \leq x \leq 5 \text{ and } 12 \leq y \leq 17, \\ 0 & \text{otherwise.} \end{cases} \quad (5.65)$$

We marginalize to find the pdfs of  $\tilde{x}$  and  $\tilde{y}$ .

$$f_{\tilde{x}|\text{small pond}}(x) = \int_{12}^{17} f_{\tilde{x}, \tilde{y}|\text{small pond}}(x, y) \, dy = \begin{cases} \frac{1}{5} & \text{if } 0 \leq x \leq 5, \\ 0 & \text{otherwise.} \end{cases} \quad (5.66)$$

$$f_{\tilde{y}|\text{small pond}}(y) = \int_0^5 f_{\tilde{x}, \tilde{y}|\text{small pond}}(x, y) \, dx = \begin{cases} \frac{1}{5} & \text{if } 12 \leq y \leq 17, \\ 0 & \text{otherwise.} \end{cases} \quad (5.67)$$

The two variables are conditionally independent because  $f_{\tilde{x}|\text{small pond}}(x) f_{\tilde{y}|\text{small pond}}(y) = f_{\tilde{x}, \tilde{y}|\text{small pond}}(x, y)$  for any values of  $x$  and  $y$ .

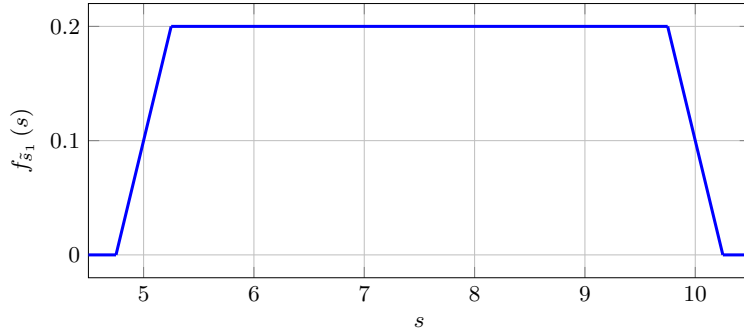
## 5.9 (Sonar)

a) Note that for fixed  $x$   $f_{\tilde{s}_1|\tilde{x}}(s|x) = 0$  if  $x < s - 0.25$  or  $x > s + 0.25$ . This implies

$$f_{\tilde{s}_1}(s) = \int_5^{10} f_{\tilde{x}}(x) f_{\tilde{s}_1|\tilde{x}}(s|x) dx \quad (5.68)$$

$$= \int_{\max\{5, s-0.25\}}^{\min\{10, s+0.25\}} \frac{2}{5} dx \quad (5.69)$$

$$= \begin{cases} 0 & \text{if } s < 4.75 \\ \frac{2(s-4.75)}{5} & \text{if } 4.75 \leq s \leq 5.25 \\ \frac{1}{5} & \text{if } 5.25 \leq s \leq 9.75 \\ \frac{2(10.25-s)}{5} & \text{if } 9.75 \leq s \leq 10.25 \\ 0 & \text{if } s > 10.25. \end{cases} \quad (5.70)$$



b) For fixed  $s_1$   $f_{\tilde{s}_1|\tilde{x}}(s_1|x) = 0$  if  $x < s_1 - 0.25$  or  $x > s_1 + 0.25$ . For fixed  $s_2$   $f_{\tilde{s}_2|\tilde{x}}(s_2|x) = 0$  if  $x < s_2 - 0.25$  or  $x > s_2 + 0.25$ , so

$$f_{\tilde{s}_1, \tilde{s}_2}(7, 7.1) = \int_5^{10} f_{\tilde{x}}(x) f_{\tilde{s}_1|\tilde{x}}(7|x) f_{\tilde{s}_2|\tilde{x}}(7.1|x) dx \quad (5.71)$$

$$= \int_{6.85}^{7.25} \frac{4}{5} dx \quad (5.72)$$

$$= 0.32. \quad (5.73)$$

$$f_{\tilde{x}|\tilde{s}_1, \tilde{s}_2}(x|7, 7.1) = \frac{f_{\tilde{x}, \tilde{s}_1, \tilde{s}_2}(x, 7, 7.1)}{f_{\tilde{s}_1, \tilde{s}_2}(7, 7.1)} \quad (5.74)$$

$$= \begin{cases} \frac{4/5}{0.32} = 2.5 & \text{if } 6.85 \leq x \leq 7.25, \\ 0 & \text{otherwise.} \end{cases} \quad (5.75)$$

c) For fixed  $s_1$   $f_{\tilde{s}_1|\tilde{x}}(s_1|x) = 0$  if  $x < s_1 - 0.25$  or  $x > s_1 + 0.25$ . For fixed  $s_2$   $f_{\tilde{s}_2|\tilde{x}}(s_2|x) = 0$  if  $x < s_2 - 0.25$  or  $x > s_2 + 0.25$ , so

$$f_{\tilde{s}_1, \tilde{s}_2}(s_1, s_2) = \int_5^{10} f_{\tilde{x}}(x) f_{\tilde{s}_1|\tilde{x}}(s_1|x) f_{\tilde{s}_2|\tilde{x}}(s_2|x) dx \quad (5.76)$$

$$= \int_{\max\{5, s_1-0.25, s_2-0.25\}}^{\min\{10, s_1+0.25, s_2+0.25\}} \frac{4}{5} dx \quad (5.77)$$

$$= 0.8 (\min\{10, s_1 + 0.25, s_2 + 0.25\} - \max\{5, s_1 - 0.25, s_2 - 0.25\})$$

if  $4.75 \leq s_1, s_2 \leq 10.25$  and  $|s_2 - s_1| \leq 0.5$ , and 0 otherwise.

The two measurements are not independent. For example  $f_{\tilde{s}_1}(5) \neq 0$  and  $f_{\tilde{s}_2}(10) \neq 0$  but  $f_{\tilde{s}_1, \tilde{s}_2}(5, 10) = 0$  so the joint pdf is not the product of the marginals. This makes sense since  $\tilde{s}_1$  provides information about  $\tilde{x}$ , which in turn provides information about  $\tilde{s}_2$ .

#### 5.10 (Rufus)

- a) By assumption, the joint probability density is equal to a constant  $c$ . It must integrate to one, so

$$\int_{\text{garden}} c \, dx \, dy = c \, \text{Area}(\text{garden}) = 1. \quad (5.78)$$

The area of the garden is equal to  $100^2 - 40^2 = 10000 - 1600 = 8400$ , which implies

$$f_{\tilde{x}, \tilde{y}}(x, y) = \begin{cases} \frac{1}{8400} & \text{if } \{(x, y) : -50 \leq x \leq 50, -50 \leq y \leq 50\} \setminus \{(x, y) : -20 \leq x \leq 20, -20 \leq y \leq 20\}, \\ 0 & \text{otherwise.} \end{cases} \quad (5.79)$$

- b) For  $y$  in  $[-50, -20]$  or  $[20, 50]$

$$f_{\tilde{y}}(y) = \int_{x=-\infty}^{\infty} f_{\tilde{x}, \tilde{y}}(x, y) \, dx \quad (5.80)$$

$$= \int_{x=-50}^{50} \frac{1}{8400}(x, y) \, dx \quad (5.81)$$

$$= \frac{50 + 50}{8400} = \frac{1}{84}. \quad (5.82)$$

For  $y$  in  $[-20, 20]$

$$f_{\tilde{y}}(y) = \int_{x=-\infty}^{\infty} f_{\tilde{x}, \tilde{y}}(x, y) \, dx \quad (5.83)$$

$$= \int_{x=-50}^{-20} \frac{1}{8400}(x, y) \, dx + \int_{x=20}^{50} \frac{1}{8400}(x, y) \, dx \quad (5.84)$$

$$= 2 \cdot \frac{50 - 20}{8400} = \frac{1}{140}. \quad (5.85)$$

For other values of  $y$ ,  $f_{\tilde{y}}(y) = 0$ . The pdf is shown in Figure 5.1.

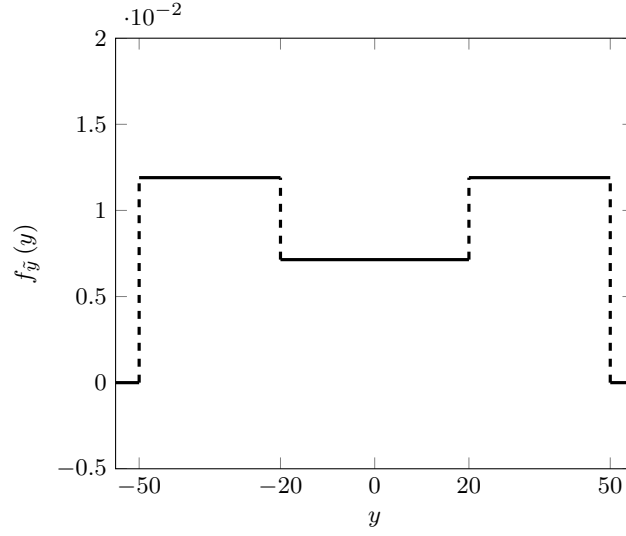
- c) In order to compute the conditional pdf of  $\tilde{x}$  given  $\tilde{y}$  we need to compute the marginal pdf of  $\tilde{y}$  by marginalizing over  $\tilde{x}$ . The conditional density is only defined for values of  $y$  such that  $f_{\tilde{y}}(y) \neq 0$ . For  $y$  in  $[-50, -20]$  or  $[20, 50]$

$$f_{\tilde{x}|\tilde{y}}(x|y) = \frac{f_{\tilde{x}, \tilde{y}}(x, y)}{f_{\tilde{y}}(y)} = \begin{cases} \frac{1}{100} & \text{if } -50 \leq x \leq 50, \\ 0 & \text{otherwise.} \end{cases} \quad (5.86)$$

For  $y$  in  $[-20, 20]$

$$f_{\tilde{x}|\tilde{y}}(x|y) = \begin{cases} \frac{1}{60} & \text{if } -50 \leq x \leq -20 \text{ or } 20 \leq x \leq 50, \\ 0 & \text{otherwise.} \end{cases} \quad (5.87)$$

The conditional pdfs are plotted in Figure 5.2.



**Figure 5.1** Pdf of the vertical position  $\tilde{y}$  of Rufus.

- d)  $\tilde{x}$  and  $\tilde{y}$  are not independent. By (5.85)  $f_{\tilde{y}}(0) = \frac{60}{8400} \neq 0$  and by symmetry  $f_{\tilde{x}}(0) = \frac{60}{8400} \neq 0$ . However  $(0,0)$  is not in the garden, so

$$f_{\tilde{x},\tilde{y}}(0,0) = 0 \neq f_{\tilde{x}}(0) f_{\tilde{y}}(0). \quad (5.88)$$

- e) The cdf of  $\tilde{y}$  equals

$$F_{\tilde{y}}(y) = \begin{cases} 0 & \text{for } y < -50, \\ \frac{y - (-50)}{84} = \frac{y+50}{84} & \text{for } -50 \leq y < -20, \\ \frac{30}{84} + \frac{y - (-20)}{140} = \frac{5}{14} + \frac{y+20}{140} & \text{for } -20 \leq y < 20, \\ \frac{30}{84} + \frac{40}{140} + \frac{y-20}{84} = \frac{9}{14} + \frac{y-20}{84} & \text{for } 20 \leq y < 50, \\ 1 & \text{for } y \geq 50. \end{cases} \quad (5.89)$$

We obtain a sample of  $\tilde{y}$  by applying the inverse transform method. Since  $0.1 < 5/14$ , we set

$$F_{\tilde{y}}(y) = \frac{y+50}{84} = 0.1, \quad (5.90)$$

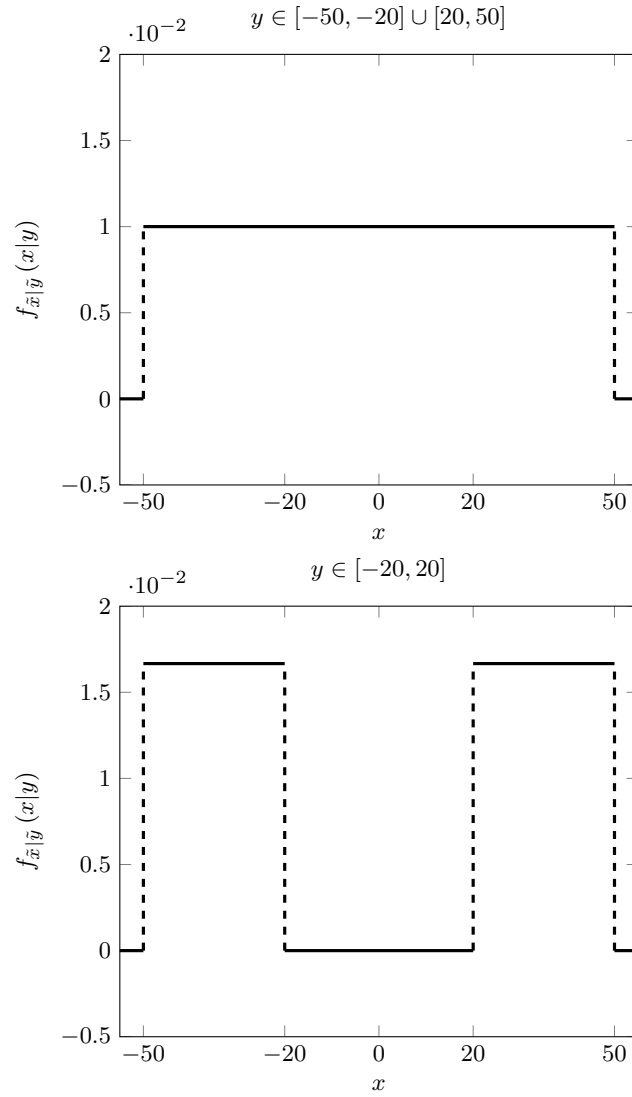
which yields  $y = 8.4 - 50 = -41.6$ . Now we consider the conditional cdf of  $\tilde{x}$  given  $\tilde{y} = -41.6$ , which is equal to

$$F_{\tilde{x}|\tilde{y}}(x | -41.6) = \begin{cases} 0 & \text{for } x < -50, \\ \frac{x+50}{100} & \text{for } -50 \leq x \leq 50, \\ 1 & \text{for } x > 50. \end{cases} \quad (5.91)$$

To obtain the sample from  $x$  we set  $F_{\tilde{x}|\tilde{y}}(x | -41.6) = 0.3$ , which yields  $x = -20$ .

#### 5.11 (Spider on a wall)

- a) Let  $\tilde{x}$  be the horizontal position and  $\tilde{y}$  be the vertical position, or height. We have the



**Figure 5.2** Conditional pdf of the position  $\tilde{x}$  of Rufus given his  $\tilde{y}$  position.

following two equations:

$$P(\text{behind the painting}) = 2P(\text{not behind the painting}), \quad (5.92)$$

$$P(\text{behind the painting}) + P(\text{not behind the painting}) = 1. \quad (5.93)$$

Under the assumptions, the probability density is constant in the region behind the painting and also in the rest of the wall. We denote these two densities by  $c_1$  and  $c_2$  respectively. Since the painting has area 4 and the rest of the wall has area 96, the

equations imply

$$4c_1 = 2 \cdot 96c_2, \quad (5.94)$$

$$4c_1 + 96c_2 = 1. \quad (5.95)$$

This implies  $c_1 = 1/6$  and  $c_2 = 1/288$ , so

$$f_{(\tilde{x}, \tilde{y})}(x, y) = \begin{cases} 1/6 & \text{if } (x, y) \in \text{painting} \\ 1/288 & \text{if } (x, y) \notin \text{painting} \end{cases} \quad (5.96)$$

b) The marginal distribution of  $\tilde{y}$  is supported on  $(0, 10)$ :

$$f_{\tilde{y}}(y) = \int f_{(\tilde{x}, \tilde{y})}(x, y) dx = \begin{cases} \int_0^{10} \frac{1}{288} dx = 5/144 & \text{if } y \notin (6, 8) , \\ \int_{(0,4) \cup (6,10)} \frac{1}{288} dx + \int_{(4,6)} 1/6 dx = 13/36 & \text{if } y \in (6, 8) , \\ 0 & \text{otherwise.} \end{cases}$$

c) There are possible three cases, the spider is either below the paint, above the paint or at the paint level:

$$F_{\tilde{y}|(\tilde{x}, \tilde{y}) \notin (4,6) \times (6,8)}(y) = P(\tilde{y} < y | (\tilde{x}, \tilde{y}) \notin (4,6) \times (6,8)) \quad (5.97)$$

$$= \frac{P(\tilde{y} < y, (\tilde{x}, \tilde{y}) \notin (4,6) \times (6,8))}{P((\tilde{x}, \tilde{y}) \notin (4,6) \times (6,8))} \quad (5.98)$$

$$= \begin{cases} 30y/288 & \text{if } y \in (0, 6) , \\ 180/288 + 24(y-6)/288 & \text{if } y \in (6, 8) , \\ 228/288 + 30(y-8)/288 & \text{if } y \in (8, 10) . \end{cases} \quad (5.99)$$

We have used the following intermediate results.

$$P((\tilde{x}, \tilde{y}) \notin (4,6) \times (6,8)) = 1 - P((\tilde{x}, \tilde{y}) \in (4,6) \times (6,8)) \quad (5.100)$$

$$= 1 - \int_{x=4}^6 \int_{y=6}^8 f_{(\tilde{x}, \tilde{y})}(x, y) dx dy \quad (5.101)$$

$$= 1 - \int_{x=4}^6 \int_{y=6}^8 \frac{1}{6} dx dy \quad (5.102)$$

$$= \frac{1}{3}. \quad (5.103)$$

If  $y \in (0, 6)$

$$P(\tilde{y} < y, (\tilde{x}, \tilde{y}) \notin (4,6) \times (6,8)) = P(\tilde{y} < y) \quad (5.104)$$

$$= \int_{u=0}^y f_{\tilde{y}}(u) du \quad (5.105)$$

$$= \frac{5y}{144}. \quad (5.106)$$

If  $y \in (6, 8)$

$$P(\tilde{y} < y, (\tilde{x}, \tilde{y}) \notin (4,6) \times (6,8)) = P(\tilde{y} < 6) + P(6 \leq \tilde{y} < y, (\tilde{x}, \tilde{y}) \notin (4,6) \times (6,8))$$

$$= \frac{30}{144} + \int_{x=0}^4 \int_{u=6}^y \frac{1}{288} dx dy + \int_{x=6}^{10} \int_{u=6}^y \frac{1}{288} dx du$$

$$= \frac{30}{144} + \frac{8(y-6)}{288}. \quad (5.107)$$



If  $y \in (8, 10)$

$$\begin{aligned} P(\tilde{y} < y, (\tilde{x}, \tilde{y}) \notin (4, 6) \times (6, 8)) &= P(\tilde{y} < 6) + P(6 \leq \tilde{y} < 8, (\tilde{x}, \tilde{y}) \notin (4, 6) \times (6, 8)) \\ &\quad + P(8 \leq \tilde{y} < y) \end{aligned} \quad (5.108)$$

$$= \frac{30}{144} + \frac{16}{288} + \int_{u=8}^y f_{\tilde{y}}(u) \, du \quad (5.109)$$

$$= \frac{76}{288} + \frac{5(y-8)}{144}. \quad (5.110)$$

- 5.12 (Simulating a constant density) There are different ways of generating the samples. Arguably the easiest is to generate two independent uniform samples  $u_1, u_2$  and only use them if the vector  $\begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$  is in the region of interest, which we denote by  $\mathcal{R}$ . Let  $\tilde{u}_1$  and  $\tilde{u}_2$  denote two independent samples from a uniform distribution in  $[0, 1]$ . The cdf of the resulting samples is

$$F_{\tilde{u}_1, \tilde{u}_2 | \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R}}(u_1, u_2) = P(\tilde{u}_1 \leq u_1, \tilde{u}_2 \leq u_2 | \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R}) \quad (5.111)$$

$$= \frac{P(\tilde{u}_1 \leq u_1, \tilde{u}_2 \leq u_2, \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R})}{P(\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R})} \quad (5.112)$$

$$= \frac{F_{\tilde{u}_1, \tilde{u}_2}(u_1, u_2)}{P(\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R})} \quad (5.113)$$

as long as  $\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R}$ . Differentiating, we obtain

$$f_{\tilde{u}_1, \tilde{u}_2 | \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R}}(u_1, u_2) = \frac{f_{\tilde{u}_1, \tilde{u}_2}(u_1, u_2)}{P(\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R})} \quad (5.114)$$

$$= \frac{f_{\tilde{u}_1}(u_1)f_{\tilde{u}_2}(u_2)}{P(\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R})} \quad (5.115)$$

$$= \frac{1}{P(\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathcal{R})}. \quad (5.116)$$

The resulting samples therefore have constant density in the region of interest.

- 5.13 (Multivariate Gaussian pdf) We denote the multivariate Gaussian pdf by  $f_{\tilde{x}}$  and its mean vector by  $\mu \in \mathbb{R}^d$ .

- a) Let  $\sigma_i^2$  be the  $i$ th diagonal entry of the covariance matrix  $\Sigma$  for  $1 \leq i \leq d$ . If the covariance matrix is diagonal, the joint pdf is equal to the product of the marginal pdfs,

$$f_{\tilde{x}}(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \quad (5.117)$$

$$= \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left(-\frac{(x[i] - \mu_i)^2}{2\sigma_i^2}\right). \quad (5.118)$$

Consequently, by Exercise 3.9

$$\int_{x \in \mathbb{R}^d} f_{\tilde{x}}(x) \, dx = \prod_{i=1}^d \int_{x[i] = -\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left(-\frac{(x[i] - \mu_i)^2}{2\sigma_i^2}\right) \, dx[i] \quad (5.119)$$

$$= \prod_{i=1}^d 1 = 1. \quad (5.120)$$

- b) Using a non-diagonal covariance matrix just rotates the joint pdf with respect to the standard axes. Consequently, its integral should not change and must therefore equal one.

More formally, by the spectral theorem,  $\Sigma = U\Lambda U^T$  for an orthogonal  $d \times d$  matrix  $U$  and a diagonal  $d \times d$  matrix  $\Lambda$ . If we perform the change of variable  $z = U^T(x - \mu)$ , then  $dz = dx$ , since  $|U| = 1$ , because  $U$  is orthogonal. Consequently,

$$\int_{x \in \mathbb{R}^d} f_{\tilde{x}}(x) dx = \int_{x \in \mathbb{R}^d} \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) dx \quad (5.121)$$

$$\begin{aligned} &= \int_{x \in \mathbb{R}^d} \frac{1}{\sqrt{(2\pi)^d |U| |\Lambda| |U|}} \exp\left(-\frac{1}{2}(x - \mu)^T U \Lambda^{-1} U^T (x - \mu)\right) dx \\ &= \int_{z \in \mathbb{R}^d} \frac{1}{\sqrt{(2\pi)^d |\Lambda|}} \exp\left(-\frac{1}{2}z^T \Lambda^{-1} z\right) dz, \end{aligned} \quad (5.122)$$

which equals one because  $\Lambda$  is diagonal, as established above.

#### 5.14 (Exotic fruit)

- a) The problem is that it runs into the curse of dimensionality. We need to condition on 4 features, but we only have 10 data points. Consequently it will be very difficult to find a relevant subset of the data, as the vast majority of possible combinations of feature values are unobserved.
- b) The maximum-likelihood estimate of the mean parameter is

$$\mu_{\text{ML}} = \frac{1}{10} \sum_{i=1}^n x_i = \begin{bmatrix} 1.53 \\ 16.5 \\ 10.6 \\ 98.6 \\ 15.4 \end{bmatrix}. \quad (5.123)$$

The maximum-likelihood estimate of the covariance-matrix parameter is

$$\Sigma_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{\text{ML}})(x_i - \mu_{\text{ML}})^T \quad (5.124)$$

$$= \begin{bmatrix} 0.51 & -1.61 & 1.45 & -7.89 & 3.97 \\ -1.61 & 7.05 & -3.9 & 0.6 & -12.25 \\ 1.45 & -3.9 & 9.44 & -52.46 & 12.62 \\ -7.89 & 0.6 & -52.46 & 616.84 & -90.18 \\ 3.97 & -12.25 & 12.62 & -90.18 & 39.01 \end{bmatrix}. \quad (5.125)$$

We represent the parameters as follows, to separate the fruit weight and the rest of the features,

$$\mu_{\text{ML}} = \begin{bmatrix} \mu_{\text{fruit}} \\ \mu_{\text{rest}} \end{bmatrix}, \quad (5.126)$$

$$\Sigma_{\text{ML}} = \begin{bmatrix} \Sigma_{\text{fruit}} & \Sigma_{\text{fruit,rest}}^T \\ \Sigma_{\text{fruit,rest}} & \Sigma_{\text{rest}} \end{bmatrix}, \quad (5.127)$$

where

$$\mu_{\text{fruit}} = 1.53, \quad \Sigma_{\text{fruit}} = 0.51, \quad (5.128)$$

$$\mu_{\text{rest}} = \begin{bmatrix} 1.53 \\ 16.5 \\ 10.6 \\ 98.6 \\ 15.4 \end{bmatrix}, \quad \Sigma_{\text{rest}} = \begin{bmatrix} 7.05 & -3.9 & 0.6 & -12.25 \\ -3.9 & 9.44 & -52.46 & 12.62 \\ 0.6 & -52.46 & 616.84 & -90.18 \\ -12.25 & 12.62 & -90.18 & 39.01 \end{bmatrix}, \quad (5.129)$$

$$\Sigma_{\text{fruit,rest}} = \begin{bmatrix} -1.61 \\ 1.45 \\ -7.89 \\ 3.97 \end{bmatrix}. \quad (5.130)$$

The conditional pdf of the fruit weight given the rest of the features is Gaussian with mean

$$\mu_{\text{cond}} = \mu_{\text{fruit}} + \Sigma_{\text{fruit,rest}}^T \Sigma_{\text{rest}}^{-1} \left( \begin{bmatrix} 15 \\ 20 \\ 120 \\ 8 \end{bmatrix} - \mu_{\text{rest}} \right) \quad (5.131)$$

$$= 1.77, \quad (5.132)$$

and variance

$$\sigma_{\text{cond}}^2 = \Sigma_{\text{cond}} \quad (5.133)$$

$$= \Sigma_{\text{fruit}} - \Sigma_{\text{fruit,rest}}^T \Sigma_{\text{rest}}^{-1} \Sigma_{\text{fruit,rest}} \quad (5.134)$$

$$= 0.037. \quad (5.135)$$

### 5.15 (Gaussian Bayesian model)

a) By the chain rule

$$f_{\tilde{\mu}, \tilde{y}}(\mu, y) = f_{\tilde{\mu}}(\mu) f_{\tilde{y} | \tilde{\mu}}(y | \mu) \quad (5.136)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right) \quad (5.137)$$

$$= \frac{1}{2\pi\sigma} \exp\left(-\frac{1}{2} \left( \frac{\mu^2}{\sigma^2} + \mu^2 - 2\mu y + y^2 \right)\right) \quad (5.138)$$

$$= \frac{1}{2\pi\sigma} \exp\left(-\frac{1}{2} \begin{bmatrix} \mu \\ y \end{bmatrix}^T \Sigma^{-1} \begin{bmatrix} \mu \\ y \end{bmatrix}\right), \quad (5.139)$$

where

$$\Sigma^{-1} = \begin{bmatrix} 1 + \frac{1}{\sigma^2} & -1 \\ -1 & 1 \end{bmatrix}. \quad (5.140)$$

This is the joint pdf of a Gaussian random vector with mean zero and covariance-matrix parameter equal to

$$\Sigma = \frac{1}{1 + \frac{1}{\sigma^2} - 1} \begin{bmatrix} 1 & 1 \\ 1 & 1 + \frac{1}{\sigma^2} \end{bmatrix} \quad (5.141)$$

$$= \begin{bmatrix} \sigma^2 & \sigma^2 \\ \sigma^2 & 1 + \sigma^2 \end{bmatrix}. \quad (5.142)$$

You can check that the determinant of  $\Sigma$  is equal to  $\sigma$ .

- b) We apply Theorem 5.24 in the notes on multiple continuous random variables, setting  $\tilde{a} := \tilde{y}$ ,  $\tilde{b} := \tilde{\mu}$ . We have  $\sigma_{\tilde{a}}^2 = 1 + \sigma^2$ ,  $\sigma_{\tilde{b}}^2 = \sigma^2$  and  $\rho = \frac{\sigma^2}{\sigma_{\tilde{a}}\sigma_{\tilde{b}}} = \frac{\sigma}{\sqrt{1+\sigma^2}}$ . By the theorem, conditioned on  $\tilde{y} = y$ ,  $\mu$  is Gaussian with mean

$$\mu_{\text{cond}} = \frac{\rho\sigma_{\tilde{b}}y}{\sigma_{\tilde{a}}} \quad (5.143)$$

$$= \frac{\sigma^2 y}{1 + \sigma^2} \quad (5.144)$$

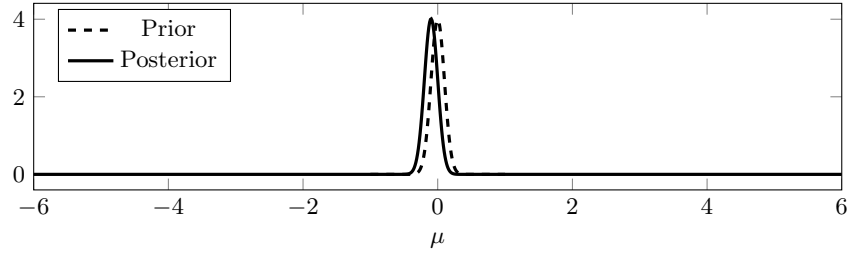
and variance

$$\sigma_{\text{cond}}^2 = (1 - \rho^2)\sigma_{\tilde{b}}^2 \quad (5.145)$$

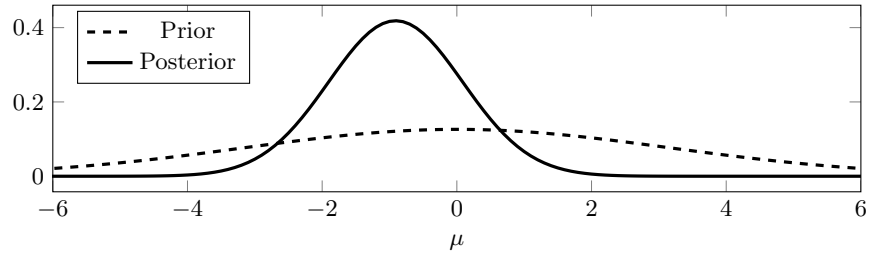
$$= \left(1 - \frac{\sigma^2}{1 + \sigma^2}\right)\sigma^2 \quad (5.146)$$

$$= \frac{\sigma^2}{1 + \sigma^2}. \quad (5.147)$$

- c) When  $y = -1$  and  $\sigma^2 := 0.01$ ,  $\mu_{\text{cond}} = -0.0099$  and  $\sigma_{\text{cond}}^2 = 0.0099$ . The scientist is very confident in the prior. The posterior and prior pdfs are very similar.



- d) When  $y = -1$  and  $\sigma^2 := 10$ ,  $\mu_{\text{cond}} = -0.909$  and  $\sigma_{\text{cond}}^2 = 0.909$ . The scientist is not confident in the prior. The posterior and prior pdfs are very different.



## Discrete and Continuous Variables

### Exercises

#### 6.1 (Shared car)

- a) Let  $\tilde{a}$  be a random variable with range  $\{1, 2, 3\}$ . If  $\tilde{a} = 1$  Carlos drives. If  $\tilde{a} = 2$  Dani drives. If  $\tilde{a} = 3$  Felix drives. We have

$$f_{\tilde{d}}(d) = \sum_{a=1}^3 p_{\tilde{a}}(a) f_{\tilde{d}|\tilde{a}}(d|a) = \begin{cases} 0 & \text{if } d < 0, \\ \frac{1}{3} \left( \frac{1}{10} + \frac{1}{20} + \frac{1}{30} \right) = \frac{11}{180} & \text{if } 0 \leq d \leq 10, \\ \frac{1}{3} \left( \frac{1}{20} + \frac{1}{30} \right) = \frac{5}{180} & \text{if } 10 \leq d \leq 20, \\ \frac{1}{90} & \text{if } 20 \leq d \leq 30, \\ 0 & \text{if } d > 30. \end{cases} \quad (6.1)$$

- b) By the chain rule of discrete and continuous random variables,

$$p_{\tilde{a}|\tilde{d}}(2|15) = \frac{p_{\tilde{a}}(2) f_{\tilde{d}|\tilde{a}}(15|2)}{f_{\tilde{d}}(15)} \quad (6.2)$$

$$= \frac{\frac{1}{3} \frac{1}{20}}{\frac{5}{180}} \quad (6.3)$$

$$= 0.6. \quad (6.4)$$

#### 6.2 (Buckets)

- a) It is easier to compute the probability that  $\tilde{s}_1 \geq 1$ , as this cannot happen if the bucket is small. Since  $f_{\tilde{s}_1|\tilde{b}}(s|b) = 1/2$  between 0 and 2,

$$\mathrm{P}(\tilde{s}_1 \geq 1) = \mathrm{P}(\tilde{s}_1 \geq 1 | \tilde{b} = 0) \mathrm{P}(\tilde{b} = 0) + \mathrm{P}(\tilde{s}_1 \geq 1 | \tilde{b} = 1) \mathrm{P}(\tilde{b} = 1) \quad (6.5)$$

$$= \mathrm{P}(\tilde{s}_1 \geq 1 | \tilde{b} = 1) \mathrm{P}(\tilde{b} = 1) \quad (6.6)$$

$$= \frac{1}{4} \int_{s=1}^2 \frac{1}{2} ds \quad (6.7)$$

$$= \frac{1}{8}. \quad (6.8)$$

We conclude that

$$\mathrm{P}(\tilde{s}_1 < 1) = 1 - \mathrm{P}(\tilde{s}_1 \geq 1) \quad (6.9)$$

$$= \frac{7}{8}. \quad (6.10)$$

b) Under the conditional independence assumption,

$$P(\tilde{b} = 1 \mid \tilde{s}_1 = 1/2, \tilde{s}_2 = 3/4) \quad (6.11)$$

$$= \frac{p_{\tilde{b}}(1)f_{\tilde{s}_1|\tilde{b}}(1/2|1)f_{\tilde{s}_2|\tilde{b}}(3/4|1)}{f_{\tilde{s}_1, \tilde{s}_2}(1/2, 3/4)} \quad (6.12)$$

$$= \frac{p_{\tilde{b}}(1)f_{\tilde{s}_1|\tilde{b}}(1/2|1)f_{\tilde{s}_2|\tilde{b}}(3/4|1)}{p_{\tilde{b}}(0)f_{\tilde{s}_1|\tilde{b}}(1/2|0)f_{\tilde{s}_2|\tilde{b}}(3/4|0) + p_{\tilde{b}}(1)f_{\tilde{s}_1|\tilde{b}}(1/2|1)f_{\tilde{s}_2|\tilde{b}}(3/4|1)} \quad (6.13)$$

$$= \frac{\frac{1}{4} \cdot \frac{1}{2} \cdot \frac{1}{2}}{\frac{3}{4} \cdot 1 \cdot 1 + \frac{1}{4} \cdot \frac{1}{2} \cdot \frac{1}{2}} \quad (6.14)$$

$$= \frac{1}{13}. \quad (6.15)$$

c) Under the assumptions  $\tilde{b}$  and  $\tilde{s}_2$  are conditionally independent given  $\tilde{s}_1$ , so

$$P(\tilde{b} = 1 \mid \tilde{s}_1 = 1/2, \tilde{s}_2 = 3/4) = P(\tilde{b} = 1 \mid \tilde{s}_1 = 1/2) \quad (6.16)$$

$$= \frac{p_{\tilde{b}}(1)f_{\tilde{s}_1|\tilde{b}}(1/2|1)}{f_{\tilde{s}_1}(1/2)} \quad (6.17)$$

$$= \frac{p_{\tilde{b}}(1)f_{\tilde{s}_1|\tilde{b}}(1/2|1)}{p_{\tilde{b}}(0)f_{\tilde{s}_1|\tilde{b}}(1/2|0) + p_{\tilde{b}}(1)f_{\tilde{s}_1|\tilde{b}}(1/2|1)} \quad (6.18)$$

$$= \frac{\frac{1}{4} \cdot \frac{1}{2}}{\frac{3}{4} \cdot 1 + \frac{1}{4} \cdot \frac{1}{2}} \quad (6.19)$$

$$= \frac{1}{7}. \quad (6.20)$$

### 6.3 (Balls)

a) Let  $\tilde{d}$  be a random variable representing the distance and  $\tilde{h}$  a Bernoulli random variable indicating whether the ball is hollow ( $\tilde{h} = 1$ ) or not ( $\tilde{h} = 0$ ). The conditional pdf  $f_{\tilde{d}|\tilde{h}}(d|1)$  of  $\tilde{d}$  given  $\tilde{h} = 1$  equals  $1/5$  if  $d$  is in  $[5, 10]$  and zero otherwise. The conditional pdf  $f_{\tilde{d}|\tilde{h}}(d|0)$  of  $\tilde{d}$  given  $\tilde{h} = 0$  equals  $1/5$  if  $d$  is in  $[1, 6]$  and zero otherwise. Consequently, the marginal pdf of  $\tilde{d}$  equals

$$f_{\tilde{d}}(d) = p_{\tilde{h}}(0)f_{\tilde{d}|\tilde{h}}(d|0) + p_{\tilde{h}}(1)f_{\tilde{d}|\tilde{h}}(d|1) \quad (6.21)$$

$$= \begin{cases} \frac{1}{20} & \text{if } 1 \leq d < 5, \\ \frac{1}{5} & \text{if } 5 \leq d < 6, \\ \frac{3}{20} & \text{if } 6 \leq d < 10, \\ 0 & \text{otherwise.} \end{cases} \quad (6.22)$$

b) Let  $\tilde{d}_1$  be the first distance and  $\tilde{d}_2$  the second. Notice that the conditional pdf  $f_{\tilde{d}|\tilde{h}}(d|1)$  of  $\tilde{d}$  given  $\tilde{h} = 1$  equals 0 if  $d = 3$ . Consequently, under the conditional independence

assumptions,

$$f_{\tilde{d}_2 | \tilde{d}_1}(d | 3) = \frac{f_{\tilde{d}_1, \tilde{d}_2}(3, d)}{f_{\tilde{d}_1}(3)} \quad (6.23)$$

$$= \frac{p_{\tilde{h}}(0)f_{\tilde{d}_1 | \tilde{h}}(3 | 0)f_{\tilde{d}_2 | \tilde{h}}(d | 0) + p_{\tilde{h}}(1)f_{\tilde{d}_1 | \tilde{h}}(3 | 1)f_{\tilde{d}_2 | \tilde{h}}(d | 1)}{p_{\tilde{h}}(0)f_{\tilde{d}_1 | \tilde{h}}(3 | 0) + p_{\tilde{h}}(1)f_{\tilde{d}_1 | \tilde{h}}(3 | 1)} \quad (6.24)$$

$$= \frac{p_{\tilde{h}}(0)f_{\tilde{d}_1 | \tilde{h}}(3 | 0)f_{\tilde{d}_2 | \tilde{h}}(d | 0)}{p_{\tilde{h}}(0)f_{\tilde{d}_1 | \tilde{h}}(3 | 0)} \quad (6.25)$$

$$= f_{\tilde{d}_2 | \tilde{h}}(d | 0) \quad (6.26)$$

$$= \frac{1}{5} \quad \text{if } 1 \leq d < 6, \quad (6.27)$$

$$(6.28)$$

and zero otherwise. Intuitively, if the first distance is three, then it the ball must be solid, so the conditional pdf of the second distance given this information is the same as the conditional pdf of the distance given that the ball is not hollow.

#### 6.4 (Computer defect)

a)

$$P(\tilde{t} > 2) = P(\tilde{t} > 2, \tilde{d} = 0) + P(\tilde{t} > 2, \tilde{d} = 1) \quad (6.29)$$

$$= P(\tilde{d} = 0)P(\tilde{t} > 2 | \tilde{d} = 0) + P(\tilde{d} = 1)P(\tilde{t} > 2 | \tilde{d} = 1) \quad (6.30)$$

$$= 0.9 \int_2^\infty \exp(-t) dt + 0.1 \int_2^\infty 2 \exp(-2t) dt \quad (6.31)$$

$$= 0.9 \exp(-2) + 0.1 \exp(-4) \quad (6.32)$$

$$= 0.124. \quad (6.33)$$

b)

$$p_{\tilde{d} | \tilde{t}}(1 | 2) = \frac{p_{\tilde{d}}(1)f_{\tilde{t} | \tilde{d}}(2 | 1)}{f_{\tilde{t}}(2)} \quad (6.34)$$

$$= \frac{p_{\tilde{d}}(1)f_{\tilde{t} | \tilde{d}}(2 | 1)}{p_{\tilde{d}}(0)f_{\tilde{t} | \tilde{d}}(2 | 0) + p_{\tilde{d}}(1)f_{\tilde{t} | \tilde{d}}(2 | 1)} \quad (6.35)$$

$$= \frac{0.1 \cdot 2 \exp(-4)}{0.9 \cdot \exp(-1) + 0.1 \cdot 2 \exp(-4)} \quad (6.36)$$

$$= 1.09\%. \quad (6.37)$$

#### 6.5 (Wolf)

$$p_{\tilde{s} | \tilde{x}} \left( 1 \left| \begin{bmatrix} -1 \\ -1 \end{bmatrix} \right. \right) = \frac{p_{\tilde{s}}(1)f_{\tilde{x} | \tilde{s}} \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix} \left| 1 \right. \right)}{p_{\tilde{s}}(1)f_{\tilde{x} | \tilde{s}} \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix} \left| 1 \right. \right) + p_{\tilde{s}}(0)f_{\tilde{x} | \tilde{s}} \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix} \left| 0 \right. \right)} \quad (6.38)$$

$$= \frac{0.25 \cdot \frac{1}{50}}{0.25 \cdot \frac{1}{50} + 0.75 \cdot \frac{1}{100}} \quad (6.39)$$

$$= 0.4. \quad (6.40)$$

## 6.6 (Potatoes)

a) We have

$$f_{\tilde{x}}(x) = \sum_{b=0}^1 \sum_{w=0}^1 p_{\tilde{b},\tilde{w}}(b,w) f_{\tilde{x}|\tilde{b},\tilde{w}}(x|b,w). \quad (6.41)$$

The pdf is therefore a mixture of Gaussians, which is not Gaussian.

b) We use the empirical pmfs

$$p_{\tilde{w}}(1) = \frac{1}{3}, \quad p_{\tilde{w}}(0) = \frac{2}{3}, \quad (6.42)$$

$$p_{\tilde{b}}(1) = \frac{2}{3}, \quad p_{\tilde{b}}(0) = \frac{1}{3}. \quad (6.43)$$

c) The random variables are independent, since the empirical joint pmf factorizes into the product of the marginal pmfs

$$p_{\tilde{b},\tilde{w}}(0,0) = \frac{2}{9} = p_{\tilde{w}}(0) p_{\tilde{b}}(0), \quad (6.44)$$

$$p_{\tilde{b},\tilde{w}}(0,1) = \frac{1}{9} = p_{\tilde{w}}(1) p_{\tilde{b}}(0), \quad (6.45)$$

$$p_{\tilde{b},\tilde{w}}(1,0) = \frac{4}{9} = p_{\tilde{w}}(0) p_{\tilde{b}}(1), \quad (6.46)$$

$$p_{\tilde{b},\tilde{w}}(1,1) = \frac{2}{9} = p_{\tilde{w}}(1) p_{\tilde{b}}(1). \quad (6.47)$$

d) We compare

$$p_{\tilde{b}|\tilde{x},\tilde{w}}(1|40,1) = \frac{p_{\tilde{b},\tilde{w}}(1,1) f_{\tilde{x}|\tilde{b},\tilde{w}}(40|1,1)}{p_{\tilde{w}}(1) f_{\tilde{x}|\tilde{w}}(40|1)} \quad (6.48)$$

$$= \frac{\frac{2}{9} \cdot 0.02}{p_{\tilde{w}}(1) f_{\tilde{x}|\tilde{w}}(40|1)} \quad (6.49)$$

with

$$p_{\tilde{b}|\tilde{x},\tilde{w}}(0|40,1) = \frac{p_{\tilde{b},\tilde{w}}(0,1) f_{\tilde{x}|\tilde{b},\tilde{w}}(40|0,1)}{p_{\tilde{w}}(1) f_{\tilde{x}|\tilde{w}}(40|1)} \quad (6.50)$$

$$= \frac{\frac{1}{9} \cdot 0.01}{p_{\tilde{w}}(1) f_{\tilde{x}|\tilde{w}}(40|1)}. \quad (6.51)$$

Consequently, according to the model, it is more likely that the beetle was present that year.

## 6.7 (Halloween parade)

a)

$$P(\tilde{w} \neq \tilde{r}) = p_{\tilde{w},\tilde{r}}(0,1) + p_{\tilde{w},\tilde{r}}(1,0) \quad (6.52)$$

$$= p_{\tilde{r}}(1) p_{\tilde{w}|\tilde{r}}(0|1) + p_{\tilde{r}}(0) p_{\tilde{w}|\tilde{r}}(1|0) \quad (6.53)$$

$$= 0.24. \quad (6.54)$$

b) It is not reasonable to assume that the forecast and the humidity are independent, even if we know that the forecast does not take the humidity into account. The reason is that both variables are linked through the rain. For example, if  $\tilde{w} = 1$  then the humidity is probably high (because it probably will rain) and if  $\tilde{w} = 0$  the humidity is probably low (because it probably won't rain). In contrast, conditioned on whether it rains or



not, it is quite reasonable to assume independence of  $\tilde{w}$  and  $\tilde{h}$ , because  $\tilde{h}$  is not used to produce the forecast.

c) By the chain rule for discrete and continuous random variables

$$p_{\tilde{r}|\tilde{w},\tilde{h}}(r|w,h) = \frac{f_{\tilde{h}|\tilde{w},\tilde{r}}(h|w,r)p_{\tilde{w}|\tilde{r}}(w|r)p_{\tilde{r}}(r)}{p_{\tilde{w}}(w)f_{\tilde{h}|\tilde{w}}(h|w)} \quad (6.55)$$

$$= \frac{f_{\tilde{h}|\tilde{r}}(h|r)p_{\tilde{w}|\tilde{r}}(w|r)p_{\tilde{r}}(r)}{\sum_{r=0}^1 f_{\tilde{h}|\tilde{r}}(h|r)p_{\tilde{w}|\tilde{r}}(w|r)p_{\tilde{r}}(r)}. \quad (6.56)$$

This expression equals to zero unless  $h$  is between 0.5 and 0.7. If  $h$  is between 0.6 and 0.7 then it is equal to 1. If  $h$  is between 0.5 and 0.6 then

$$p_{\tilde{r}|\tilde{w},\tilde{h}}(1|1,h) = \frac{0.2 \cdot 0.8 \cdot 5}{0.2 \cdot 0.8 \cdot 5 + 0.8 \cdot 0.25 \cdot 2} \quad (6.57)$$

$$= 0.667, \quad (6.58)$$

$$p_{\tilde{r}|\tilde{w},\tilde{h}}(1|0,h) = \frac{0.2 \cdot 0.2 \cdot 5}{0.2 \cdot 0.2 \cdot 5 + 0.8 \cdot 0.75 \cdot 2} \quad (6.59)$$

$$= 0.143. \quad (6.60)$$

The conditional pmf is only well defined for  $0.1 \leq h \leq 0.7$ ,

$$p_{\tilde{r}|\tilde{w},\tilde{h}}(1|w,h) = \begin{cases} 0 & \text{if } 0.1 \leq h \leq 0.5 \\ 0.667 & \text{if } 0.5 \leq h \leq 0.6 \text{ and } w = 1 \\ 0.143 & \text{if } 0.5 \leq h \leq 0.6 \text{ and } w = 0 \\ 1 & \text{if } 0.6 \leq h \leq 0.7. \end{cases} \quad (6.61)$$

We predict no rain (a) if the humidity is between 0.1 and 0.5 or (b) if  $0.5 \leq h \leq 0.6$  and  $w = 0$ . Otherwise, we predict rain.

d) According to the model, if the humidity is between 0.1 and 0.5 then we are always right. Similarly, if the humidity is between 0.6 and 0.7 we are also always right. For  $h$  between 0.5 and 0.6 we make a mistake under two scenarios: (1)  $w = 0$  and it rains and (2)  $w = 1$  and it does not rain. The probability of these two events happening equals

$$\begin{aligned} P(\text{Error}) &= P(\tilde{r} = 0, w = 1, 0.5 \leq \tilde{h} \leq 0.6) + P(\tilde{r} = 1, w = 0, 0.5 \leq \tilde{h} \leq 0.6) \\ &= \int_{h=0.5}^{0.6} f_{\tilde{h}|\tilde{r}}(h|0)p_{\tilde{w}|\tilde{r}}(1|0)p_{\tilde{r}}(0) dh + \int_{h=0.5}^{0.6} f_{\tilde{h}|\tilde{r}}(h|1)p_{\tilde{w}|\tilde{r}}(0|1)p_{\tilde{r}}(1) dh \\ &= \int_{h=0.5}^{0.6} 0.8 \cdot 2 \cdot 0.25 dh + \int_{h=0.5}^{0.6} 0.2 \cdot 5 \cdot 0.2 dh = 0.06. \end{aligned}$$

6.8 (Fish)

a) The position is a uniform random variable since

$$f_{\tilde{x}}x = p_{\tilde{s}}(a)f_{\tilde{x}|\tilde{s}}(x|a) + p_{\tilde{s}}(b)f_{\tilde{x}|\tilde{s}}(x|b) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (6.62)$$

b) By the chain rule for discrete and continuous random variables

$$p_{\tilde{s}|\tilde{x}}(b|0.25) = \frac{p_{\tilde{s}}(b)f_{\tilde{x}|\tilde{s}}(0.25|b)}{f_{\tilde{x}}(0.25)} \quad (6.63)$$

$$= 0.25. \quad (6.64)$$

c) Integrating the conditional pdfs we obtain

$$F_{\tilde{x}|\tilde{s}}(x|a) = \begin{cases} 0 & \text{if } x < 0, \\ 2x - x^2 & \text{if } 0 \leq x \leq 1, \\ 1 & \text{if } x > 1. \end{cases} \quad (6.65)$$

$$F_{\tilde{x}|\tilde{s}}(x|b) = \begin{cases} 0 & \text{if } x < 0, \\ x^2 & \text{if } 0 \leq x \leq 1, \\ 1 & \text{if } x > 1, \end{cases} \quad (6.66)$$

d) Let  $u_1$  and  $u_2$  be the two uniform samples. We use the first one to sample from  $\tilde{s}$  by assigning  $\tilde{s} = a$  if  $u_1$  is below 0.5, and  $\tilde{s} = b$  otherwise. In this case,  $u_1 = 0.8$  so we obtain  $\tilde{s} = b$ . Then we sample from the conditional distribution of  $\tilde{x}$  given  $\tilde{s} = b$  using  $u_2$  via the inverse transform method. We have  $F_{\tilde{x}|\tilde{s}}(0.8|b) = 0.64$  (computed by solving for  $x^2 = 0.64$ ), so the sample of  $\tilde{x}$  equals 0.8.

6.9 (Samples and conditional independence) In (1) if we only consider the samples for which  $c = 0$ , the density of the samples in the vertical direction looks the same for all values of  $a$  between 0 and 1, and is zero otherwise, which suggests that the conditional pdf of  $\tilde{b}$  given  $\tilde{a} = a$  and  $\tilde{c} = 0$  is the same for all values of  $a$ . Similarly, when  $c = 1$ , the vertical density again looks the same for all values of  $a$  between 1 and 2, and is zero otherwise, which suggests that the conditional pdf of  $\tilde{b}$  given  $\tilde{a} = a$  and  $\tilde{c} = 1$  is the same for all values of  $a$ . Since  $\tilde{c}$  only has these two possible values, this would imply that  $\tilde{a}$  and  $\tilde{b}$  are conditionally independent given  $\tilde{c}$ .

In (2)  $\tilde{a}$  and  $\tilde{b}$  again look conditionally independent given  $\tilde{c} = 0$ , but not given  $\tilde{c} = 1$ , because for those samples the vertical density for  $a$  between 0.5 and 1 is clearly different to the density for  $a$  between 1 and 1.5, or between 1.5 and 2. Consequently,  $\tilde{a}$  and  $\tilde{b}$  are not conditionally independent given  $\tilde{c}$ .

In (3) if we only consider the samples for which  $c = 0$ , the vertical density for  $a$  between 0 and 1 is clearly different to the density for  $a$  between 1 and 2. The same occurs for the samples for which  $c = 1$ . Consequently,  $\tilde{a}$  and  $\tilde{b}$  are not conditionally independent given  $\tilde{c}$ .

6.10 (Chad)

a) The kernel density estimate is of the form

$$f_{\tilde{t}|\tilde{c}}(t|0) = \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{1}{2} \Pi\left(\frac{t - d_{0,i}}{2}\right), \quad (6.67)$$

$$f_{\tilde{t}|\tilde{c}}(t|1) = \frac{1}{n_1} \sum_{i=1}^{n_1} \frac{1}{2} \Pi\left(\frac{t - d_{1,i}}{2}\right), \quad (6.68)$$

$$(6.69)$$

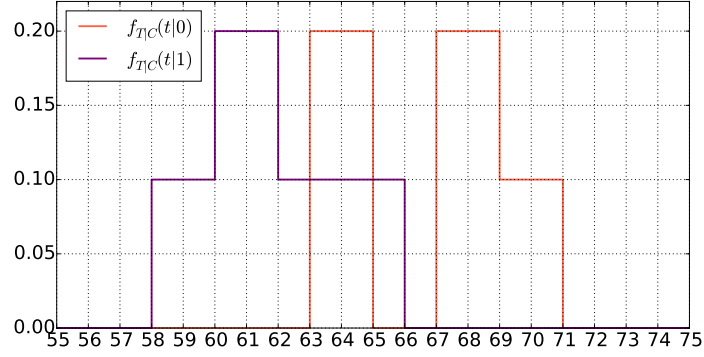
where  $\Pi$  is a rectangular kernel with unit width,  $d_{1,0}, \dots, d_{n_1,0}$  the temperatures when Chad is not there and  $d_{1,1}, \dots, d_{n_1,1}$  the temperatures when he is there. The estimate is shown in Figure 6.1.

b) We have

$$f_{\tilde{t}|\tilde{c}}(68|0) = 0.2 > 0 = f_{\tilde{t}|\tilde{c}}(68|1), \quad (6.70)$$

so the ML estimate is that Chad is not at the office.

c) That probability does not make sense because the parameter representing the presence of Chad is deterministic.



**Figure 6.1** Kernel density estimate for Exercise 6.10.

d) The empirical pmf is

$$p_{\bar{c}}(0) = \frac{5}{15} = \frac{1}{3}, \quad (6.71)$$

$$p_{\bar{c}}(1) = \frac{10}{15} = \frac{2}{3}. \quad (6.72)$$

e) Applying Bayes' rule,

$$p_{\bar{c}|\tilde{t}}(0|64) = \frac{p_{\bar{c}}(0) f_{\tilde{t}|\bar{c}}(64|0)}{p_{\bar{c}}(0) f_{\tilde{t}|\bar{c}}(64|0) + p_{\bar{c}}(1) f_{\tilde{t}|\bar{c}}(64|1)} \quad (6.73)$$

$$= \frac{\frac{1}{3}0.2}{\frac{1}{3}0.2 + \frac{2}{3}0.1} \quad (6.74)$$

$$= \frac{1}{2}, \quad (6.75)$$

$$p_{\bar{c}|\tilde{t}}(1|64) = 1 - p_{\bar{c}|\tilde{t}}(0|64) \quad (6.76)$$

$$= \frac{1}{2}. \quad (6.77)$$

According to the posterior pmf there is a 50 % chance that Chad is there.

f) Both  $f_{\tilde{t}|\bar{c}}(57|0)$  and  $f_{\tilde{t}|\bar{c}}(68|1)$  are zero so both the ML estimate and Bayesian posterior are inconclusive. If we use a parametric distribution such as a Gaussian to fit the data, then  $f_{\tilde{t}|\bar{c}}(57|0)$  and  $f_{\tilde{t}|\bar{c}}(68|1)$  would not be set to zero as long as the distribution has nonzero values on all of the real line (as is the case for a Gaussian pdf). This would allow us to apply MAP or ML estimation. A nonparametric solution would be to use a kernel with a larger width.

6.11 (Mixture model with fixed mean) We have

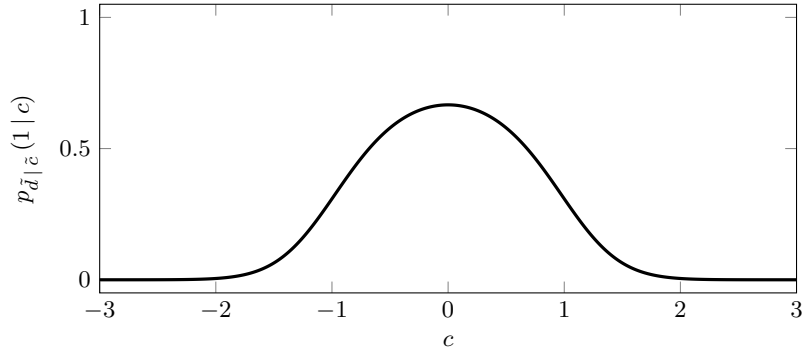
$$p_{\tilde{s}|\tilde{c}}(1|c) = \frac{p_{\tilde{s}}(1) f_{\tilde{c}|\tilde{s}}(c|1)}{p_{\tilde{s}}(0) f_{\tilde{c}|\tilde{s}}(c|0) + p_{\tilde{s}}(1) f_{\tilde{c}|\tilde{s}}(c|1)} \quad (6.78)$$

$$= \frac{\frac{\theta}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{c^2}{2\sigma_1^2}\right)}{\frac{\theta}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{c^2}{2\sigma_1^2}\right) + \frac{1-\theta}{\sigma_0 \sqrt{2\pi}} \exp\left(-\frac{c^2}{2\sigma_0^2}\right)} \quad (6.79)$$

$$= \frac{1}{1 + \frac{(1-\theta)\sigma_1}{\theta\sigma_0} \exp\left(\frac{c^2}{2} \left(\frac{1}{\sigma_1^2} - \frac{1}{\sigma_0^2}\right)\right)} \quad (6.80)$$

$$(6.81)$$

For  $\theta := 0.5$ ,  $\sigma_0 := 1$ , and  $\sigma_1 := 0.5$ , the function looks like this:



6.12 (K-means and Gaussian mixture models) There are two main differences. First, when fitting a Gaussian mixture model we make *soft* assignments to each cluster, represented by the membership probabilities, as opposed to the *hard* assignment in Lloyd's algorithm, where each data point is assigned to a single cluster. Second, the Gaussian mixture model estimates the covariance matrix of each cluster, as well as its mean, whereas in *k*-means the covariance structure of each cluster is implicitly assumed to be the same.

6.13 (Alternative model for coin flip)

- a) The model makes sense because we suspect that the coin that is more prone to be heads than tails, but we are uncertain about the exact bias so we just assume it is uniform between  $1/2$  (fair coin) and  $1$  (coin that always lands heads). To compute the probability of heads, we integrate over all the possible values of the Bernoulli parameter. Let  $\tilde{r}$  be the result of the coin flip and  $\tilde{\theta}$  the parameter of the Bernoulli. Since  $\tilde{r}$  is uniform between  $1/2$  and  $1$   $f_{\tilde{\theta}}(\theta)$  is equal to  $2$  for  $1/2 \leq \theta \leq 1$  and zero otherwise.

$$P(\text{heads}) = p_{\tilde{r}}(1) = \int_{\theta=1/2}^1 f_{\tilde{\theta}}(\theta) p_{\tilde{r}|\tilde{\theta}}(1|\theta) d\theta = \int_{\theta=1/2}^1 2\theta d\theta = 1 - \frac{1}{4} = \frac{3}{4}, \quad (6.82)$$

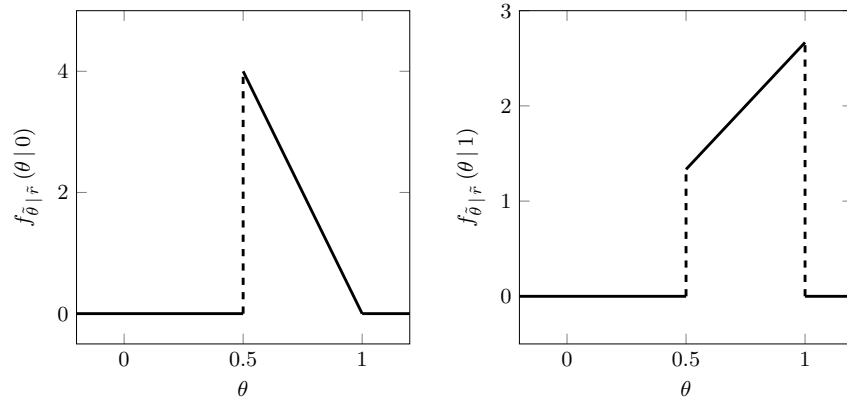
$$P(\text{tails}) = 1 - P(\text{heads}) = \frac{1}{4}. \quad (6.83)$$

- b) The conditional pdf on the bias of the coin flip conditioned on tails equals

$$f_{\tilde{\theta}|\tilde{r}}(\theta|0) = \frac{f_{\tilde{\theta}}(\theta) p_{\tilde{r}|\tilde{\theta}}(0)}{p_{\tilde{r}}(0)} \quad (6.84)$$

$$= \frac{f_{\tilde{\theta}}(\theta)(1-\theta)}{1/4} = \begin{cases} 8(1-\theta) & 1/2 \leq \theta \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (6.85)$$

Note that the conditional pdf is more skewed towards 1/2, as shown on Figure 6.2. It makes sense because the coin flip is tails. Intuitively the model should be adjusted towards the coin flip being less biased.



**Figure 6.2** Conditional pdf of the bias of the coin flip given tails (left) and heads (right).

The conditional pdf on the bias of the coin flip conditioned on heads equals

$$f_{\tilde{\theta}|\tilde{r}}(\theta|1) = \frac{f_{\tilde{\theta}}(\theta) p_{\tilde{r}|\tilde{\theta}}(1)}{p_{\tilde{r}}(1)} = \frac{f_{\tilde{\theta}}(\theta) \theta}{3/4} = \begin{cases} \frac{8\theta}{3} & 1/2 \leq \theta \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (6.86)$$

In this case the conditional pdf is more skewed towards 1/2, as shown on Figure 6.2. Intuitively, we are adjusting the model by incorporating some evidence that supports that the coin might be biased towards heads.

- c) The data seem to indicate that the probability of heads is very small. However, no matter what evidence we observe, our model will always assign zero probability density to any value of the bias below 0.5, so we should probably consider changing the prior.

6.14 (Two coin flips)

a)

$$p_{\tilde{x}_1, \tilde{x}_2}(1, 1) = \int_{\theta=-\infty}^{\infty} f_{\tilde{\theta}}(\theta) p_{\tilde{x}_1 | \tilde{\theta}}(1 | \theta) p_{\tilde{x}_2 | \tilde{\theta}}(1 | \theta) d\theta \quad (6.87)$$

$$= \int_0^1 \theta^2 d\theta = \frac{1}{3}, \quad (6.88)$$

$$p_{\tilde{x}_1, \tilde{x}_2}(0, 1) = \int_0^1 \theta(1 - \theta) d\theta = \frac{1}{6}, \quad (6.89)$$

$$p_{\tilde{x}_1, \tilde{x}_2}(1, 0) = \int_0^1 \theta(1 - \theta) d\theta = \frac{1}{6}, \quad (6.90)$$

$$p_{\tilde{x}_1, \tilde{x}_2}(0, 0) = \int_0^1 (1 - \theta)^2 d\theta = \frac{1}{3}. \quad (6.91)$$

b) The posterior pdf equals

$$f_{\tilde{\theta} | \tilde{x}_1, \tilde{x}_2}(\theta | 1, 1) = \frac{f_{\tilde{\theta}}(\theta) p_{\tilde{x}_1 | \tilde{\theta}}(1 | \theta) p_{\tilde{x}_2 | \tilde{\theta}}(1 | \theta)}{p_{\tilde{x}_1, \tilde{x}_2}(1, 1)} \quad (6.92)$$

$$= \frac{\theta^2}{1/3} = 3\theta^2. \quad (6.93)$$

The conditional probability that  $\tilde{\theta} < 1/2$  equals

$$P\left(\tilde{\theta} < \frac{1}{2}\right) = \int_{\theta=0}^{1/2} f_{\tilde{\theta} | \tilde{x}_1, \tilde{x}_2}(\theta | 1, 1) d\theta \quad (6.94)$$

$$= \int_{\theta=0}^{1/2} 3\theta^2 d\theta \quad (6.95)$$

$$= \frac{1}{8}. \quad (6.96)$$

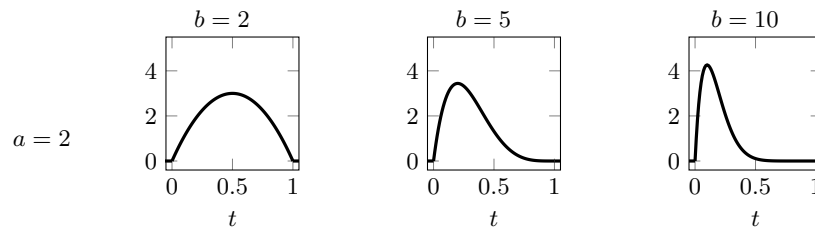
6.15 (PT Cruiser) Let  $\tilde{t}$  denote the number of times the car is used. Under the assumptions, given  $\tilde{\theta} = \theta$ , the conditional distribution of  $\tilde{t}$  is geometric with parameter  $\theta$ . The posterior distribution given  $\tilde{t} = t$  is therefore

$$f_{\tilde{\theta} | \tilde{t}}(\theta | t) = \frac{f_{\tilde{\theta}}(\theta) p_{\tilde{t} | \tilde{\theta}}(t | \theta)}{p_{\tilde{t}}(t)} \quad (6.97)$$

$$= \frac{f_{\tilde{\theta}}(\theta) p_{\tilde{t} | \tilde{\theta}}(t | \theta)}{\int_u f_{\tilde{\theta}}(u) p_{\tilde{t} | \tilde{\theta}}(t | u) du} \quad (6.98)$$

$$= \frac{\theta(1 - \theta)^{t-1}}{\int_{u=0}^1 u(1 - u)^{t-1} du}, \quad (6.99)$$

which is a beta distribution with parameters  $a = 2$  and  $b = t$ . In order to understand how it changes with  $t$ , we need to look at the shape of the beta distribution for  $a = 2$  and different values of  $b$ :



For small  $t$  the posterior is very flat (see  $a = 2$ ,  $b = 1$ ). As  $t$  grows, the pdf is increasingly more skewed towards zero and more concentrated, indicating that the car is likely to be in good condition after all.

---

## Averaging

### Exercises

#### 7.1 (Race)

a) By the formula for the mean of a function of two random variables,

$$\mathbb{E} [\tilde{d} - \tilde{a}] = \sum_{d \in \{5, 10, 20\}} \sum_{a \in \{5, 10, 20\}} (d - a) p_{\tilde{d}, \tilde{a}}(d, a) \quad (7.1)$$

$$= -5 \cdot 0.1 + 5 \cdot 0.1 + 15 \cdot 0.1 + 10 \cdot 0.2 \quad (7.2)$$

$$= 3.5, \quad (7.3)$$

so the mean difference is 3.5 km.

b) The mean equals

$$\mathbb{E} [\tilde{a}] = \sum_{d \in \{5, 10, 20\}} \sum_{a \in \{5, 10, 20\}} a p_{\tilde{d}, \tilde{a}}(d, a) \quad (7.4)$$

$$= 5 (0.2 + 0.1 + 0.1) + 10 (0.1 + 0.2 + 0.2) + 20 \cdot 0.1 \quad (7.5)$$

$$= 2 + 5 + 2 = 9. \quad (7.6)$$

The mean square equals

$$\mathbb{E} [\tilde{a}^2] = \sum_{d \in \{5, 10, 20\}} \sum_{a \in \{5, 10, 20\}} a^2 p_{\tilde{d}, \tilde{a}}(d, a) \quad (7.7)$$

$$= 25 (0.2 + 0.1 + 0.1) + 100 (0.1 + 0.2 + 0.2) + 400 \cdot 0.1 \quad (7.8)$$

$$= 10 + 50 + 40 = 100. \quad (7.9)$$

Consequently the variance equals

$$\text{Var} [\tilde{a}] = \mathbb{E} [\tilde{a}^2] - \mathbb{E} [\tilde{a}]^2 \quad (7.10)$$

$$= 100 - 81 = 19. \quad (7.11)$$



## 7.2 (Mean of a function)

$$\mathbb{E}[\tilde{b}] := \sum_{b \in B} b p_{\tilde{b}}(b) \quad (7.12)$$

$$= \sum_{b \in B} b P(\tilde{b} = b) \quad (7.13)$$

$$= \sum_{b \in B} b \sum_{\{a: h(a)=b\}} P(\tilde{a} = a) \quad (7.14)$$

$$= \sum_{b \in B} \sum_{\{a: h(a)=b\}} h(a) P(\tilde{a} = a) \quad (7.15)$$

$$= \sum_{a \in A} h(a) p_{\tilde{a}}(a). \quad (7.16)$$

Where the last step follows from the fact that each  $a \in A$  is in exactly one of the sets  $\{a : h(a) = b\}$ .

- 7.3 (Presents) We define a random variable  $\tilde{a}_i$  that is equal to one when kid  $i$  gets the present bought by their own parents, and to zero otherwise. The probability of the event of kid  $i$  getting their own present is  $1/n$  so

$$p_{\tilde{a}_i}(0) = 1 - \frac{1}{n}, \quad (7.17)$$

$$p_{\tilde{a}_i}(1) = \frac{1}{n}. \quad (7.18)$$

By linearity of expectation

$$\mathbb{E}[\text{Number of kids that get their own present}] = \mathbb{E}\left[\sum_{i=1}^n \tilde{a}_i\right] \quad (7.19)$$

$$= \sum_{i=1}^n \mathbb{E}[\tilde{a}_i] \quad (7.20)$$

$$= \sum_{i=1}^n p_{\tilde{a}_i}(1) \quad (7.21)$$

$$= 1. \quad (7.22)$$

On average one kid gets their own present.

- 7.4 (Cats and dogs) The mean is equal to 1.05 (from the example) and the mean square equals

$$\mathbb{E}[(\tilde{c} + \tilde{d})^2] = \sum_{c=0}^3 \sum_{d=0}^2 (c+d)^2 p_{\tilde{c}, \tilde{d}}(c, d) \quad (7.23)$$

$$= 0.15 + 4 \cdot 0.1 + 9 \cdot 0.05 + 0.2 + 4 \cdot 0.05 + 9 \cdot 0.03 + 4 \cdot 0.05 + 9 \cdot 0.02 \\ = 2.05. \quad (7.24)$$

The standard deviation therefore equals

$$\sqrt{\text{Var}[\tilde{c} + \tilde{d}]} = \sqrt{\mathbb{E}[(\tilde{c} + \tilde{d})^2] - \mathbb{E}[\tilde{c} + \tilde{d}]^2} \quad (7.25)$$

$$= \sqrt{2.05 - 1.05^2} = 0.973. \quad (7.26)$$

7.5 (Mean squared difference)

$$\mathbb{E}[(\tilde{a} - \tilde{b})^2] = \mathbb{E}[\tilde{a}^2 + \tilde{b}^2 - 2\tilde{a}\tilde{b}] \quad (7.27)$$

$$= \mathbb{E}[\tilde{a}^2] + \mathbb{E}[\tilde{b}^2] - 2\mathbb{E}[\tilde{a}\tilde{b}] \quad \text{by linearity} \quad (7.28)$$

$$= \mathbb{E}[\tilde{a}^2] + \mathbb{E}[\tilde{b}^2] - 2\mathbb{E}[\tilde{x}_1] \mathbb{E}[\tilde{b}] \quad \text{by independence} \quad (7.29)$$

$$= \mathbb{E}[\tilde{a}^2] - \mu^2 + \mathbb{E}[\tilde{b}^2] - \mu^2 \quad (7.30)$$

$$= 2\sigma^2. \quad (7.31)$$

7.6 (Basketball player) Let  $\tilde{x}$  be a random variable that represents the points scored by the player. We have

$$P(\tilde{x} = 2) = P(\{\text{attempts 2 point shot}\} \cap \{\text{makes shot}\}) \quad (7.32)$$

$$= P(\{\text{attempts 2 point shot}\})P(\{\text{makes shot}\} \mid \{\text{attempts 2 point shot}\}) \\ = 0.5\alpha, \quad (7.33)$$

$$P(\tilde{x} = 3) = P(\{\text{attempts 3 point shot}\} \cap \{\text{makes shot}\}) \quad (7.34)$$

$$= P(\{\text{attempts 3 point shot}\})P(\{\text{makes shot}\} \mid \{\text{attempts 3 point shot}\}) \\ = 0.4(1 - \alpha), \quad (7.35)$$

$$P(\tilde{x} = 0) = 1 - 0.5\alpha - 0.4(1 - \alpha). \quad (7.36)$$

The mean, mean square and variance equal

$$\mathbb{E}[\tilde{x}] = 2P(\tilde{x} = 2) + 3P(\tilde{x} = 3) \quad (7.37)$$

$$= \alpha + 1.2 - 1.2\alpha \quad (7.38)$$

$$= 1.2 - 0.2\alpha, \quad (7.39)$$

$$\mathbb{E}[\tilde{x}^2] = 4P(\tilde{x} = 2) + 9P(\tilde{x} = 3) \quad (7.40)$$

$$= 2\alpha + 3.6 - 3.6\alpha \quad (7.41)$$

$$= 3.6 - 1.6\alpha, \quad (7.42)$$

$$\text{Var}[\tilde{x}] = \mathbb{E}[\tilde{x}^2] - \mathbb{E}[\tilde{x}]^2 \quad (7.43)$$

$$= 3.6 - 1.6\alpha - (1.2 - 0.2\alpha)^2 \quad (7.44)$$

$$= -0.04\alpha^2 - 1.12\alpha + 2.16. \quad (7.45)$$

Therefore, the standard deviation is equal to  $\sqrt{-0.04\alpha^2 - 1.12\alpha + 2.16}$ . The derivative of the standard deviation with respect to  $\alpha$  equals

$$\frac{-0.08\alpha - 1.12}{2\sqrt{-0.04\alpha^2 - 1.12\alpha + 2.16}}, \quad (7.46)$$

which is negative for  $\alpha > 0$ . Consequently, for  $0 \leq \alpha \leq 1$  the standard deviation decreases as we increase  $\alpha$ . It is minimized at  $\alpha = 1$  (the player only attempts 2-point shots), where it is equal to 1. It is maximized at  $\alpha = 0$  (the player only attempts 3-point shots), where it is equal to 1.47.

7.7 (Computer defect)

- a) The random variable  $\tilde{d}$  is discrete, so the conditional mean of  $\tilde{t}$  is also discrete (it is a function of  $\tilde{d}$ ). The mean of an exponential random variable with parameter  $\lambda$  equal  $\frac{1}{\lambda}$ .

As a result, we have  $\mu_{\tilde{t}|\tilde{d}}(1) = 0.5$  and  $\mu_{\tilde{t}|\tilde{d}}(0) = 1$ , so the pmf of  $\mu_{\tilde{t}|\tilde{d}}(\tilde{d})$  is

$$p_{\mu_{\tilde{t}|\tilde{d}}(\tilde{d})}(0.5) = P(\tilde{d} = 1) \quad (7.47)$$

$$= 0.1, \quad (7.48)$$

$$p_{\mu_{\tilde{t}|\tilde{d}}(\tilde{d})}(1) = P(\tilde{d} = 0) \quad (7.49)$$

$$= 0.9, \quad (7.50)$$

$$p_{\tilde{a}}(a) = 0, \quad \text{for } \tilde{a} \notin \{0.5, 1\}. \quad (7.51)$$

b) From the previous question,

$$\mathbb{E}[\tilde{t}] = P(\tilde{d} = 0)\mu_{\tilde{t}|\tilde{d}}(0) + P(\tilde{d} = 1)\mu_{\tilde{t}|\tilde{d}}(1) \quad (7.52)$$

$$= 0.95. \quad (7.53)$$

The mean and variance of an exponential random variable with parameter  $\lambda$  equal  $\frac{1}{\lambda}$  and  $\frac{1}{\lambda^2}$  respectively, so its mean square is  $\frac{2}{\lambda^2}$ . Therefore,

$$\mathbb{E}[\tilde{t}^2 | \tilde{d} = 0] = 2, \quad (7.54)$$

$$\mathbb{E}[\tilde{t}^2 | \tilde{d} = 1] = \frac{1}{2}. \quad (7.55)$$

By iterated expectation,

$$\mathbb{E}[\tilde{t}^2] = P(\tilde{d} = 0)\mu_{\tilde{t}^2|\tilde{d}}(0) + P(\tilde{d} = 1)\mu_{\tilde{t}^2|\tilde{d}}(1) \quad (7.56)$$

$$= 1.85. \quad (7.57)$$

We conclude that

$$\text{Var}[\tilde{t}] = \mathbb{E}[\tilde{t}^2] - \mathbb{E}[\tilde{t}]^2 \quad (7.58)$$

$$= 0.9475. \quad (7.59)$$

c) The number of computers that break down during the first year is binomial with parameters  $n = 100$  and  $\theta$ , which represents the probability that  $\tilde{t} \leq 1$ :

$$\theta = P(\tilde{t} \leq 1) = \int_0^1 f_{\tilde{t}}(t) dt \quad (7.60)$$

$$= \int_0^1 (0.9 \exp(-t) + 0.1 \exp(-2t)) dt \quad (7.61)$$

$$= -0.9 \exp(-t)]_0^1 - 0.1 \exp(-2t)]_0^1 \quad (7.62)$$

$$= 0.9(1 - \exp(-1)) + 0.1(1 - \exp(-2)) \quad (7.63)$$

$$= 0.655. \quad (7.64)$$

The binomial therefore has mean  $n\theta = 65.5$  and variance  $n\theta(1 - \theta) = 22.6$ .

## 7.8 (Hen breeds)

a) Let  $\tilde{b}_1, \dots, \tilde{b}_4$  represent the eggs laid by the bantam hens and  $\tilde{s}_1, \tilde{s}_2, \tilde{s}_3$  the ones laid by the Sussex hens. Since it's always warm, the mean of the eggs laid by each bantam

chicken equals,

$$\mu_{\text{bantam}} = \sum_{i=4}^6 i \mathbf{P}(\text{eggs} = i) \quad (7.65)$$

$$= 4 \cdot 0.25 + 5 \cdot 0.5 + 6 \cdot 0.25 \quad (7.66)$$

$$= 5, \quad (7.67)$$

and the mean of the eggs laid by each Sussex hen is

$$\mu_{\text{sussex}} = \sum_{i=3}^6 i \mathbf{P}(\text{eggs} = i) \quad (7.68)$$

$$= 3 \cdot 0.25 + 4 \cdot 0.25 + 5 \cdot 0.25 + 6 \cdot 0.25 \quad (7.69)$$

$$= 4.5. \quad (7.70)$$

By linearity of expectation, which does not require independence assumptions, the mean of the total number of eggs  $\tilde{t}$  is

$$\mathbb{E}[\tilde{t}] = \mathbb{E}\left[\sum_{i=1}^4 \tilde{b}_i + \sum_{i=1}^3 \tilde{s}_i\right] \quad (7.71)$$

$$= \sum_{i=1}^4 \mathbb{E}[\tilde{b}_i] + \sum_{i=1}^3 \mathbb{E}[\tilde{s}_i] \quad (7.72)$$

$$= 4\mu_{\text{bantam}} + 3\mu_{\text{sussex}} \quad (7.73)$$

$$= 33.5 \text{ eggs}. \quad (7.74)$$

- b) Let  $\tilde{w}$  represent the weather,  $\tilde{b}$  the eggs laid by a bantam hen and  $\tilde{s}$  the eggs laid by a Sussex hen. The conditional mean function of  $\tilde{b}$  given  $\tilde{w} = \text{cold}$  is

$$\mu_{\tilde{b}|\tilde{w}}(\text{cold}) = \sum_{i=0}^1 i \cdot p_{\tilde{b}|\tilde{w}}(i|\text{cold}) \quad (7.75)$$

$$= 0.25. \quad (7.76)$$

The conditional mean function of  $\tilde{s}$  given  $\tilde{w} = \text{cold}$  is

$$\mu_{\tilde{s}|\tilde{w}}(\text{cold}) = \sum_{i=0}^2 i \cdot p_{\tilde{s}|\tilde{w}}(i|\text{cold}) \quad (7.77)$$

$$= 0.5 + 2 \cdot 0.25 = 1. \quad (7.78)$$

The conditional mean function of  $\tilde{b}$  given  $\tilde{w} = \text{warm}$  is 5, and conditional mean function of  $\tilde{s}$  given  $\tilde{w} = \text{warm}$  is 4.5 as derived above. By iterated expectation,

$$\mathbb{E}[\tilde{b}] = \mathbb{E}\left[\mu_{\tilde{b}|\tilde{w}}(\tilde{w})\right] \quad (7.79)$$

$$= \mathbf{P}(\tilde{w} = \text{cold}) \mu_{\tilde{b}|\tilde{w}}(\text{cold}) + \mathbf{P}(\tilde{w} = \text{warm}) \mu_{\tilde{b}|\tilde{w}}(\text{warm}) \quad (7.80)$$

$$= 0.25(1 - \alpha) + 5\alpha = 0.25 + 4.75\alpha. \quad (7.81)$$

Similarly

$$\mathbb{E}[\tilde{s}] = \mathbb{E}[\mu_{\tilde{s}|\tilde{w}}(\tilde{w})] \quad (7.82)$$

$$= P(\tilde{w} = \text{cold})\mu_{\tilde{s}|\tilde{w}}(\text{cold}) + P(\tilde{w} = \text{warm})\mu_{\tilde{s}|\tilde{w}}(\text{warm}) \quad (7.83)$$

$$= 1 - \alpha + 4.5\alpha = 1 + 3.5\alpha. \quad (7.84)$$

The Sussex hen has a higher mean if

$$0.25 + 4.75\alpha \leq 1 + 3.5\alpha, \quad (7.85)$$

which is equivalent to  $\alpha \leq 0.75/1.25 = 0.6$ .

7.9 (Hotel questionnaire)

a) By iterated expectation

$$\mathbb{E}[s] = p_{\tilde{r}}(0)\mu_{\tilde{s}|\tilde{r}}(0) + p_{\tilde{r}}(1)\mu_{\tilde{s}|\tilde{r}}(1) \quad (7.86)$$

$$= 0.8 \cdot 4 + 0.2 \cdot 2 = 3.6. \quad (7.87)$$

b) Given the conditional independence assumptions, the conditional pmf of  $\tilde{s}$  given  $\tilde{q} = 1$  equals

$$p_{\tilde{s}|\tilde{q}}(s|1) = \frac{p_{\tilde{s},\tilde{q}}(s,1)}{p_{\tilde{q}}(1)} \quad (7.88)$$

$$= \frac{\sum_{r=0}^1 p_{\tilde{r},\tilde{s},\tilde{q}}(r,s,1)}{\sum_{r=0}^1 p_{\tilde{r},\tilde{q}}(r,1)} \quad (7.89)$$

$$= \frac{\sum_{r=0}^1 p_{\tilde{r}}(r)p_{\tilde{q}|\tilde{r}}(q|1)p_{\tilde{s}|\tilde{r}}(s|r)}{\sum_{r=0}^1 p_{\tilde{r}}(r)p_{\tilde{q}|\tilde{r}}(q|1)} \quad (7.90)$$

$$= \frac{0.8 \cdot 0.25 \cdot p_{\tilde{s}|\tilde{r}}(s|0) + 0.2 \cdot 0.9 \cdot p_{\tilde{s}|\tilde{r}}(s|1)}{0.8 \cdot 0.25 + 0.2 \cdot 0.9} \quad (7.91)$$

$$= \frac{0.2p_{\tilde{s}|\tilde{r}}(s|0) + 0.18p_{\tilde{s}|\tilde{r}}(s|1)}{0.38}. \quad (7.92)$$

The conditional mean given  $\tilde{q} = 1$  therefore equals

$$\mu_{\tilde{s}|\tilde{q}}(1) = \sum_{s=1}^5 s p_{\tilde{s}|\tilde{q}}(s|1) \quad (7.93)$$

$$= \sum_{s=1}^5 s \left( \frac{0.2p_{\tilde{s}|\tilde{r}}(s|0) + 0.18p_{\tilde{s}|\tilde{r}}(s|1)}{0.38} \right) \quad (7.94)$$

$$= \frac{0.2}{0.38} \sum_{s=1}^5 s p_{\tilde{s}|\tilde{r}}(s|0) + \frac{0.18}{0.38} \sum_{s=1}^5 s p_{\tilde{s}|\tilde{r}}(s|1) \quad (7.95)$$

$$= \frac{0.2}{0.38} \mu_{\tilde{s}|\tilde{r}}(0) + \frac{0.18}{0.38} \mu_{\tilde{s}|\tilde{r}}(1) \quad (7.96)$$

$$= \frac{0.2 \cdot 4}{0.38} + \frac{0.18 \cdot 2}{0.38} \quad (7.97)$$

$$= 3.05. \quad (7.98)$$

The observed mean is 0.55 less than the true mean.

7.10 (Another basketball player) Let  $\tilde{x}_1$ ,  $\tilde{x}_2$  and  $\tilde{x}_3$  represent the number of free throws, 2-point shots, 3-point shots, and free throws respectively. We denote the number of made

free throws, 2-point shots, and 3-point shots by  $\tilde{m}_1$ ,  $\tilde{m}_2$  and  $\tilde{m}_3$  respectively, and set  $p_1 = 0.8$ ,  $p_2 = 0.5$  and  $p_3 = 0.3$ . The number of points  $\tilde{y}$  is equal to

$$\tilde{y} = \tilde{m}_1 + 2\tilde{m}_2 + 3\tilde{m}_3. \quad (7.99)$$

By linearity of expectation (which does not require independence)

$$\mathbb{E}[\tilde{y}] = \mathbb{E}[\tilde{m}_1] + 2\mathbb{E}[\tilde{m}_2] + 3\mathbb{E}[\tilde{m}_3]. \quad (7.100)$$

We now apply iterated expectation to compute each term. For this we require the conditional means. Conditioned on  $\tilde{x}_i = x$ ,  $\tilde{m}_i = \sum_{j=1}^x \tilde{s}_i$  where  $\tilde{s}_i$  is a Bernoulli random variable that equals one if the  $i$ th shot is made. By linearity of expectation

$$\mu_{\tilde{m}_i | \tilde{x}_i}(x) = \mathbb{E} \left[ \sum_{j=1}^x \tilde{s}_i \right] \quad (7.101)$$

$$= \sum_{j=1}^x \mathbb{E}[\tilde{s}_i] \quad (7.102)$$

$$= xp_i. \quad (7.103)$$

By iterated expectation,

$$\mathbb{E}[\tilde{m}_i] = \mathbb{E}[\mu_{\tilde{m}_i | \tilde{x}_i}(\tilde{x}_i)] \quad (7.104)$$

$$= \mathbb{E}[p_i \tilde{x}_i] \quad (7.105)$$

$$= p_i \mathbb{E}[\tilde{x}_i] \quad (7.106)$$

$$= 2p_i, \quad (7.107)$$

where we have used the fact that

$$\mathbb{E}[\tilde{x}_i] = \sum_{j=0}^4 jp_x(j) \quad (7.108)$$

$$= \frac{1 + 2 + 3 + 4}{5} \quad (7.109)$$

$$= 2. \quad (7.110)$$

We conclude that

$$\mathbb{E}[\tilde{y}] = \mathbb{E}[\tilde{m}_1] + 2\mathbb{E}[\tilde{m}_2] + 3\mathbb{E}[\tilde{m}_3] \quad (7.111)$$

$$= 2(p_1 + 2 \cdot p_2 + 3 \cdot p_3) \quad (7.112)$$

$$= 5.4 \text{ points.} \quad (7.113)$$

7.11 (Life expectancy) We represent the age of death as a random variable  $\tilde{x}$  and define a Bernoulli random variable  $\tilde{d}$  with parameter 0.25 to represent whether a person died ( $\tilde{d} = 1$ ) in the first year or not ( $\tilde{d} = 0$ ). We are interested in the conditional mean function  $\mu_{\tilde{x} | \tilde{d}}$  evaluated at 0. If  $\tilde{d} = 1$ , then  $\tilde{x} = 0$ , so  $\mu_{\tilde{x} | \tilde{d}}(1) = 0$ . By iterated expectation,

$$37.5 = \mathbb{E}[\tilde{x}] = \mathbb{E}[\mu_{\tilde{x} | \tilde{d}}(\tilde{d})] \quad (7.114)$$

$$= 0 \cdot \mathbb{P}(\tilde{d} = 1) + \mu_{\tilde{x} | \tilde{d}}(0)\mathbb{P}(\tilde{d} = 0) \quad (7.115)$$

$$= 0.75\mu_{\tilde{x} | \tilde{d}}(0). \quad (7.116)$$

We conclude that  $\mu_{\tilde{x} | \tilde{d}}(0) = 50$ .

7.12 (Buckets) Let 0 indicate the small bucket and 1 the large one. The mean of  $\tilde{s}_1$  equals

$$\mathbb{E}[\tilde{s}_1] = \sum_{b=0}^1 \int_{s=-\infty}^{\infty} s p_{\tilde{b}}(b) f_{\tilde{s}_1|\tilde{b}}(s|b) ds \quad (7.117)$$

$$= p_{\tilde{b}}(0) \int_{s=0}^1 s ds + p_{\tilde{b}}(1) \int_{s=0}^2 \frac{s}{2} ds \quad (7.118)$$

$$= \frac{3}{4} \cdot \frac{1}{2} + \frac{1}{4} \cdot 1 \quad (7.119)$$

$$= \frac{5}{8}. \quad (7.120)$$

The mean square equals

$$\mathbb{E}[\tilde{s}_1^2] = \sum_{b=0}^1 \int_{s=-\infty}^{\infty} s^2 p_{\tilde{b}}(b) f_{\tilde{s}_1|\tilde{b}}(s|b) ds \quad (7.121)$$

$$= p_{\tilde{b}}(0) \int_{s=0}^1 s^2 ds + p_{\tilde{b}}(1) \int_{s=0}^2 \frac{s^2}{2} ds \quad (7.122)$$

$$= \frac{3}{4} \cdot \frac{1}{3} + \frac{1}{4} \cdot \frac{2^3}{3} \quad (7.123)$$

$$= \frac{11}{12}. \quad (7.124)$$

The variance equals

$$\text{Var}[\tilde{s}_1] = \mathbb{E}[\tilde{s}_1^2] - \mathbb{E}[\tilde{s}_1]^2 \quad (7.125)$$

$$= 0.526. \quad (7.126)$$

7.13 (Potatoes)

- a) The minimum MSE estimate is  $\mu_{\tilde{x}|\tilde{b},\tilde{w}}(1,1) = 50$  tons (from the corresponding graph).
- b) Since we know the conditional distribution of  $\tilde{x}$  given the weather and the presence of the beetle, we express the desired conditional mean function including  $\tilde{b}$ :

$$\mu_{\tilde{x}|\tilde{w}}(0) = \sum_{b=0}^1 p_{\tilde{b}|\tilde{w}}(b|0) \int_{x=-\infty}^{\infty} x f_{\tilde{x}|\tilde{w},\tilde{b}} dx \quad (7.127)$$

$$= \sum_{b=0}^1 p_{\tilde{b}|\tilde{w}}(b|0) \mu_{\tilde{x}|\tilde{w},\tilde{b}}(0,b). \quad (7.128)$$

The conditional distribution of  $\tilde{x}$  given  $\tilde{b}$  and  $\tilde{w}$  is Gaussian, so their means are equal to their mode, right at the center of symmetry of the pdf. From the graphs we have  $\mu_{\tilde{x}|\tilde{w},\tilde{b}}(0,0) = 40$  and  $\mu_{\tilde{x}|\tilde{w},\tilde{b}}(0,1) = 30$ . In addition, from the table  $p_{\tilde{b}|\tilde{w}}(1|0) = 20/30 = 2/3$  and  $p_{\tilde{b}|\tilde{w}}(0|0) = 10/30 = 1/3$ . We conclude that

$$\mu_{\tilde{x}|\tilde{w}}(0) = \frac{40}{3} + \frac{2 \cdot 30}{3} \quad (7.129)$$

$$= \frac{100}{3}. \quad (7.130)$$

7.14 (Wolf)

a)

$$\mu_{\tilde{x}[2]|\tilde{s}}(0) = \int_{a=-5}^5 \int_{b=-5}^5 \frac{b}{100} da db \quad (7.131)$$

$$= 0, \quad (7.132)$$

$$\mu_{\tilde{x}[2]|\tilde{s}}(1) = \int_{a=-5}^5 \int_{b=-5}^0 \frac{b}{50} da db \quad (7.133)$$

$$= \frac{10}{50} \left( 0 - \frac{(-5)^2}{2} \right) \quad (7.134)$$

$$= -2.5. \quad (7.135)$$

By iterated expectation,

$$\mathbb{E}[\tilde{x}[2]] = \mathbb{E}[\mu_{\tilde{x}[2]|\tilde{s}}(\tilde{s})] \quad (7.136)$$

$$= 0.25 \cdot (-2.5) \quad (7.137)$$

$$= -0.625. \quad (7.138)$$

- b) The minimum MSE estimator of  $\tilde{x}[1]$  given  $\tilde{x}[2]$  is the conditional mean function of  $\tilde{x}[1]$  given  $\tilde{x}[2]$ . We first need to compute the conditional pdf of  $\tilde{x}[1]$  given  $\tilde{x}[2]$ .

If  $0 \leq x_2 \leq 5$ ,

$$f_{\tilde{x}[1]|\tilde{x}[2]}(x_1 | x_2) = \frac{f_{\tilde{x}}(x_1, x_2)}{f_{\tilde{x}[2]}(x_2)} \quad (7.139)$$

$$= \frac{\sum_{s=0}^1 p_{\tilde{s}}(s) f_{\tilde{x}|\tilde{s}}(x_1, x_2 | s)}{\int_{x_1=-\infty}^{\infty} \sum_{s=0}^1 p_{\tilde{s}}(s) f_{\tilde{x}|\tilde{s}}(x_1, x_2 | s) dx_1} \quad (7.140)$$

$$= \frac{0.75 \cdot \frac{1}{100}}{\int_{x_1=-5}^5 0.75 \cdot \frac{1}{100} dx_1} \quad (7.141)$$

$$= \frac{1}{10}, \quad (7.142)$$

for  $-5 \leq x_1 \leq 5$  and 0 otherwise.

If  $-5 \leq x_2 \leq 0$ ,

$$f_{\tilde{x}[1]|\tilde{x}[2]}(x_1 | x_2) = \frac{f_{\tilde{x}}(x_1, x_2)}{f_{\tilde{x}[2]}(x_2)} \quad (7.143)$$

$$= \frac{\sum_{s=0}^1 p_{\tilde{s}}(s) f_{\tilde{x}|\tilde{s}}(x_1, x_2 | s)}{\int_{x_1=-\infty}^{\infty} \sum_{s=0}^1 p_{\tilde{s}}(s) f_{\tilde{x}|\tilde{s}}(x_1, x_2 | s) dx_1} \quad (7.144)$$

$$= \frac{0.75 \cdot \frac{1}{100} + 0.25 \cdot \frac{1}{50}}{\int_{x_1=-5}^5 (0.75 \cdot \frac{1}{100} + 0.25 \cdot \frac{1}{50}) dx_1} \quad (7.145)$$

$$= \frac{1}{10}, \quad (7.146)$$

for  $-5 \leq x_1 \leq 5$  and 0 otherwise.

It turns out that the two entries of  $\tilde{x}$  are independent (you can check that the marginal pdf is also equal to  $1/10$ ). As a result, the best estimate is just the mean of  $\tilde{x}[1]$ , no



matter what the value of  $\tilde{x}[2]$  is. For any  $-5 \leq x_2 \leq 5$ ,

$$\mu_{\tilde{x}[1] | \tilde{x}[2]}(x_2) = \int_{x_1=-\infty}^{\infty} x_1 f_{\tilde{x}[1] | \tilde{x}[2]}(x_1 | x_2) dx_1 \quad (7.147)$$

$$= \int_{x_1=-5}^5 \frac{x_1}{10} dx_1 \quad (7.148)$$

$$= 0. \quad (7.149)$$

### 7.15 (Law of conditional variance)

a) Recall that

$$p_{\tilde{b} | \tilde{a}}(1 | 1) = \frac{1}{7}, \quad p_{\tilde{b} | \tilde{a}}(2 | 1) = \frac{4}{7}, \quad p_{\tilde{b} | \tilde{a}}(3 | 1) = \frac{2}{7} \quad (7.150)$$

$$p_{\tilde{b} | \tilde{a}}(1 | 2) = \frac{2}{7}, \quad p_{\tilde{b} | \tilde{a}}(2 | 2) = \frac{1}{7}, \quad p_{\tilde{b} | \tilde{a}}(3 | 2) = \frac{4}{7} \quad (7.151)$$

$$p_{\tilde{b} | \tilde{a}}(1 | 3) = \frac{1}{3}, \quad p_{\tilde{b} | \tilde{a}}(2 | 3) = \frac{1}{3}, \quad p_{\tilde{b} | \tilde{a}}(3 | 3) = \frac{1}{3}, \quad (7.152)$$

and the conditional mean function of  $\tilde{b}$  given  $\tilde{a}$  is

$$\mu_{\tilde{b} | \tilde{a}}(1) = \frac{15}{7}, \quad (7.153)$$

$$\mu_{\tilde{b} | \tilde{a}}(2) = \frac{16}{7}, \quad (7.154)$$

$$\mu_{\tilde{b} | \tilde{a}}(3) = 2. \quad (7.155)$$

We have

$$\nu_{\tilde{b} | \tilde{a}}(1) = \sum_{b=1}^3 \left(b - \frac{15}{7}\right)^2 p_{\tilde{b} | \tilde{a}}(b | 1) \quad (7.156)$$

$$= \frac{20}{49}, \quad (7.157)$$

$$\nu_{\tilde{b} | \tilde{a}}(2) = \sum_{b=1}^3 \left(b - \frac{16}{7}\right)^2 p_{\tilde{b} | \tilde{a}}(b | 2) \quad (7.158)$$

$$= \frac{38}{49}, \quad (7.159)$$

$$\nu_{\tilde{b} | \tilde{a}}(3) = \sum_{b=1}^3 (b - 2)^2 p_{\tilde{b} | \tilde{a}}(b | 3) \quad (7.160)$$

$$= \frac{2}{3}. \quad (7.161)$$

b)

$$p_{\nu_{\tilde{b} | \tilde{a}}(\tilde{a})}\left(\frac{20}{49}\right) = P\left(\mu_{\tilde{b} | \tilde{a}}(\tilde{a}) = \frac{20}{49}\right) \quad (7.162)$$

$$= P(\tilde{a} = 1) \quad (7.163)$$

$$= 0.35. \quad (7.164)$$

Similarly,

$$p_{\mu_{\tilde{b} | \tilde{a}}(\tilde{a})}\left(\frac{38}{49}\right) = 0.35, \quad (7.165)$$

$$p_{\mu_{\tilde{b} | \tilde{a}}(\tilde{a})}\left(\frac{2}{3}\right) = 0.3. \quad (7.166)$$

- c) We denote the conditional mean square function of  $\tilde{b}$  given  $\tilde{a} = a$  by  $\mu_{\tilde{b}^2|\tilde{a}}(a)$ . By iterated expectation  $\mathbb{E}[\mu_{\tilde{b}^2|\tilde{a}}(\tilde{a})] = \mathbb{E}[\tilde{b}^2]$ . We have

$$\text{Var}[\mu_{\tilde{b}|\tilde{a}}(\tilde{a})] = \mathbb{E}[\mu_{\tilde{b}|\tilde{a}}(\tilde{a})^2] - \mathbb{E}[\mu_{\tilde{b}|\tilde{a}}(\tilde{a})]^2 \quad (7.167)$$

$$= \mathbb{E}[\mu_{\tilde{b}|\tilde{a}}(\tilde{a})^2] - \mathbb{E}[\tilde{b}]^2, \quad (7.168)$$

$$\mathbb{E}[\nu_{\tilde{b}|\tilde{a}}(\tilde{a})] = \mathbb{E}[\mu_{\tilde{b}^2|\tilde{a}}(\tilde{a}) - \mu_{\tilde{b}|\tilde{a}}(\tilde{a})^2] \quad (7.169)$$

$$= \mathbb{E}[\tilde{b}^2] - \mathbb{E}[\mu_{\tilde{b}|\tilde{a}}(\tilde{a})^2]. \quad (7.170)$$

Consequently,

$$\mathbb{E}[\nu_{\tilde{b}|\tilde{a}}(\tilde{a})] + \text{Var}[\mu_{\tilde{b}|\tilde{a}}(\tilde{a})] = \mathbb{E}[\tilde{b}^2] - \mathbb{E}[\tilde{b}]^2 \quad (7.171)$$

$$= \text{Var}[\tilde{b}]. \quad (7.172)$$

- d) The pmf of the conditional mean is

$$p_{\mu_{\tilde{b}|\tilde{a}}(\tilde{a})}\left(\frac{15}{7}\right) = 0.35, \quad (7.173)$$

$$p_{\mu_{\tilde{b}|\tilde{a}}(\tilde{a})}\left(\frac{16}{7}\right) = 0.35, \quad (7.174)$$

$$p_{\mu_{\tilde{b}|\tilde{a}}(\tilde{a})}(2) = 0.3. \quad (7.175)$$

The mean square of  $\mu_{\tilde{b}|\tilde{a}}(\tilde{a})$  therefore equals

$$\mathbb{E}[\mu_{\tilde{b}|\tilde{a}}(\tilde{a})^2] = \left(\frac{15}{7}\right)^2 \cdot 0.35 + \left(\frac{16}{7}\right)^2 \cdot 0.35 + 2^2 \cdot 0.3 \quad (7.176)$$

$$= 4.636, \quad (7.177)$$

and its mean is 2.15, so the variance of the conditional mean equals

$$\text{Var}[\mu_{\tilde{b}|\tilde{a}}(\tilde{a})] = 4.636 - 2.15^2 = 0.0132. \quad (7.178)$$

The mean of the conditional variance is

$$\mathbb{E}[\nu_{\tilde{b}|\tilde{a}}(\tilde{a})] = \sum_{a=1}^3 \nu_{\tilde{b}|\tilde{a}}(a) p_{\tilde{a}}(a) \quad (7.179)$$

$$= 0.35 \cdot \frac{20}{49} + 0.35 \cdot \frac{38}{49} + 0.3 \cdot \frac{2}{3} \quad (7.180)$$

$$= 0.6143. \quad (7.181)$$

By the law of conditional variance the variance of  $\tilde{b}$  should equal  $0.0132 + 0.6143 = 0.6275$ . We verify this using the marginal pmf of  $\tilde{b}$ . The mean square of  $\tilde{b}$  equals

$$\mathbb{E}[\tilde{b}^2] = \sum_{b=1}^3 b^2 p_{\tilde{b}}(b) \quad (7.182)$$

$$= 1 \cdot 0.25 + 4 \cdot 0.35 + 9 \cdot 0.4 \quad (7.183)$$

$$= 5.25, \quad (7.184)$$

so the variance equals  $5.25 - 2.15^2 = 0.6275$ .

7.16 (Faulty database)

- a) Let  $\tilde{c}$  be a random variable that represents whether the entry is corrupted, we have

$$\tilde{b} = \begin{cases} -x & \text{if } \tilde{c} = -1, \\ \tilde{a} & \text{if } \tilde{c} = 0, \\ x & \text{if } \tilde{c} = 1. \end{cases} \quad (7.185)$$

The conditional mean function of  $\tilde{b}$  given  $\tilde{c} = c$  equals

$$\mu_{\tilde{b}|\tilde{c}}(c) = \begin{cases} -x & \text{if } c = -1, \\ 0 & \text{if } c = 0, \\ x & \text{if } c = 1, \end{cases} \quad (7.186)$$

so by iterated expectation

$$\mathbb{E}[\tilde{b}] = \mathbb{E}[\mu_{\tilde{b}|\tilde{c}}(\tilde{c})] \quad (7.187)$$

$$= 0.1(-x) + 0.1x = 0. \quad (7.188)$$

The conditional mean function of  $\tilde{b}^2$  given  $\tilde{c} = c$  equals

$$\mu_{\tilde{b}^2|\tilde{c}}(c) = \begin{cases} x^2 & \text{if } c = -1, \\ \sigma^2 & \text{if } c = 0, \\ x^2 & \text{if } c = 1, \end{cases} \quad (7.189)$$

so by iterated expectation

$$\mathbb{E}[\tilde{b}^2] = \mathbb{E}[\mu_{\tilde{b}^2|\tilde{c}}(\tilde{c})] \quad (7.190)$$

$$= 0.2x^2 + 0.8\sigma^2. \quad (7.191)$$

We conclude that  $\text{Var}[\tilde{b}] = \mathbb{E}[\tilde{b}^2] - \mathbb{E}[\tilde{b}]^2 = 0.2x^2 + 0.8\sigma^2$ .

- b) The minimum mean-squared error estimator of  $\tilde{a}$  given  $\tilde{b} = b$  is the conditional mean  $\mu_{\tilde{a}|\tilde{b}}$  of  $\tilde{a}$  given  $\tilde{b} = b$ . Conditioned on  $\tilde{b} = -x$  or  $\tilde{b} = x$ , the distribution of  $\tilde{a}$  is Gaussian with zero mean and variance  $\sigma^2$  because the database failure is independent from the content of the database entry. Therefore  $\mu_{\tilde{a}|\tilde{b}}(-x) = \mu_{\tilde{a}|\tilde{b}}(x) = 0$ . Conditioned on  $\tilde{b} = b \notin \{-x, x\}$ ,  $\tilde{a}$  is equal to  $b$  with probability one, so  $\mu_{\tilde{a}|\tilde{b}}(b) = b$  for all  $b \notin \{-x, x\}$ .

#### 7.17 (Coffee and life expectancy)

- a) Let  $\tilde{\ell}$ ,  $\tilde{c}$  and  $\tilde{s}$  represent life span, smoking and coffee respectively. Since have access to the conditional mean given both smoking and coffee  $\mu_{\tilde{\ell}|\tilde{c},\tilde{s}}$ ,

$$\mu_{\tilde{\ell}|\tilde{c}}(1) = \int_{\ell=-\infty}^{\infty} \ell f_{\tilde{\ell}|\tilde{c}}(\ell|1) \, d\ell \quad (7.192)$$

$$= \int_{\ell=-\infty}^{\infty} \ell \sum_{s=0}^1 p_{\tilde{s}|\tilde{c}}(s|1) f_{\tilde{\ell}|\tilde{c}}(\ell|1, s) \, d\ell \quad (7.193)$$

$$= \sum_{s=0}^1 p_{\tilde{s}|\tilde{c}}(s|1) \mu_{\tilde{\ell}|\tilde{c},\tilde{s}}(1, s) \quad (7.194)$$

$$= \frac{80}{120} \cdot 68 + \frac{40}{120} \cdot 82 \quad (7.195)$$

$$= 72.7. \quad (7.196)$$

Similarly,

$$\mu_{\tilde{\ell}|\tilde{c}}(0) = \sum_{s=0}^1 p_{\tilde{s}}(s) \mu_{\tilde{\ell}|\tilde{c},\tilde{s}}(0, s) \quad (7.197)$$

$$= \frac{20}{80} \cdot 70 + \frac{60}{80} \cdot 80 \quad (7.198)$$

$$= 77.5. \quad (7.199)$$

The observed ATE is therefore  $72.7 - 77.5 = -4.8$  years.

- b) From the tables  $p_{\tilde{s}}(0) = 100/200 = 1/2$  and  $p_{\tilde{s}}(1) = 1/2$ . By the formula for the adjusted ATE given the conditional independence assumption,

$$\text{ATE} = \sum_{s=0}^1 p_{\tilde{s}}(s) \mu_{\tilde{\ell}|\tilde{c},\tilde{s}}(1, s) - \sum_{s=0}^1 p_{\tilde{s}}(s) \mu_{\tilde{\ell}|\tilde{c},\tilde{s}}(0, s) \quad (7.200)$$

$$= \frac{1}{2} \cdot 68 + \frac{1}{2} \cdot 82 - \frac{1}{2} \cdot 70 + \frac{1}{2} \cdot 80 \quad (7.201)$$

$$= 0. \quad (7.202)$$

The adjusted ATE is zero, which indicates that coffee has no effect on life span.

#### 7.18 (Weight-loss supplement)

- a) The unadjusted ATE equals

unadjusted ATE

$$\begin{aligned} &= \frac{0.5 - 4.6 - 7.5 + 1.1 - 3.4 - 1.4 + 3.2 - 1.9 - 3.6 + 2.4 + 1.2 + 4.2 - 2.5 - 3.5 + 4.7}{15} \\ &\quad - \frac{1.6 - 6.7 + 2.3 - 4.5 + 3.4 - 0.6 + 2.3 + 1.5 + 1.7 + 2.6 + 0.5 + 1.2 - 3.6 + 2.1 - 1.2 + 0.4}{16} \\ &= -0.93. \end{aligned} \quad (7.203)$$

- b) We model the treatment as a random variable  $\tilde{t}$ , the weight loss as a random variable  $\tilde{w}$  and the exercise as a random variable  $\tilde{x}$ . In order to adjust for the effect of exercise, we compute

$$p_{\tilde{x}}(1) = \frac{13}{31} = 0.42, \quad (7.204)$$

$$\mu_{\tilde{w}|\tilde{t},\tilde{x}}(0, 0) = \frac{1.6 + 2.3 - 0.6 + 2.3 + 1.5 + 1.7 + 2.6 + 0.5 + 1.2 + 2.1 - 1.2 + 0.4}{12} \quad (7.205)$$

$$= 1.2, \quad (7.206)$$

$$\mu_{\tilde{w}|\tilde{t},\tilde{x}}(0, 1) = \frac{-6.7 - 4.5 + 3.4 - 3.6}{4} = -2.85, \quad (7.207)$$

$$\mu_{\tilde{w}|\tilde{t},\tilde{x}}(1, 0) = \frac{0.5 + 1.1 - 1.9 + 2.4 + 1.2 + 4.7}{6} = 1.33, \quad (7.208)$$

$$\mu_{\tilde{w}|\tilde{t},\tilde{x}}(1, 1) = \frac{-4.6 - 7.5 - 3.4 - 1.4 + 3.2 - 3.6 + 4.2 - 2.5 - 3.5}{9} = -2.12, \quad (7.209)$$

The adjusted ATE equals

$$\text{adjusted ATE} = \sum_{x=0}^1 p_{\tilde{x}}(x) \mu_{\tilde{w}|\tilde{t},\tilde{x}}(1, x) - \sum_{x=0}^1 p_{\tilde{x}}(x) \mu_{\tilde{w}|\tilde{t},\tilde{x}}(0, x) \quad (7.210)$$

$$= (0.58 \cdot 1.33 - 0.42 \cdot 2.12) - (0.58 \cdot 1.2 - 0.42 \cdot 2.85) = 0.382. \quad (7.211)$$

This approximates the true ATE if the potential outcomes of the weight gain/loss for treatment and control are conditionally independent of the treatment given the exercise.

### 7.19 (Chocolate bar)

- a) We model taking the chocolate bar as a treatment  $\tilde{t}$ , the cholesterol as a random variable  $\tilde{y}$  and the sex of the participants as another random variable  $\tilde{s}$ . From the information in the problem,

$$\mu_{\tilde{y}|\tilde{t},\tilde{s}}(1, \text{man}) = 150, \quad (7.212)$$

$$\mu_{\tilde{y}|\tilde{t},\tilde{s}}(0, \text{man}) = 140, \quad (7.213)$$

$$\mu_{\tilde{y}|\tilde{t},\tilde{s}}(1, \text{woman}) = 130, \quad (7.214)$$

$$\mu_{\tilde{y}|\tilde{t},\tilde{s}}(0, \text{woman}) = 120. \quad (7.215)$$

and

$$p_{\tilde{s}|\tilde{t}}(\text{man} | 1) = \frac{p_{\tilde{s}}(\text{man})p_{\tilde{t}|\tilde{s}}(1 | \text{man})}{p_{\tilde{s}}(\text{man})p_{\tilde{t}|\tilde{s}}(1 | \text{man}) + p_{\tilde{s}}(\text{woman})p_{\tilde{t}|\tilde{s}}(1 | \text{woman})} \quad (7.216)$$

$$= \frac{\alpha}{\alpha + 0.5}, \quad (7.217)$$

$$p_{\tilde{s}|\tilde{t}}(\text{man} | 0) = \frac{1 - \alpha}{1 - \alpha + 0.5} \quad (7.218)$$

The average cholesterol in the treatment group is

$$\mu_{\tilde{y}|\tilde{t}}(1) = \sum_{s \in \{\text{man}, \text{woman}\}} p_{\tilde{s}|\tilde{t}}(s | 1) \int_{y=-\infty}^{\infty} y f_{\tilde{y}|\tilde{t},\tilde{s}}(y | 1, s) \, dy \quad (7.219)$$

$$= p_{\tilde{s}|\tilde{t}}(\text{man} | 1) \mu_{\tilde{y}|\tilde{t},\tilde{s}}(1, \text{man}) + p_{\tilde{s}|\tilde{t}}(\text{woman} | 1) \mu_{\tilde{y}|\tilde{t},\tilde{s}}(1, \text{woman}) \quad (7.220)$$

$$= \frac{150\alpha}{\alpha + 0.5} + \frac{0.5 \cdot 130}{\alpha + 0.5}. \quad (7.221)$$

By the same argument, the average cholesterol in the control group is

$$\mu_{\tilde{y}|\tilde{t}}(0) = p_{\tilde{s}|\tilde{t}}(\text{man} | 0) \mu_{\tilde{y}|\tilde{t},\tilde{s}}(0, \text{man}) + p_{\tilde{s}|\tilde{t}}(\text{woman} | 0) \mu_{\tilde{y}|\tilde{t},\tilde{s}}(0, \text{woman}) \quad (7.222)$$

$$= \frac{(1 - \alpha)140}{1.5 - \alpha} + \frac{0.5 \cdot 120}{1.5 - \alpha}. \quad (7.223)$$

The observed ATE equals

$$\text{observed ATE}(\alpha) = \frac{150\alpha + 65}{\alpha + 0.5} - \frac{200 - 140\alpha}{1.5 - \alpha}. \quad (7.224)$$

b)

$$\text{observed ATE}(0.05) = -1.29, \quad (7.225)$$

$$\text{observed ATE}(0.5) = 10, \quad (7.226)$$

$$\text{observed ATE}(0.95) = 21.3. \quad (7.227)$$

## Correlation

### Exercises

- 8.1 (Three variables) Yes, it is possible. Let  $\tilde{a}$ ,  $\tilde{b}$  and  $\tilde{c}$  be random variables with zero mean and unit variance, which are pairwise uncorrelated. We define:

$$\tilde{w} := \tilde{a} + \alpha\tilde{b}, \quad (8.1)$$

$$\tilde{y} := \tilde{b} + \beta\tilde{c}, \quad (8.2)$$

$$\tilde{z} := \gamma\tilde{a} + \tilde{c}, \quad (8.3)$$

where  $\alpha$ ,  $\beta$  and  $\gamma$  are real-valued constants. By linearity of expectation

$$\mathbb{E}[\tilde{w}] = \mathbb{E}[\tilde{a}] + \alpha\mathbb{E}[\tilde{b}] = 0, \quad (8.4)$$

$$\mathbb{E}[\tilde{y}] = \mathbb{E}[\tilde{b}] + \beta\mathbb{E}[\tilde{c}] = 0, \quad (8.5)$$

$$\mathbb{E}[\tilde{z}] = \gamma\mathbb{E}[\tilde{a}] + \mathbb{E}[\tilde{c}] = 0. \quad (8.6)$$

Since  $\tilde{a}$ ,  $\tilde{b}$  and  $\tilde{c}$  have zero mean and are pairwise uncorrelated,

$$\mathbb{E}[\tilde{a}\tilde{b}] = \mathbb{E}[\tilde{a}\tilde{c}] = \mathbb{E}[\tilde{b}\tilde{c}] = 0. \quad (8.7)$$

Consequently, again by linearity of expectation,

$$\mathbb{E}[\tilde{w}\tilde{y}] = \mathbb{E}[\tilde{a}\tilde{b}] + \beta\mathbb{E}[\tilde{a}\tilde{c}] + \alpha\mathbb{E}[\tilde{b}^2] + \alpha\beta\mathbb{E}[\tilde{b}\tilde{c}] \quad (8.8)$$

$$= \alpha, \quad (8.9)$$

$$\mathbb{E}[\tilde{w}\tilde{z}] = \gamma\mathbb{E}[\tilde{a}^2] + \mathbb{E}[\tilde{a}\tilde{c}] + \gamma\alpha\mathbb{E}[\tilde{a}\tilde{b}] + \alpha\mathbb{E}[\tilde{b}\tilde{c}] \quad (8.10)$$

$$= \gamma, \quad (8.11)$$

$$\mathbb{E}[\tilde{y}\tilde{z}] = \gamma\mathbb{E}[\tilde{a}\tilde{b}] + \mathbb{E}[\tilde{b}\tilde{c}] + \gamma\beta\mathbb{E}[\tilde{a}\tilde{c}] + \beta\mathbb{E}[\tilde{c}^2] \quad (8.12)$$

$$= \beta. \quad (8.13)$$

If we choose  $\alpha$  and  $\beta$  to be positive and  $\gamma$  negative, then  $\tilde{w}$  and  $\tilde{y}$  are positively correlated, and so are  $\tilde{y}$  and  $\tilde{z}$ , but  $\tilde{w}$  and  $\tilde{z}$  are negatively correlated.

- 8.2 (Properties of covariance)

a) By the definition of covariance and linearity of expectation

$$\text{Cov}[\beta\tilde{a} + \alpha, \tilde{b}] := \mathbb{E}[(\beta\tilde{a} + \alpha - \mathbb{E}[\beta\tilde{a} + \alpha])(\tilde{b} - \mathbb{E}[\tilde{b}])] \quad (8.14)$$

$$= \beta\mathbb{E}[(\tilde{a} - \mathbb{E}[\tilde{a}])(\tilde{b} - \mathbb{E}[\tilde{b}])] \quad (8.15)$$

$$= \beta\text{Cov}[\tilde{a}, \tilde{b}]. \quad (8.16)$$

- b) Since  $-1 \leq \rho_{\tilde{a}, \tilde{b}} \leq 1$  and  $\rho_{\tilde{a}, \tilde{b}} = \frac{\text{Cov}[\tilde{a}, \tilde{b}]}{\sqrt{\text{Var}[\tilde{a}] \text{Var}[\tilde{b}]}}$ , it follows that

$$-\sqrt{\text{Var}[\tilde{a}] \text{Var}[\tilde{b}]} \leq \text{Cov}[\tilde{a}, \tilde{b}] \leq \sqrt{\text{Var}[\tilde{a}] \text{Var}[\tilde{b}]} \quad (8.17)$$

- c) In the proof of Theorem 8.34 we show that the covariance is symmetric and that for any random variables  $\tilde{x}_1$ ,  $\tilde{x}_2$  and  $\tilde{y}$ ,  $\text{Cov}[\tilde{x}_1 + \tilde{x}_2, \tilde{y}] = \text{Cov}[\tilde{x}_1, \tilde{y}] + \text{Cov}[\tilde{x}_2, \tilde{y}]$ . Consequently,

$$\text{Cov}[\tilde{a} + \tilde{b}, \tilde{a} - \tilde{b}] = \text{Cov}[\tilde{a}, \tilde{a} - \tilde{b}] + \text{Cov}[\tilde{b}, \tilde{a} - \tilde{b}] \quad (8.18)$$

$$= \text{Cov}[\tilde{a}, \tilde{a}] - \text{Cov}[\tilde{a}, \tilde{b}] + \text{Cov}[\tilde{b}, \tilde{a}] - \text{Cov}[\tilde{b}, \tilde{b}] \quad (8.19)$$

$$= \text{Var}[\tilde{a}] - \text{Var}[\tilde{b}], \quad (8.20)$$

because  $\text{Cov}[\tilde{a}, \tilde{a}] = \mathbb{E}[(\tilde{a} - \mathbb{E}[\tilde{a}])^2] = \text{Var}[\tilde{a}]$ .

- 8.3 (Dunning-Kruger effect) Since the means of all variables are zero,

$$\text{Cov}[\tilde{t}, \tilde{d}] = \mathbb{E}[\tilde{t}(\tilde{s} - \tilde{t})] \quad (8.21)$$

$$= \mathbb{E}[\tilde{t}, \tilde{s}] - \mathbb{E}[\tilde{t}^2] \quad (8.22)$$

$$= -\text{Var}[\tilde{t}] = 1, \quad (8.23)$$

and by the uncorrelation assumption

$$\text{Var}[\tilde{d}] = \text{Var}[\tilde{s} - \tilde{t}] \quad (8.24)$$

$$= \text{Var}[\tilde{s}] + \text{Var}[-\tilde{t}] \quad (8.25)$$

$$= \text{Var}[\tilde{s}] + \text{Var}[\tilde{t}] \quad (8.26)$$

$$= 2. \quad (8.27)$$

Consequently,

$$\rho_{\tilde{t}, \tilde{d}} = \frac{\text{Cov}[\tilde{t}, \tilde{d}]}{\sqrt{\text{Var}[\tilde{t}] \text{Var}[\tilde{d}]}} \quad (8.28)$$

$$= \frac{-1}{\sqrt{2}} = -0.707. \quad (8.29)$$

The study finds a strong negative correlation, even though the true and self-evaluated competences are uncorrelated!

- 8.4 (Cross) The conditional mean  $\mu_{\tilde{y}|\tilde{x}}(x)$  is zero for all values of  $x$  with nonzero density, because all the conditional distributions are uniform and centered at zero. By iterated expectation  $\mathbb{E}[\tilde{y}] = \mathbb{E}[\mu_{\tilde{y}|\tilde{x}}(\tilde{x})] = 0$ . Also,

$$\mu_{\tilde{x}\tilde{y}|\tilde{x}}(x) = x\mu_{\tilde{y}|\tilde{x}}(x) \quad (8.30)$$

$$= 0 \quad (8.31)$$

for all  $x$  with nonzero density. By iterated expectation  $E[\tilde{x}\tilde{y}] = \mathbb{E}[\mu_{\tilde{x}\tilde{y}|\tilde{x}}(\tilde{x})] = 0$ . The variables are uncorrelated.

- 8.5 (Rufus) The mean of  $\tilde{x}$  equals

$$\mathbb{E}[\tilde{x}] = \int_{y=-50}^{50} \left( \int_{x=-50}^{-20} cx \, dx + \int_{x=20}^{50} cx \, dx \right) dy \quad (8.32)$$

$$+ \int_{x=-20}^{20} cx \, dx \left( \int_{y=-50}^{-20} dy + \int_{y=20}^{50} dy \right), \quad (8.33)$$

where  $c = 1/8400$ . Since by a change of variables  $t = -x$   $\int_{x=-50}^{-20} cx \, dx = -\int_{x=20}^{50} cx \, dx$ , and  $\int_{x=-20}^{20} cx \, dx = 0$ , we have  $\mathbb{E}[\tilde{x}] = 0$ .

Since the mean of  $\tilde{x}$  is zero, to compute the covariance, we only need to compute  $\mathbb{E}[\tilde{x}\tilde{y}]$ . By iterated expectation, we can take the mean of  $\mu_{\tilde{x}\tilde{y}}(\tilde{y})$ , which is computed by first deriving the conditional mean function  $\mu_{\tilde{x}\tilde{y}}(y)$  for all possible values of  $y$  and plugging in  $\tilde{y}$ . For  $y$  in  $[-50, -20]$  or  $[20, 50]$ , we have

$$\mu_{\tilde{x}\tilde{y}}(y) = \int_{x=-50}^{50} \frac{yx}{100} \, dx \quad (8.34)$$

$$= \frac{yx}{100} \int_{x=-50}^{50} x \, dx \quad (8.35)$$

$$= 0. \quad (8.36)$$

For  $y$  in  $[-20, 20]$ , similarly

$$\mu_{\tilde{x}\tilde{y}}(y) = \frac{1}{60} \left( \int_{x=-50}^{-20} x \, dx + \int_{x=20}^{50} x \, dx \right) \quad (8.37)$$

$$= 0. \quad (8.38)$$

The conditional mean of  $\tilde{x}\tilde{y}$  is always zero for values of  $\tilde{y}$  that have nonzero probability. By iterated expectation this implies that the covariance is zero, so the variables are uncorrelated.

- 8.6 (Bernoulli random variables) Let  $\tilde{a}$  and  $\tilde{b}$  be two uncorrelated Bernoulli random variables with parameters  $\theta_{\tilde{a}}$  and  $\theta_{\tilde{b}}$  defined in the same probability space. Their means equal  $\theta_{\tilde{a}}$  and  $\theta_{\tilde{b}}$  respectively, so

$$\text{Cov}[\tilde{a}, \tilde{b}] = \mathbb{E}[\tilde{a}\tilde{b}] - \theta_{\tilde{a}}\theta_{\tilde{b}}. \quad (8.39)$$

Since

$$\mathbb{E}[\tilde{a}\tilde{b}] = \sum_{a \in \{0,1\}} \sum_{b \in \{0,1\}} ab p_{\tilde{a}, \tilde{b}}(a, b) \quad (8.40)$$

$$= p_{\tilde{a}, \tilde{b}}(1, 1), \quad (8.41)$$

so

$$p_{\tilde{a}, \tilde{b}}(1, 1) = \theta_{\tilde{a}}\theta_{\tilde{b}} = p_{\tilde{a}}(1)p_{\tilde{b}}(1). \quad (8.42)$$

In addition,

$$\theta_{\tilde{a}} = p_{\tilde{a}}(1) = p_{\tilde{a}, \tilde{b}}(1, 0) + p_{\tilde{a}, \tilde{b}}(1, 1) \quad (8.43)$$

$$= p_{\tilde{a}, \tilde{b}}(1, 0) + \theta_{\tilde{a}}\theta_{\tilde{b}}, \quad (8.44)$$

so

$$p_{\tilde{a}, \tilde{b}}(1, 0) = \theta_{\tilde{a}}(1 - \theta_{\tilde{b}}) = p_{\tilde{a}}(1)p_{\tilde{b}}(0) \quad (8.45)$$

and by the same argument  $p_{\tilde{a}, \tilde{b}}(0, 1) = p_{\tilde{a}}(0)p_{\tilde{b}}(1)$ . Finally,

$$1 - \theta_{\tilde{a}} = p_{\tilde{a}}(0) = p_{\tilde{a}, \tilde{b}}(0, 0) + p_{\tilde{a}, \tilde{b}}(0, 1) \quad (8.46)$$

$$= p_{\tilde{a}, \tilde{b}}(0, 0) + (1 - \theta_{\tilde{a}})\theta_{\tilde{b}}, \quad (8.47)$$



so

$$p_{\tilde{a}, \tilde{b}}(0, 0) = (1 - \theta_{\tilde{a}})(1 - \theta_{\tilde{b}}) = p_{\tilde{a}}(0)p_{\tilde{b}}(0). \quad (8.48)$$

We conclude that uncorrelated Bernoulli random variables are independent.

- 8.7 (Sample mean and variance of standardized data) We define  $X := \{x_1, x_2, \dots, x_n\}$  and  $Y := \{y_1, y_2, \dots, y_n\}$ .

$$m(S_X) = \frac{1}{n} \sum_{i=1}^n s(x_i) \quad (8.49)$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{x_i - m(X)}{\sqrt{v(X)}} \quad (8.50)$$

$$= \frac{1}{\sqrt{v(X)}} \left( \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{n} \sum_{i=1}^n m(X) \right) \quad (8.51)$$

$$= \frac{m(X) - m(X)}{\sqrt{v(X)}} = 0. \quad (8.52)$$

By the same argument,  $m(Y) = 0$ .

$$v(S_X) = \frac{1}{n-1} \sum_{i=1}^n (s(x_i) - m(S_X))^2 \quad (8.53)$$

$$= \frac{1}{n-1} \sum_{i=1}^n s(x_i)^2 \quad (8.54)$$

$$= \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - m(X)}{\sqrt{v(X)}} \right)^2 \quad (8.55)$$

$$= \frac{1}{v(X)} \frac{1}{n-1} \sum_{i=1}^n (x_i - m(X))^2 \quad (8.56)$$

$$= 1. \quad (8.57)$$

By the same argument,  $v(Y) = 1$ .

$$c(S_X, S_Y) = \frac{1}{n-1} \sum_{i=1}^n (s(x_i) - m(S_X))(s(y_i) - m(S_Y)) \quad (8.58)$$

$$= \frac{1}{n-1} \sum_{i=1}^n s(x_i)s(y_i) \quad (8.59)$$

$$= \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - m(X))(y_i - m(Y))}{\sqrt{v(X)v(Y)}} \quad (8.60)$$

$$= \rho_{X,Y}. \quad (8.61)$$

- 8.8 (Properties of the sample correlation coefficient)

a) By the results in Exercise 8.7 and the formula for the OLS estimator,

$$\frac{1}{n-1} \sum_{i=1}^n r_i^2 = \frac{1}{n-1} \sum_{i=1}^n \left( y_i - \sqrt{v(Y)} \rho_{X,Y} \left( \frac{x - m(X)}{\sqrt{v(X)}} \right) - m(Y) \right)^2 \quad (8.62)$$

$$= \frac{v(Y)}{n-1} \sum_{i=1}^n (s(y_i) - \rho_{X,Y} s(x_i))^2 \quad (8.63)$$

$$= \frac{v(Y)}{n-1} \left( \sum_{i=1}^n s(y_i)^2 + \rho_{X,Y}^2 \sum_{i=1}^n s(x_i)^2 - 2\rho_{X,Y} \sum_{i=1}^n s(x_i)s(y_i) \right) \quad (8.64)$$

$$= v(Y) (1 - \rho_{X,Y}^2). \quad (8.65)$$

b) By (8.65),

$$(1 - \rho_{X,Y}^2) v(Y) \geq \frac{1}{n-1} \sum_{i=1}^n r_i^2 \geq 0, \quad (8.66)$$

which implies that  $\rho_{X,Y}^2 \leq 1$ , unless  $v(Y)$  is zero.

c) By (8.65), if

$$\rho_{X,Y} = \pm 1 \quad (8.67)$$

then

$$\frac{1}{n-1} \sum_{i=1}^n r_i^2 = v(Y) (1 - \rho_{X,Y}^2) = 0, \quad (8.68)$$

which implies that  $r_i = 0$  for  $1 \leq i \leq n$ , so that

$$y_i = \ell_{\text{OLS}}(x_i) \quad (8.69)$$

$$= \sqrt{v(Y)} \rho_{X,Y} \left( \frac{x - m(X)}{\sqrt{v(X)}} \right) + m(Y). \quad (8.70)$$

## 8.9 (Decomposition of sample variance and coefficient of determination)

a) The residual has mean zero,

$$m(R) = \frac{1}{n} \sum_{i=1}^n \left( y_i - \sqrt{v(Y)} \rho_{X,Y} \left( \frac{x - m(X)}{\sqrt{v(X)}} \right) - m(Y) \right) \quad (8.71)$$

$$= \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n m(Y) - \frac{\sqrt{v(Y)} \rho_{X,Y}}{\sqrt{v(X)}} \left( \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{n} \sum_{i=1}^n m(X) \right) \quad (8.72)$$

$$= m(Y) - m(Y) - \frac{\sqrt{v(Y)} \rho_{X,Y}}{\sqrt{v(X)}} (m(X) - m(X)) \quad (8.73)$$

$$= 0. \quad (8.74)$$

Therefore, by the results in Exercise 8.7 and the formula for the OLS estimator,

$$c(X, R) = \frac{1}{n-1} \sum_{i=1}^n (x_i - m(X)) \left( y_i - \sqrt{v(Y)} \rho_{X,Y} \left( \frac{x - m(X)}{\sqrt{v(X)}} \right) - m(Y) \right) \quad (8.75)$$

$$= \frac{\sqrt{v(X)v(Y)}}{n-1} \sum_{i=1}^n s(x_i) (s(y_i) - \rho_{X,Y} s(x_i)) \quad (8.76)$$

$$= \sqrt{v(X)v(Y)} \left( \frac{1}{n-1} \sum_{i=1}^n s(x_i) s(y_i) - \rho_{X,Y} \frac{1}{n-1} \sum_{i=1}^n s(x_i)^2 \right) \quad (8.77)$$

$$= \sqrt{v(X)v(Y)} (\rho_{X,Y}^2 - \rho_{X,Y}^2) = 0. \quad (8.78)$$

b) The sample mean is a linear operation,

$$m(M) = \frac{1}{n} \sum_{i=1}^n (a_i + b_i) \quad (8.79)$$

$$= \frac{1}{n} \sum_{i=1}^n a_i + \frac{1}{n} \sum_{i=1}^n b_i \quad (8.80)$$

$$= m(A) + m(B), \quad (8.81)$$

so

$$v(M) = \frac{1}{n-1} \sum_{i=1}^n (a_i + b_i - m(A+B))^2 \quad (8.82)$$

$$= \frac{1}{n-1} \sum_{i=1}^n (a_i - m(A) + b_i - m(B))^2 \quad (8.83)$$

$$= \frac{1}{n-1} \sum_{i=1}^n (a_i - m(A))^2 + \frac{1}{n-1} \sum_{i=1}^n (b_i - m(B))^2 \quad (8.84)$$

$$+ \frac{2}{n-1} \sum_{i=1}^n (a_i - m(A)) (b_i - m(B)) \quad (8.85)$$

$$= v(A) + v(B) + 2c(A, B). \quad (8.86)$$

c) The sample mean of the scaled dataset satisfies,

$$m(A_{\beta,\alpha}) = \frac{1}{n} \sum_{i=1}^n (\beta a_i + \alpha) \quad (8.87)$$

$$= \beta m(A) + \alpha, \quad (8.88)$$

so

$$c(A_{\beta,\alpha}, B) = \frac{1}{n-1} \sum_{i=1}^n (\beta a_i + \alpha - \beta m(A) - \alpha) (b_i - m(B)) \quad (8.89)$$

$$= \frac{\beta}{n-1} \sum_{i=1}^n (a_i - m(A)) (b_i - m(B)) \quad (8.90)$$

$$= \beta c(A, B). \quad (8.91)$$

By (8.78)  $c(X, R) = 0$ , so  $c(L, R) = 0$  because  $\ell_{\text{OLS}}(x_i)$  is an affine function of  $x_i$ . Therefore, by (8.86)

$$v(Y) = v(L) + v(R). \quad (8.92)$$

d) By (8.74) the sample of mean of the residual is zero, so by (8.65)

$$v(R) = \frac{1}{n-1} \sum_{i=1}^n r_i^2 \quad (8.93)$$

$$= v(Y) \left(1 - \rho_{X,Y}^2\right). \quad (8.94)$$

By (8.92) we conclude that

$$R^2 := \frac{v(L)}{v(Y)} \quad (8.95)$$

$$= \frac{v(Y) - v(R)}{v(Y)} \quad (8.96)$$

$$= \rho_{X,Y}^2. \quad (8.97)$$

#### 8.10 (Noisy measurement)

a) We have

$$\text{Var}[\tilde{y}] = \text{Var}[\tilde{x}] + \text{Var}[\tilde{z}] \quad (8.98)$$

$$= 1 + \sigma^2, \quad (8.99)$$

$$\mathbb{E}[\tilde{x}\tilde{y}] = \mathbb{E}[\tilde{x}^2] + \mathbb{E}[\tilde{x}\tilde{z}] \quad (8.100)$$

$$= 1, \quad (8.101)$$

$$\rho_{\tilde{x},\tilde{y}} = \frac{\mathbb{E}[\tilde{x}\tilde{y}]}{\sqrt{1 + \sigma^2}} \quad (8.102)$$

$$= \frac{1}{\sqrt{1 + \sigma^2}}. \quad (8.103)$$

The linear MMSE estimate of  $\tilde{x}$  given  $\tilde{y} = y$  equals

$$\ell_{\text{MMSE}}(y) = \frac{\rho_{\tilde{x},\tilde{y}} \sigma_{\tilde{x}} y}{\sigma_{\tilde{y}}} \quad (8.104)$$

$$= \frac{y}{1 + \sigma^2}. \quad (8.105)$$

b)

$$\mathbb{E}[(\tilde{x} - \ell_{\text{MMSE}}(\tilde{y}))^2] = (1 - \rho_{\tilde{x},\tilde{y}}^2) \sigma_{\tilde{x}}^2 \quad (8.106)$$

$$= \frac{\sigma^2}{1 + \sigma^2}. \quad (8.107)$$

- c) If  $\sigma \rightarrow 0$  the estimate equals  $y$  and the error is zero, which makes sense since in that case the noise is zero and hence  $\tilde{y} = \tilde{x}$  with probability one.
- d) If  $\sigma \rightarrow \infty$  the estimate equals zero and the error is 1, which makes sense since in that case we are better off just ignoring the measurement and setting the estimate equal to the mean of  $\tilde{x}$ . The corresponding error is consequently the variance of  $\tilde{x}$ .

- 8.11 (Affine function) Since  $\tilde{c}$  is an affine function of  $\tilde{b}$ , the standardized random variables  $s(\tilde{c})$  and  $s(\tilde{b})$  are the same, so  $\rho_{\tilde{a},\tilde{c}} = \rho_{s(\tilde{a}),s(\tilde{c})} = \rho_{s(\tilde{a}),s(\tilde{b})} = 0.25$ . The mean and variance of  $\tilde{c}$  equal

$$\mu_{\tilde{c}} = \mathbb{E}(2\tilde{b} + 3) \quad (8.108)$$

$$= 2\mathbb{E}(\tilde{b}) + 3 \quad (8.109)$$

$$= -1, \quad (8.110)$$

$$\sigma_{\tilde{c}}^2 = \text{Var}(2\tilde{b} + 3) \quad (8.111)$$

$$= 4\text{Var}(\tilde{b}) \quad (8.112)$$

$$= 36. \quad (8.113)$$

The linear minimum-mean-squared-error estimate of  $\tilde{c}$  given  $\tilde{a}$  is therefore

$$\frac{\sigma_{\tilde{c}}\rho_{\tilde{a},\tilde{c}}(a - \mu_{\tilde{a}})}{\sigma_{\tilde{a}}} + \mu_{\tilde{c}} = \frac{6 \cdot 0.25(a - 2)}{2} - 1 \quad (8.114)$$

$$= 0.75a - 2.5. \quad (8.115)$$

- 8.12 (Averaging noisy data)

- a) Let  $\tilde{w} := \sum_{i=1}^n \tilde{y}_i = n\tilde{x} + \sum_{i=1}^n \tilde{z}_i$ . We need to find the linear MMSE estimator of  $\tilde{x}$  given  $\tilde{w}$ . We have

$$\text{Var}[\tilde{w}] = \text{Var}\left[n\tilde{x} + \sum_{i=1}^n \tilde{z}_i\right] \quad (8.116)$$

$$= n^2 + n\sigma^2, \quad (8.117)$$

$$\mathbb{E}[\tilde{x}\tilde{w}] = \mathbb{E}[n\tilde{x}^2] + \sum_{i=1}^n \mathbb{E}[\tilde{x}\tilde{z}_i] \quad (8.118)$$

$$= n, \quad (8.119)$$

$$\rho_{\tilde{x},\tilde{w}} = \frac{\mathbb{E}[\tilde{x}\tilde{w}]}{\sigma_{\tilde{x}}\sigma_{\tilde{w}}} \quad (8.120)$$

$$= \frac{n}{\sqrt{n^2 + n\sigma^2}} \quad (8.121)$$

$$= \sqrt{\frac{n}{n + \sigma^2}}. \quad (8.122)$$

The linear MMSE estimate of  $\tilde{x}$  given  $\tilde{w} = w$  equals

$$\ell_{\text{MMSE}}(w) = \frac{\rho_{\tilde{x},\tilde{w}}\sigma_{\tilde{x}}w}{\sigma_{\tilde{w}}} \quad (8.123)$$

$$= \frac{w}{n + \sigma^2}. \quad (8.124)$$

Therefore the optimal  $\alpha$  equals  $\frac{1}{n + \sigma^2}$ .

- b) As  $\sigma^2 \rightarrow 0$ , the estimator is just the average of the measurements. As  $\sigma^2 \rightarrow \infty$  the estimator tends to zero: the measurements are ignored because they are just noise.  
c) The MSE equals

$$\sigma_{\tilde{x}}^2(1 - \rho_{\tilde{x},\tilde{w}}^2) = 1 - \frac{n}{n + \sigma^2} \quad (8.125)$$

$$= \frac{\sigma^2}{n + \sigma^2}, \quad (8.126)$$

which tends to zero as  $n \rightarrow \infty$ .

### 8.13 (Interference)

a) Since  $\tilde{w}$  and  $\tilde{a}$  are independent

$$\mathbb{E}[\tilde{y}] = \mathbb{E}[\tilde{w}\tilde{a}] \quad (8.127)$$

$$= \mathbb{E}[\tilde{w}] \mathbb{E}[\tilde{a}] \quad (8.128)$$

$$= 0 \quad (8.129)$$

because  $\mathbb{E}[\tilde{w}] = -0.5 + 0.5 = 0$ . Similarly,

$$\mathbb{E}[\tilde{a}\tilde{y}] = \mathbb{E}[\tilde{w}\tilde{a}^2] \quad (8.130)$$

$$= \mathbb{E}[\tilde{w}] \mathbb{E}[\tilde{a}^2] \quad (8.131)$$

$$= 0. \quad (8.132)$$

Therefore,

$$\text{Cov}[\tilde{a}, \tilde{y}] = \mathbb{E}[\tilde{a}\tilde{y}] - \mathbb{E}[\tilde{a}] \mathbb{E}[\tilde{y}] \quad (8.133)$$

$$= 0. \quad (8.134)$$

As a result, the linear MMSE estimator of  $\tilde{a}$  given  $\tilde{y} = y$  is

$$\ell_{\text{MMSE}}(y) = \mathbb{E}[\tilde{a}] \quad (8.135)$$

$$= \mu. \quad (8.136)$$

b) The MSE equals

$$\text{MSE} := \mathbb{E}[(\tilde{a} - \ell_{\text{MMSE}}(\tilde{y}))^2] \quad (8.137)$$

$$= \mathbb{E}[(\tilde{a} - \mu)^2] \quad (8.138)$$

$$= \sigma^2. \quad (8.139)$$

c) Consider the nonlinear estimator  $|\tilde{y}|$ . Since  $\tilde{a}$  is nonnegative,  $|y| = \tilde{a}$ , so the MSE is zero.

### 8.14 (Exam)

- a) The MMSE estimator is the conditional mean, which we estimate using the sample conditional mean, depicted by the circle markers in the graph below. The advantage of this estimator is that it is very flexible, it indicates that grades increase a lot proportionally to the number of hours up to a certain point, but then there are diminishing returns (if any at all). The main disadvantage is that it is very noisy due to insufficient data, which results in overfitting. Also it does not provide an estimate at 4, 7 or 9.
- b) The least-squares simple-linear-regression estimator that approximates the linear MMSE estimator is

$$\ell(h) = \sqrt{v_{\text{grade}} \rho_{\text{grade, hours}}} \left( \frac{h - m_{\text{hours}}}{\sqrt{v_{\text{hours}}}} \right) + m_{\text{grade}} \quad (8.140)$$

$$= \frac{c_{\text{grade, hours}}}{v_{\text{hours}}} h + m_{\text{grade}} - \frac{c_{\text{grade, hours}}}{v_{\text{hours}}} m_{\text{hours}}, \quad (8.141)$$

where the sample means of the hours and grades equal

$$m_{\text{grade}} = \frac{85 + 90 + 30 + 80 + 90 + 80 + 60 + 40 + 90 + 85}{10} \quad (8.142)$$

$$= 73, \quad (8.143)$$

$$m_{\text{hours}} = \frac{5 + 6 + 2 + 6 + 5 + 5 + 3 + 3 + 6 + 8}{10} \quad (8.144)$$

$$= 4.9, \quad (8.145)$$

the sample variance of the hours is

$$v_{\text{hours}} = \frac{(5 - 4.9)^2 + (6 - 4.9)^2 + (2 - 4.9)^2 + (6 - 4.9)^2 + (5 - 4.9)^2}{9} + \frac{(5 - 4.9)^2 + (3 - 4.9)^2 + (3 - 4.9)^2 + (6 - 4.9)^2 + (9 - 4.9)^2}{9} \quad (8.146)$$

$$= 4.01, \quad (8.147)$$

and the sample covariance is

$$\begin{aligned} c_{\text{grade, hours}} &= \frac{(85 - 73)(5 - 4.9) + (90 - 73)(6 - 4.9) + (30 - 73)(2 - 4.9) + (80 - 73)(6 - 4.9)}{9} \\ &+ \frac{(90 - 73)(5 - 4.9) + (80 - 73)(5 - 4.9) + (60 - 73)(3 - 4.9) + (40 - 73)(3 - 4.9)}{9} \\ &+ \frac{(90 - 73)(6 - 4.9) + (85 - 73)(8 - 4.9)}{9} \end{aligned} \quad (8.148)$$

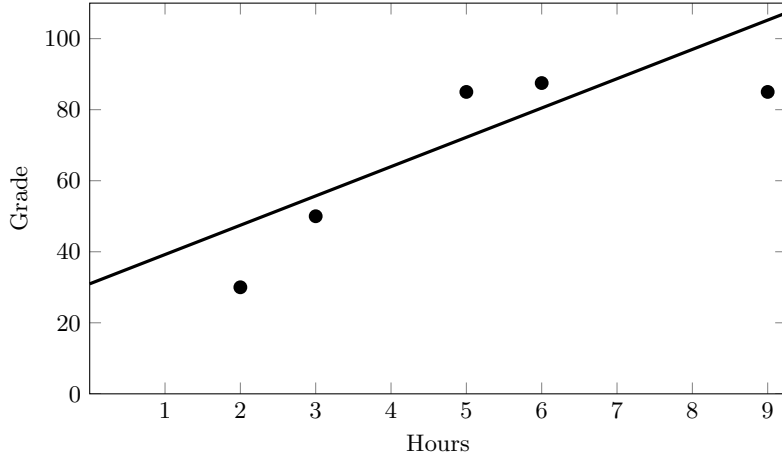
$$= 33.11. \quad (8.149)$$

The least-squares estimator equals

$$\ell(h) = \frac{33.1}{4.01}h + 73 - \frac{34.4}{4.01}4.9 \quad (8.150)$$

$$= 8.25h + 30.97, \quad (8.151)$$

The estimator is depicted as a straight line on the graph below. Its advantage is that it is less noisy than the nonlinear sample conditional mean estimator, and produces estimates for any number of hours. However it underfits the data, as it does not capture the fact that the grade increase eventually saturates when the number of hours exceed 5.



8.15 (Ice cream sales) The correlation coefficient equals

$$\rho = \frac{\text{Cov}[\text{adv}, \text{sales}]}{\sqrt{\text{Var}[\text{adv}] \text{Var}[\text{adv}]}} \quad (8.152)$$

$$= \frac{960}{\sqrt{900 \cdot 1600}} \quad (8.153)$$

$$= 0.8. \quad (8.154)$$

Sales and advertising are highly correlated, but this does not necessarily imply that advertising causes higher sales. From the table, it looks like temperature acts as a confounding factor. When it's warmer the company spends more on advertising, and people also buy more ice creams. Consequently, the argument of the marketing department is not very compelling.

8.16 (Nonconfounding variable)

- a) By the theorem, the covariance between the observed outcome  $\tilde{y}$  and the treatment  $\tilde{t}$  equals  $\beta_{\text{treat}} + \beta_{\text{conf}}\sigma_{\tilde{t}, \tilde{c}}$ . The mean of  $\tilde{y}$  is zero by linearity of expectation, since  $\tilde{t}$ ,  $\tilde{c}$  and  $\tilde{z}$  all have zero mean. Since  $\tilde{z}$  is independent from  $\tilde{t}$  and  $\tilde{c}$ , again by linearity of expectation, under the assumptions of the theorem, the variance of  $\tilde{y}$  equals

$$\text{Var}[\tilde{y}] \quad (8.155)$$

$$= \mathbb{E}[\tilde{y}^2] \quad (8.156)$$

$$= \mathbb{E}[(\beta_{\text{treat}}\tilde{t} + \beta_{\text{conf}}\tilde{c} + \tilde{z})^2] \quad (8.157)$$

$$= \mathbb{E}[\beta_{\text{treat}}^2\tilde{t}^2 + \beta_{\text{conf}}^2\tilde{c}^2 + \tilde{z}^2 + 2(\beta_{\text{treat}}\tilde{t} + \beta_{\text{conf}}\tilde{c})\tilde{z} + 2\beta_{\text{treat}}\beta_{\text{conf}}\tilde{t}\tilde{c}] \quad (8.158)$$

$$= \beta_{\text{treat}}^2\mathbb{E}[\tilde{t}^2] + \beta_{\text{conf}}^2\mathbb{E}[\tilde{c}^2] + \mathbb{E}[\tilde{z}^2] + 2\mathbb{E}[\beta_{\text{treat}}\tilde{t} + \beta_{\text{conf}}\tilde{c}]\mathbb{E}[\tilde{z}] + 2\beta_{\text{treat}}\beta_{\text{conf}}\mathbb{E}[\tilde{t}\tilde{c}]$$

$$= \beta_{\text{treat}}^2 + \beta_{\text{conf}}^2 + \sigma_{\tilde{z}}^2 + 2\beta_{\text{treat}}\beta_{\text{conf}}\sigma_{\tilde{t}, \tilde{c}}, \quad (8.159)$$

since  $\tilde{t}$  and  $\tilde{c}$  are standardized, so their mean square is one. Consequently, the correlation



coefficient  $\rho_{\tilde{y}, \tilde{t}}$  equals

$$\rho_{\tilde{y}, \tilde{t}} := \frac{\sigma_{\tilde{t}, \tilde{c}}}{\sqrt{\text{Var}[\tilde{y}] \text{Var}[\tilde{t}]}} \quad (8.160)$$

$$= \frac{\beta_{\text{treat}} + \beta_{\text{conf}} \sigma_{\tilde{t}, \tilde{c}}}{\sqrt{\sigma_{\tilde{z}}^2 + \beta_{\text{treat}}^2 + \beta_{\text{conf}}^2 + 2\beta_{\text{treat}}\beta_{\text{conf}}\sigma_{\tilde{t}, \tilde{c}}}}. \quad (8.161)$$

- b) The limit of the correlation coefficient as  $\sigma_{\tilde{z}} \rightarrow \infty$  equals zero. In that case, the non-confounding variable dominates the observed outcome, rendering it uncorrelated with the treatment, because the treatment variable and the nonconfounding variable are uncorrelated.
- c) In Example 8.38,  $\beta_{\text{treat}} := 0$  and  $\beta_{\text{conf}} := 1$ , so setting  $\sigma_{\tilde{z}}^2 := 1$ ,

$$\rho_{\tilde{y}, \tilde{t}} = \frac{\sigma_{\tilde{t}, \tilde{c}}}{\sqrt{2}}. \quad (8.162)$$

For  $\sigma_{\tilde{t}, \tilde{c}} := 0.8$ ,  $\rho_{\tilde{y}, \tilde{t}} = 0.57$ . For  $\sigma_{\tilde{t}, \tilde{c}} := -0.8$ ,  $\rho_{\tilde{y}, \tilde{t}} = -0.57$ . For  $\sigma_{\tilde{t}, \tilde{c}} := 0$ ,  $\rho_{\tilde{y}, \tilde{t}} = 0$ .

## Estimation Of Population Parameters

### Exercises

9.1 (Markov's and Chebyshev's inequalities are tight)

- a) Consider a random variable  $\tilde{a}$  that is equal to  $c$  with probability  $\theta$ , and equal to zero otherwise. Then  $P(\tilde{a} \geq c) = \theta$  and

$$\mathbb{E}[\tilde{a}] = cP(\tilde{a} = c) \quad (9.1)$$

$$= c\theta. \quad (9.2)$$

- b) Consider a random variable  $\tilde{b}$  with pmf  $p_{\tilde{b}}$  such that

$$p_{\tilde{b}}(\mu - c) = \frac{\theta}{2}, \quad (9.3)$$

$$p_{\tilde{b}}(\mu) = 1 - \theta, \quad (9.4)$$

$$p_{\tilde{b}}(\mu + c) = \frac{\theta}{2}. \quad (9.5)$$

Then  $P(|\tilde{b} - \mu| \geq c) = \theta$ ,

$$\mathbb{E}[\tilde{b}] = \frac{\theta}{2}(\mu - c) + \frac{\theta}{2}(\mu + c) + (1 - \theta)\mu \quad (9.6)$$

$$= \mu, \quad (9.7)$$

and

$$\text{Var}[\tilde{b}] = \mathbb{E}[(\tilde{b} - \mu)^2] \quad (9.8)$$

$$= c^2P(\tilde{a} = \mu + c) + c^2P(\tilde{a} = \mu - c) \quad (9.9)$$

$$= c^2\theta. \quad (9.10)$$

9.2 (Online poll)

- a) The proportion of Democrat voters in the whole population is  $\theta_{\text{tot}} = \alpha\theta_{\text{young}} + (1 - \alpha)\theta_{\text{old}}$ . The fraction of young voters in the poll that intend to vote Democrat  $y/n_1$  is a reasonable estimate for  $\theta_{\text{young}}$ . The fraction of old voters in the poll that intend to vote Democrat  $o/n_2$  is a reasonable estimate for  $\theta_{\text{old}}$ . The estimator for  $\theta_{\text{tot}} = \alpha\theta_{\text{young}} + (1 - \alpha)\theta_{\text{old}}$  therefore equals

$$e(o, y) := \frac{\alpha y}{n_1} + \frac{(1 - \alpha)o}{n_2}. \quad (9.11)$$

- b) We have  $y = 50$ ,  $n_1 = 70$ ,  $o = 10$ ,  $n_2 = 30$ ,  $\alpha = 0.25$ , so

$$e(10, 50) = \frac{1}{4} \cdot \frac{50}{70} + \frac{3}{4} \cdot \frac{10}{30} = 42.9\%. \quad (9.12)$$

- c) Let us assume that the young people in the population are sampled independently and uniformly at random with replacement, and so are the old people. Let  $\tilde{y}$  denote the number of young voters in the poll that intend to vote Democrat, and let  $\tilde{o}$  denote the number of old voters in the poll that intend to vote Democrat. By the notes,  $\tilde{y}/n_1$  and  $\tilde{o}/n_2$  are unbiased estimators of  $\theta_{\text{young}}$  and  $\theta_{\text{old}}$  respectively. Consequently, by linearity of expectation

$$\mathbb{E}[e(\tilde{o}, \tilde{y})] := \mathbb{E}\left[\frac{\alpha\tilde{y}}{n_1} + \frac{(1-\alpha)\tilde{o}}{n_2}\right] \quad (9.13)$$

$$= \frac{\alpha\mathbb{E}[\tilde{y}]}{n_1} + \frac{(1-\alpha)\mathbb{E}[\tilde{o}]}{n_2} \quad (9.14)$$

$$= \alpha\theta_{\text{young}} + (1-\alpha)\theta_{\text{old}} = \theta_{\text{tot}}, \quad (9.15)$$

so the estimator is unbiased (under our assumptions).

- d) Assuming independent and uniform random sampling with replacement,

$$\text{Var}\left[\frac{\tilde{y}}{n_1}\right] = \frac{\theta_{\text{young}}(1-\theta_{\text{young}})}{n_1}, \quad (9.16)$$

$$\text{Var}\left[\frac{\tilde{o}}{n_2}\right] = \frac{\theta_{\text{old}}(1-\theta_{\text{old}})}{n_2}. \quad (9.17)$$

Assuming  $\tilde{y}$  and  $\tilde{o}$  are independent, the variance of our estimator is

$$\text{Var}[e(\tilde{o}, \tilde{y})] = \text{Var}\left[\frac{\alpha\tilde{y}}{n_1} + \frac{(1-\alpha)\tilde{o}}{n_2}\right] \quad (9.18)$$

$$= \alpha^2 \text{Var}\left[\frac{\tilde{y}}{n_1}\right] + (1-\alpha)^2 \text{Var}\left[\frac{\tilde{o}}{n_2}\right] \quad (9.19)$$

$$= \frac{\alpha^2\theta_{\text{young}}(1-\theta_{\text{young}})}{n_1} + \frac{(1-\alpha)^2\theta_{\text{old}}(1-\theta_{\text{old}})}{n_2}. \quad (9.20)$$

Since the estimator is unbiased, by Chebyshev's inequality, for any  $\epsilon$

$$\mathbb{P}(|e(\tilde{o}, \tilde{y}) - \theta_{\text{tot}}| > \epsilon) \leq \frac{\text{Var}[e(\tilde{o}, \tilde{y})]}{\epsilon^2}, \quad (9.21)$$

which converges to zero as  $n_1 \rightarrow \infty$  and  $n_2 \rightarrow \infty$ , so the estimator converges to  $\theta_{\text{tot}}$  in probability and is therefore consistent.

### 9.3 (Participation)

- a) A reasonable estimator for  $\alpha_R$  is the fraction of people in the poll who are certain to vote and will vote Democrat. The estimator is unbiased if the people are selected independently with replacement, because the probability that each of them is in the category of interest is exactly  $\alpha_R$ . Following the same reasoning as in Example 9.14, the standard error is  $\sqrt{(\alpha_R(1-\alpha_R))/n} = \sqrt{(\alpha_R(1-\alpha_R))/150}$ . For the observed data, the estimator equals  $30/150 = 0.2$ .
- b) Let  $\alpha_D$  denote the fraction of people who are certain to vote and will vote Democrat,  $\beta_R$  the fraction of people in the population who are likely to vote and will vote Republican, and  $\beta_D$  the fraction of people in the population who are likely to vote and will vote Democrat. According to the study, the fraction of people who will actually vote equals  $0.8(\alpha_D + \alpha_R) + 0.1(\beta_D + \beta_R)$ . Out of them, the ones who will vote Republican are  $0.8\alpha_R + 0.1\beta_R$ . As explained above, we can obtain unbiased estimators for  $\alpha_R$ ,  $\alpha_D$ ,  $\beta_D$  and  $\beta_R$  by computing the corresponding fractions from the poll, under the assumption

that the poll participants are selected independently with replacement. This yields the following estimate for the fraction of people who will vote Republican:

$$\frac{0.8 \frac{30}{150} + 0.1 \frac{10}{150}}{0.8 \left( \frac{20}{150} + \frac{30}{150} \right) + 0.1 \left( \frac{40}{150} + \frac{10}{150} \right)} = \frac{0.8 \cdot 3 + 0.1 \cdot 1}{0.8 \cdot 5 + 0.1 \cdot 5} \quad (9.22)$$

$$= 55.6\%. \quad (9.23)$$

9.4 (The sample covariance is unbiased)

a)

$$\sum_{i=1}^n (\tilde{x}_i - \tilde{m}_1)(\tilde{y}_i - \tilde{m}_2) = \sum_{i=1}^n \tilde{x}_i \tilde{y}_i - \tilde{m}_1 \sum_{i=1}^n \tilde{y}_i - \tilde{m}_2 \sum_{i=1}^n \tilde{x}_i + n \tilde{m}_1 \tilde{m}_2 \quad (9.24)$$

$$= \sum_{i=1}^n \tilde{x}_i \tilde{y}_i - n \tilde{m}_1 \tilde{m}_2. \quad (9.25)$$

b) For any  $1 \leq i \leq n$

$$\text{Cov} [\tilde{x}_i, \tilde{y}_i] := \mathbb{E} [(\tilde{x}_i - \mathbb{E} [\tilde{x}_i])(\tilde{y}_i - \mathbb{E} [\tilde{y}_i])] \quad (9.26)$$

$$= \mathbb{E} [(\tilde{x}_i - \mu_a)(\tilde{y}_i - \mu_b)] \quad (9.27)$$

$$= \sum_{k=1}^N (a_k - \mu_a)(b_k - \mu_b) p_{\tilde{k}_i}^-(k) \quad (9.28)$$

$$= \frac{1}{N} \sum_{k=1}^N (a_k - \mu_a)(b_k - \mu_b) \quad (9.29)$$

$$= c_{\text{pop}}. \quad (9.30)$$

Since the means of  $\tilde{m}_1$  and  $\tilde{m}_2$  are  $\mu_a$  and  $\mu_b$  respectively, the result follows.

c) By linearity of the covariance operator

$$\mathbb{E} [\tilde{m}_1 \tilde{m}_2] = \text{Cov} [\tilde{m}_1, \tilde{m}_2] + \mu_a \mu_b \quad (9.31)$$

$$= \text{Cov} \left[ \frac{1}{n} \sum_{i=1}^n \tilde{x}_i, \frac{1}{n} \sum_{j=1}^n \tilde{y}_j \right] + \mu_a \mu_b \quad (9.32)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{Cov} [\tilde{x}_i, \tilde{y}_j] + \mu_a \mu_b \quad (9.33)$$

$$= \frac{c_{\text{pop}}}{n} + \mu_a \mu_b. \quad (9.34)$$

d) Combining the above results

$$\mathbb{E}[\tilde{c}] = \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (\tilde{x}_i - \tilde{m}_1)(\tilde{y}_i - \tilde{m}_2)\right] \quad (9.35)$$

$$= \mathbb{E}\left[\frac{1}{n-1} \left(\sum_{i=1}^n \tilde{x}_i \tilde{y}_i - n\tilde{m}_1 \tilde{m}_2\right)\right] \quad (9.36)$$

$$= \frac{1}{n-1} \left(\sum_{i=1}^n \mathbb{E}[\tilde{x}_i \tilde{y}_i] - n\mathbb{E}[\tilde{m}_1 \tilde{m}_2]\right) \quad (9.37)$$

$$= \frac{1}{n-1} (nc_{\text{pop}} + n\mu_a \mu_b - c_{\text{pop}} - n\mu_a \mu_b) \quad (9.38)$$

$$= c_{\text{pop}}. \quad (9.39)$$

### 9.5 (Herding)

a) Note that

$$\tilde{h}_3 = \frac{\tilde{m}_1}{4} + \frac{\frac{\tilde{m}_1 + \tilde{m}_2}{2}}{4} + \frac{\tilde{m}_3}{2} \quad (9.40)$$

$$= \frac{3\tilde{m}_1}{8} + \frac{\tilde{m}_2}{8} + \frac{\tilde{m}_3}{2}. \quad (9.41)$$

Let  $\theta_{\text{pop}}$  denote the population proportion, such that

$$\mathbb{E}[\tilde{m}_1] = \mathbb{E}[\tilde{m}_2] = \mathbb{E}[\tilde{m}_3] = \theta_{\text{pop}}. \quad (9.42)$$

By linearity of expectation and (9.41)

$$\mathbb{E}[\tilde{h}_2] = \frac{\mathbb{E}[\tilde{m}_1] + \mathbb{E}[\tilde{m}_2]}{2} \quad (9.43)$$

$$= \frac{\theta_{\text{pop}}}{2}, \quad (9.44)$$

$$\mathbb{E}[\tilde{h}_3] = \frac{3\mathbb{E}[\tilde{m}_1]}{8} + \frac{\mathbb{E}[\tilde{m}_2]}{8} + \frac{4\mathbb{E}[\tilde{m}_3]}{8} \quad (9.45)$$

$$= \theta_{\text{pop}}. \quad (9.46)$$

Both herded estimators are unbiased. In this situation, herding does not produce bias.

b) The polls are independent, so by (9.41)

$$\text{Var}[\tilde{h}_2] = \frac{\text{Var}[\tilde{m}_1] + \text{Var}[\tilde{m}_2]}{4} \quad (9.47)$$

$$= \frac{s^2}{2}, \quad (9.48)$$

$$\text{Var}[\tilde{h}_3] = \frac{9\text{Var}[\tilde{m}_1]}{64} + \frac{\text{Var}[\tilde{m}_2]}{64} + \frac{\text{Var}[\tilde{m}_3]}{4} \quad (9.49)$$

$$= \frac{13s^2}{32}. \quad (9.50)$$

Consequently the standard errors equal

$$\text{se}[\tilde{h}_2] = \frac{s}{\sqrt{2}} = 0.71s, \quad (9.51)$$

$$\text{se}[\tilde{h}_3] = \sqrt{\frac{13}{32}}s = 0.64s. \quad (9.52)$$

Herding improves each individual estimator.

c) Note that by (9.41)

$$\tilde{h}_{\text{all}} = \frac{\tilde{m}_1}{3} + \frac{\frac{\tilde{m}_1 + \tilde{m}_2}{2}}{3} + \frac{\tilde{h}_3}{3} \quad (9.53)$$

$$= \left(\frac{1}{3} + \frac{1}{6} + \frac{3}{24}\right) \tilde{m}_1 + \left(\frac{1}{6} + \frac{1}{24}\right) \tilde{m}_2 + \frac{\tilde{m}_3}{6} \quad (9.54)$$

$$= \frac{15\tilde{m}_3}{24} + \frac{5\tilde{m}_3}{24} + \frac{\tilde{m}_3}{6}. \quad (9.55)$$

By the independence assumption,

$$\text{Var} [\tilde{m}_{\text{all}}] = \frac{\text{Var} [\tilde{m}_1] + \text{Var} [\tilde{m}_2] + \text{Var} [\tilde{m}_3]}{9} \quad (9.56)$$

$$= \frac{s^2}{3}, \quad (9.57)$$

$$\text{Var} [\tilde{h}_{\text{all}}] = \frac{13^2 \text{Var} [\tilde{m}_1]}{24^2} + \frac{5^2 \text{Var} [\tilde{m}_2]}{24^2} + \frac{\text{Var} [\tilde{m}_3]}{36} \quad (9.58)$$

$$= \frac{133s^2}{288}. \quad (9.59)$$

Consequently the standard errors equal

$$\text{se} [\tilde{m}_{\text{all}}] = \frac{s}{\sqrt{3}} = 0.58s, \quad (9.60)$$

$$\text{se} [\tilde{h}_{\text{all}}] = \sqrt{\frac{133}{288}}s = 0.68s. \quad (9.61)$$

The herded estimator  $\tilde{h}_{\text{all}}$  has a larger standard error than the non-herded estimator  $\tilde{m}_{\text{all}}$ , so herding is problematic because it degrades the aggregated estimator. This happens because the first poll (and to a lesser extent the second poll) are overrepresented in the herded estimators.

## 9.6 (Multiplicative noise)

a) The estimator is

$$\tilde{h} := \frac{\sum_{i=1}^n x_i}{\mu}. \quad (9.62)$$

As long as  $\mu$  is nonzero, it is unbiased, since by linearity of expectation

$$\mathbb{E} [\tilde{h}] = \mathbb{E} \left[ \frac{\frac{1}{n} \sum_{i=1}^n \tilde{x}_i}{\mu} \right] \quad (9.63)$$

$$= \frac{\gamma}{\mu} \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\tilde{a}_i] \quad (9.64)$$

$$= \frac{n\mu}{n\mu} \gamma = \gamma. \quad (9.65)$$

b) The variance of the estimator equals

$$\text{Var} [\tilde{h}] = \text{Var} \left[ \frac{\frac{1}{n} \sum_{i=1}^n \tilde{x}_i}{\mu} \right] \quad (9.66)$$

$$= \frac{\gamma^2}{\mu^2} \frac{1}{n^2} \sum_{i=1}^n \text{Var} [\tilde{a}_i] \quad (9.67)$$

$$= \frac{\gamma^2 \sigma^2}{\mu^2 n}, \quad (9.68)$$

so the standard error equals

$$\text{se}[\tilde{h}] = \frac{\gamma\sigma}{\mu\sqrt{n}}. \quad (9.69)$$

c) Since the estimator is unbiased, by Chebyshev's inequality and the previous question,

$$\mathbb{P}(|\tilde{h} - \gamma| > \epsilon) \leq \frac{\text{Var}[\tilde{h}]}{\epsilon^2} \quad (9.70)$$

$$\leq \frac{\gamma^2\sigma^2}{\mu^2n\epsilon}, \quad (9.71)$$

which converges to zero as  $n \rightarrow \infty$ , so the estimator is consistent.

#### 9.7 (Consistency of sample median)

- a) If  $n$  is odd, then for the median to be greater than  $\gamma + \epsilon$  at least  $(n+1)/2 \geq n/2$  elements have to be greater than  $\gamma + \epsilon$ . If  $n$  is even, then  $n/2$  have to be greater than  $\gamma + \epsilon$ .
- b) By the assumption that  $\gamma$  is the only point that satisfies  $F_{\tilde{x}_i}(\gamma) = 1/2$ , there exists a constant  $\epsilon' > 0$  such that for any  $i$  the probability that  $\tilde{x}_i > \gamma + \epsilon$  is

$$\mathbb{P}(\tilde{x}_i > \gamma + \epsilon) = 1/2 - \epsilon' := \theta. \quad (9.72)$$

The random variable  $\tilde{b}$  is binomial with parameters  $n$  and  $\theta$ . As a result, we have

$$\mathbb{P}\left(\tilde{b} \geq \frac{n}{2}\right) = \mathbb{P}\left(\tilde{b} - n\theta \geq \frac{n}{2} - n\theta\right) \quad (9.73)$$

$$\leq \mathbb{P}(|\tilde{b} - n\theta| \geq n\epsilon') \quad (9.74)$$

$$\leq \frac{\text{Var}[\tilde{b}]}{(n\epsilon')^2} \quad \text{by Chebyshev's inequality} \quad (9.75)$$

$$= \frac{n\theta(1-\theta)}{n^2(\epsilon')^2} \quad (9.76)$$

$$= \frac{\theta(1-\theta)}{n(\epsilon')^2}, \quad (9.77)$$

which converges to zero as  $n \rightarrow \infty$ . This completes the proof, combined with the first part.

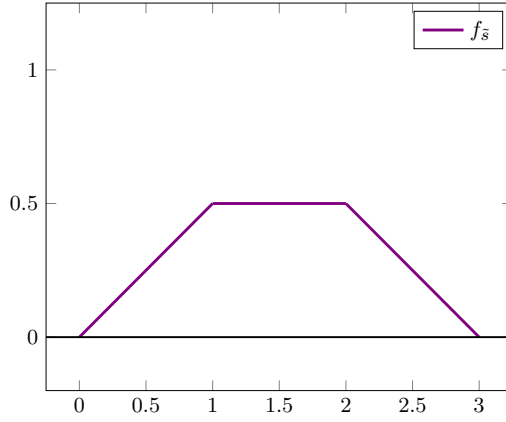
- 9.8 (Watering a plant) We represent the water poured by the roommates using the random variables  $\tilde{a}$  and  $\tilde{b}$ . We are interested in the pdf of the sum  $\tilde{s} := \tilde{a} + \tilde{b}$ . It equals

$$f_{\tilde{s}}(s) = \int_{u=-\infty}^{\infty} f_{\tilde{a}}(s-u) f_{\tilde{b}}(u) du \quad (9.78)$$

$$= \frac{1}{2} \int_{u=0}^2 f_{\tilde{a}}(s-u) du \quad (9.79)$$

$$= \begin{cases} \frac{1}{2} \int_{u=0}^s du = \frac{s}{2} & \text{if } s \leq 1, \\ \frac{1}{2} \int_{u=s-1}^s du = \frac{1}{2} & \text{if } 1 \leq s \leq 2, \\ \frac{1}{2} \int_{u=s-1}^2 du = \frac{3-s}{2} & \text{if } 2 \leq s \leq 3. \end{cases} \quad (9.80)$$

The pdf of  $\tilde{s}$  is shown below.



## 9.9 (Confidence intervals)

- a) 72 out of 80 intervals contain  $\mu$ . A reasonable estimate for  $1 - \alpha$  is therefore 0.9, so the estimate for  $\alpha$  is 0.1.
- b) •  $X = B$  and  $Y = A$ . The central limit theorem produces much narrower confidence intervals.
  - $X = B$  and  $Y = A$ . More data results in narrower confidence intervals (see equation in the notes).
  - $X = A$  and  $Y = B$ . Higher variance results in wider confidence intervals (see equation in the notes).

## 9.10 (Cholesterol)

- a) Let  $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_{200}$  be samples from the men, and  $\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_{100}$  the sample from the women. If each individual is selected independently with replacement, then  $\text{Var}[\tilde{x}_i] = \sigma_{\text{men}}^2$  and  $\text{Var}[\tilde{y}_i] = \sigma_{\text{women}}^2$ . As a result, by independence

$$\text{Var} \left[ \frac{1}{200} \sum_{i=1}^{200} \tilde{x}_i - \frac{1}{100} \sum_{i=1}^{100} \tilde{y}_i \right] = \text{Var} \left[ \frac{1}{200} \sum_{i=1}^{200} \tilde{x}_i \right] + \text{Var} \left[ -\frac{1}{100} \sum_{i=1}^{100} \tilde{y}_i \right] \quad (9.81)$$

$$= \text{Var} \left[ \frac{1}{200} \sum_{i=1}^{200} \tilde{x}_i \right] + \text{Var} \left[ \frac{1}{100} \sum_{i=1}^{100} \tilde{y}_i \right] \quad (9.82)$$

$$= \frac{\sigma_{\text{men}}^2}{200} + \frac{\sigma_{\text{women}}^2}{100}. \quad (9.83)$$

Therefore the standard error is  $\sqrt{\frac{\sigma_{\text{men}}^2}{200} + \frac{\sigma_{\text{women}}^2}{100}}$ .

- b) By the central limit theorem, if  $\mu_{\text{men}}$  and  $\mu_{\text{women}}$  are the population means for men and women respectively, then  $\frac{1}{200} \sum_{i=1}^{200} \tilde{x}_i$  is approximately Gaussian with mean  $\mu_{\text{men}}$  and variance  $\frac{\sigma_{\text{men}}^2}{200}$ , and  $\frac{1}{100} \sum_{i=1}^{100} \tilde{y}_i$  is approximately Gaussian with mean  $\mu_{\text{women}}$  and variance  $\frac{\sigma_{\text{women}}^2}{100}$ . As a result,  $\frac{\sigma_{\text{men}}^2}{200} - \frac{1}{100} \sum_{i=1}^{100} \tilde{y}_i$  is approximately Gaussian with mean  $\mu_{\text{men}} - \mu_{\text{women}}$  and variance  $\frac{\sigma_{\text{men}}^2}{200} + \frac{\sigma_{\text{women}}^2}{100} = 1.5$ . By Lemma 9.42 the following is a 0.95 confidence interval for  $\mu_{\text{men}} - \mu_{\text{women}}$ :

$$\left[ 20 - 1.96\sqrt{1.5}, 20 + 1.96\sqrt{1.5} \right] = [17.6, 22.4]. \quad (9.84)$$

## 9.11 (Difference of proportions)



- a) The variance of the first proportion equals  $\theta_A(1 - \theta_A)/n_A$  because it is a binomial with parameters  $n_A$  and  $\theta_A$  multiplied by  $1/n_A$ . By the same reasoning the variance of the second proportion is  $\theta_B(1 - \theta_B)/n_B$ . Assuming independence between the two sets, the variance of the difference is equal to the sum of the variances (the difference is equal to the sum of the proportions after multiplying one by -1, which does not modify the variance). We conclude that the standard error equals

$$\text{se}_{\text{diff}} = \sqrt{\frac{\theta_A(1 - \theta_A)}{n_A} + \frac{\theta_B(1 - \theta_B)}{n_B}}. \quad (9.85)$$

- b) By the Gaussian approximation to the binomial distribution and the properties of Gaussian random variables, the difference between the proportions is approximately Gaussian with mean  $\theta_A - \theta_B$  and variance equal to the standard error  $\text{se}_{\text{diff}}$ . To compute the corresponding confidence interval we apply Lemma 9.42, setting  $\theta_A \approx k_A/n_A$  and  $\theta_B \approx k_B/n_B$  to obtain an approximation to the standard error that we denote by  $\text{se}_{\text{approx}}$ . The confidence interval equals

$$\left[ \frac{k_A}{n_A} - \frac{k_B}{n_B} - c_\alpha \text{se}_{\text{approx}}, \frac{k_A}{n_A} - \frac{k_B}{n_B} + c_\alpha \text{se}_{\text{approx}} \right], \quad c_\alpha := F_{\bar{z}}^{-1} \left( 1 - \frac{\alpha}{2} \right), \quad (9.86)$$

where  $F_{\bar{z}}$  denotes the cdf of a standard Gaussian with zero mean and unit variance.

- c) For the first trial, setting  $k_A := 52$ ,  $n_A := 100$ ,  $k_B := 30$  and  $n_B := 100$  we obtain a standard error  $\text{se}_{\text{approx}}$  equal to 0.0678 and a 0.95 confidence interval equal to

$$[0.52 - 0.3 - 1.96 \cdot 0.0678, 0.52 - 0.3 + 1.96 \cdot 0.0678] = [0.0871, 0.353]. \quad (9.87)$$

For the second trial, setting  $k_A := 30650$ ,  $n_A := 100000$ ,  $k_B := 30000$  and  $n_B := 100000$  we obtain a standard error equal to 0.00206 and a 0.95 confidence interval equal to

$$[0.3065 - 0.3 - 1.96 \cdot 0.00206, 0.3065 - 0.3 + 1.96 \cdot 0.00206] = [0.0025, 0.0105]. \quad (9.88)$$

- d) We can approximate the standard error via bootstrapping following Definition 9.50 as follows:

- Generate  $T$  batches of  $n$  bootstrap samples from Dataset  $A$ ,  $a_j^{[t]}$ ,  $1 \leq j \leq n_A$ ,  $1 \leq t \leq T$ , following Definition 9.49, where  $T$  is a very large integer. Let  $k_A^{[t]}$  denote the number of ones in batch  $t$ .
- Generate  $T$  batches of  $n$  bootstrap samples from Dataset  $B$ ,  $b_j^{[t]}$ ,  $1 \leq j \leq n_B$ ,  $1 \leq t \leq T$ , following Definition 9.49. Let  $k_B^{[t]}$  denote the number of ones in batch  $t$ .
- Compute the bootstrap difference in proportions for the batches:

$$d_t := \frac{k_A^{[t]}}{n_A} - \frac{k_B^{[t]}}{n_B}. \quad (9.89)$$

- Compute the sample standard deviation of the bootstrap differences  $d_1, \dots, d_T$  to obtain the bootstrap standard error  $\text{se}_{\text{bs}}$ .

- e) Option 1: We plug the bootstrap standard error  $\text{se}_{\text{bs}}$  into Lemma 9.42 to obtain the  $1 - \alpha$  bootstrap Gaussian confidence interval,

$$\left[ \frac{k_A}{n_A} - \frac{k_B}{n_B} - c_\alpha \text{se}_{\text{bs}}, \frac{k_A}{n_A} - \frac{k_B}{n_B} + c_\alpha \text{se}_{\text{bs}} \right], \quad c_\alpha := F_{\bar{z}}^{-1} \left( 1 - \frac{\alpha}{2} \right). \quad (9.90)$$

Option 2: We compute the  $\alpha/2$  and  $1 - \alpha/2$  quantiles  $q_{\alpha/2}$  and  $q_{1-\alpha/2}$  of the bootstrap differences  $d_1, \dots, d_T$ . The  $1 - \alpha$  bootstrap percentile confidence interval is  $[q_{\alpha/2}, q_{1-\alpha/2}]$ .

By consistency of the sample variance, the bootstrap standard error of the difference is approximately equal to the true standard error of the difference. Similarly, our approximation to the standard error is a consistent estimator of the true standard error (again because of consistency of the sample variance), so the bootstrap Gaussian confidence interval will be very similar to the confidence interval above. The bootstrap percentile confidence interval will also be similar, because the standard errors are approximately the same and the difference is approximately Gaussian as explained above, so we can apply the same argument as in the case of the sample mean (see beginning of Section 9.9.2).

9.12 (Exact confidence interval)

a) Since the variance of  $\tilde{m}$  equals  $\sigma^2/n$ ,

$$P\left(\mu \in \left[\tilde{m} - \frac{b}{\sqrt{\alpha n}}, \tilde{m} + \frac{b}{\sqrt{\alpha n}}\right]\right) = 1 - P\left(|\tilde{m} - \mu| > \frac{b}{\sqrt{\alpha n}}\right) \quad (9.91)$$

$$\begin{aligned} &\geq 1 - \frac{\alpha n \text{Var}[\tilde{m}]}{b^2} \quad \text{by Chebyshev's inequality} \\ &= 1 - \frac{\alpha \sigma^2}{b^2} \end{aligned} \quad (9.92)$$

$$\geq 1 - \alpha. \quad (9.93)$$

b) We can set  $b = 1$  because the parameter of interest is a ratio bounded by one. If  $\alpha = 0.05$ , setting the half-width equal to 0.01, yields

$$\frac{1}{\sqrt{0.05 n}} = 0.01 \quad (9.94)$$

so  $n$  must be greater than 200,000, which is 20 times more tests than required by the approximate confidence interval.

9.13 (Modified Gaussian bootstrap confidence interval) Let  $\tau$  denote the monotonic transformation and  $h(X)$  the observed value of the estimator of interest. Since the transformed parameter is approximately Gaussian, we use bootstrap samples to estimate the standard error  $\text{se}_{\text{bs}}$  of the transformed estimator (for each bootstrap batch we compute the estimator and apply the transformation, then we compute the sample standard deviation of the transformed values). We can use the bootstrap standard error to construct a confidence interval for the transformed parameter:

$$[\ell, u] = [\tau(h(X)) - c_\alpha \text{se}_{\text{bs}}, \tau(h(X)) + c_\alpha \text{se}_{\text{bs}}], \quad (9.95)$$

where  $c_\alpha$  is set to obtain a  $1-\alpha$  confidence interval under the Gaussian assumption. Now, if the transformed parameter of interest  $\tau(\tilde{w})$  is in  $[\ell, u]$  with probability  $1 - \alpha$ , then the parameter of interest  $\tilde{w}$  is in the interval  $[\tau^{-1}(\ell), \tau^{-1}(u)]$  also with probability  $1 - \alpha$ , because  $\tau^{-1}$  is a monotonic function (the inverse of a monotonic function is monotonic). Consequently, the interval

$$\mathcal{I}_{1-\alpha}^{[\tau]} := \left[\tau^{-1}(\tau(h(X)) - c_\alpha \text{se}_{\text{bs}}), \tau^{-1}(\tau(h(X)) + c_\alpha \text{se}_{\text{bs}})\right] \quad (9.96)$$

is a  $1-\alpha$  confidence interval for the parameter of interest.

---

## Hypothesis Testing

### Exercises

#### 10.1 (Test choice)

- a) Test 2. We want to control false positives at all costs, in order to avoid declaring the water safe when it is not and get the kids poisoned. The price to pay is low power, i.e. failure to declare water safe when it actually is.
- b) Test 1. Even if this means high significance level (so high probability of false positives), we want as much power as possible to not miss the gold!

#### 10.2 (Scout)

- a) The scout can apply a hypothesis test, to ensure that the workout is not consistent with the player having a 3-point percentage smaller or equal to 0.4. The parameter of interest  $\theta$  is the 3-point percentage of the player. The null hypothesis is that  $\theta \leq 0.4$ . The test statistic is the number of 3-pointers the player makes before missing. If the p value is below the desired significance level, she will recruit the player.
- b) We assume that the probability that the player makes each 3-point shot is  $\theta$  and that the shots are independent. Let  $\tilde{t}_\theta$  be a random variable that represents the number of made 3-pointers if the 3-point percentage is  $\theta$ . The probability that the player makes  $t$  3-point shots or more is the probability that they make the first  $t$  shots, so  $P(\tilde{t}_\theta \geq t) = \theta^t$ . Therefore, under our assumptions, the p-value function equals

$$\text{pv}(t) := \sup_{\theta \in [0, 0.4]} P(\tilde{t}_\theta \geq t) \quad (10.1)$$

$$= \sup_{\theta \in [0, 0.4]} \theta^t \quad (10.2)$$

$$= 0.4^t. \quad (10.3)$$

We need

$$\text{pv}(t) = 0.4^t \leq 0.05, \quad (10.4)$$

so  $t \geq \log 0.05 / \log 0.4 = 3.26$ . The player needs to make at least 4 3-pointers in a row to be recruited.

- c) Since the player is recruited if they make at least 4 3-pointers in a row, the probability that they are not recruited is  $1 - 0.5^4 = 0.9375$ . The scout is being very conservative and will definitely miss some good shooters.

#### 10.3 (Road renovation)

- a) The possible values of the test statistic are the integers between 0 and 4. Since the probability of an accident occurring in the *dangerous* section is  $1/4$  and the accidents

are independent, the pmf of the test statistic under the null hypothesis  $t_{\text{null}}$  is binomial with parameters  $n := 4$  and  $\theta := 1/4$ :

$$p_{\tilde{t}_{\text{null}}}(0) = \left(1 - \frac{1}{4}\right)^4 = \frac{81}{256}, \quad (10.5)$$

$$p_{\tilde{t}_{\text{null}}}(1) = 4 \left(\frac{1}{4}\right)^3 \left(1 - \frac{1}{4}\right) = \frac{108}{256}, \quad (10.6)$$

$$p_{\tilde{t}_{\text{null}}}(2) = 1 - \left(\frac{1}{4}\right)^4 = \frac{54}{256}, \quad (10.7)$$

$$p_{\tilde{t}_{\text{null}}}(3) = 4 \left(\frac{1}{4}\right)^3 \left(1 - \frac{1}{4}\right) = \frac{12}{256}, \quad (10.8)$$

$$p_{\tilde{t}_{\text{null}}}(4) = \left(\frac{1}{4}\right)^4 = \frac{1}{256}. \quad (10.9)$$

The p value function is

$$\text{pv}(0) := P(\tilde{t}_{\text{null}} \geq 0) = 1, \quad (10.10)$$

$$\text{pv}(1) := P(\tilde{t}_{\text{null}} \geq 1) = \frac{175}{256}, \quad (10.11)$$

$$\text{pv}(2) := P(\tilde{t}_{\text{null}} \geq 2) = \frac{67}{256}, \quad (10.12)$$

$$\text{pv}(3) := P(\tilde{t}_{\text{null}} \geq 3) = \frac{13}{256}, \quad (10.13)$$

$$\text{pv}(4) := P(\tilde{t}_{\text{null}} \geq 4) = \frac{1}{256}. \quad (10.14)$$

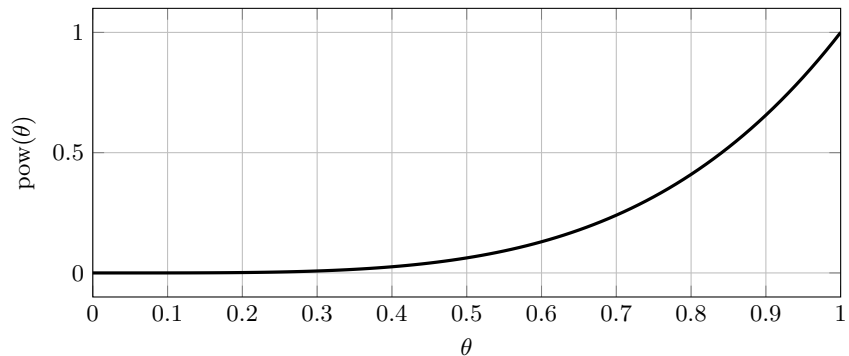
- b) The null hypothesis is rejected if  $\text{pv}(t) \leq \alpha := 0.05$ , which only holds for  $t = 4$  ( $\text{pv}(3) = 0.0508$ ). The probability of a false positive is therefore the probability of  $\tilde{t}_{\text{null}} = 4$ , which equals  $1/256 = 3.9 \cdot 10^{-3}$ .
- c) The null hypothesis is rejected if the test statistic equals 4. Under the assumption that accidents occur independently, the test statistic  $\tilde{t}_{\theta}$  is binomial with parameters  $n := 4$  and  $\theta$ . The power function consequently is

$$\text{pow}(\theta) := P(\text{pv}(\tilde{t}_{\theta}) \leq \alpha) \quad (10.15)$$

$$= P(\tilde{t}_{\text{null}} = 4) \quad (10.16)$$

$$= \theta^4. \quad (10.17)$$

Here is the graph:



- d) We need  $\theta^4 \geq 0.5$ , so  $\theta \geq 0.5^{1/4} = 0.84$ .

- 10.4 (Computer component) We denote the exponential random variables representing the data by  $\tilde{x}_1, \dots, \tilde{x}_n$ . The null hypothesis is that the mean of the exponential  $1/\lambda$  is smaller than one, or equivalently that  $\lambda \geq 1$ . For any  $\lambda$ ,

$$P(\tilde{t}_\lambda \geq t) = P\left(\min_{1 \leq i \leq n} \tilde{x}_i \geq t\right) \quad (10.18)$$

$$= P(\cap_{1 \leq i \leq n} \{\tilde{x}_i \geq t\}) \quad (10.19)$$

$$= \prod_{i=1}^n P(\tilde{x}_i \geq t) \quad (10.20)$$

$$= \prod_{i=1}^n \exp(-\lambda t) \quad (10.21)$$

$$= \exp(-\lambda n t). \quad (10.22)$$

Consequently, the p value function is

$$\text{pv}(t) := \sup_{0 \leq \lambda \leq 1} P(\tilde{t}_\lambda \geq t) \quad (10.23)$$

$$= \sup_{0 \leq \lambda \leq 1} \exp(-\lambda n t) \quad (10.24)$$

$$= \exp(-n t). \quad (10.25)$$

The null hypothesis is rejected when  $\text{pv}(t) = \exp(-n t) \leq \alpha$ , or equivalently

$$t \geq \frac{1}{n} \log\left(\frac{1}{\alpha}\right). \quad (10.26)$$

By (10.22), the power function then equals

$$\text{pow}(\lambda) := P(\text{pv}(\tilde{t}_\lambda) \leq \alpha) \quad (10.27)$$

$$= P\left(\tilde{t}_\lambda \geq \frac{1}{n} \log\left(\frac{1}{\alpha}\right)\right) \quad (10.28)$$

$$= \exp\left(-\lambda n \cdot \frac{1}{n} \log\left(\frac{1}{\alpha}\right)\right) \quad (10.29)$$

$$= \alpha^\lambda. \quad (10.30)$$

The power at  $\lambda = 1$  is  $\alpha$ . As  $\lambda \rightarrow 0$  it approaches one.

- 10.5 (Test statistic in two-sample test) Recall that  $\tilde{h}_i$  and  $\tilde{a}_i$  are independent Bernoulli random variables representing the  $i$ th free throw attempted at home and away respectively, for  $1 \leq i \leq n$ , so their parameters equal  $\theta_{\text{home}}$  and  $\theta_{\text{away}}$ , respectively. Consequently,  $\sum_{i=1}^{n_{\text{home}}} \tilde{h}_i$  is binomial with parameters  $n_{\text{home}}$  and  $\theta_{\text{home}}$ , as explained in the proof of Lemma 7.19. By the Gaussian approximation to the binomial distribution (Definition 9.39) and Theorem 3.31, this implies that  $\frac{1}{n_{\text{home}}} \sum_{i=1}^{n_{\text{home}}} \tilde{h}_i$  is approximately Gaussian with mean  $\theta_{\text{home}}$  and variance  $\frac{\theta_{\text{home}}(1-\theta_{\text{home}})}{n_{\text{home}}}$ . By the same argument,  $-\frac{1}{n_{\text{away}}} \sum_{i=1}^{n_{\text{away}}} \tilde{a}_i$  is approximately Gaussian with mean  $\theta_{\text{away}}$  and variance  $\frac{\theta_{\text{away}}(1-\theta_{\text{away}})}{n_{\text{away}}}$ . By Theorem 9.36,

$$\tilde{t}_\theta := \frac{1}{n_{\text{home}}} \sum_{i=1}^{n_{\text{home}}} \tilde{h}_i - \frac{1}{n_{\text{away}}} \sum_{i=1}^{n_{\text{away}}} \tilde{a}_i, \quad (10.31)$$

is therefore approximately Gaussian with mean  $\theta_{\text{home}} - \theta_{\text{away}}$  and variance

$$\frac{\theta_{\text{home}}(1-\theta_{\text{home}})}{n_{\text{home}}} + \frac{\theta_{\text{away}}(1-\theta_{\text{away}})}{n_{\text{away}}}. \quad (10.32)$$

## 10.6 (Tom Brady and hurricanes)

- a) In the  $n_W := 7$  years that Brady won, there were  $k_W := 4$  hurricanes. In the  $n_L := 13$  years that Brady didn't win, there were also  $k_L := 4$  hurricanes. The test statistic of the one-tailed test is

$$t_{\text{data}} := \frac{k_W}{n_W} - \frac{k_L}{n_L} \quad (10.33)$$

$$= \frac{4}{7} - \frac{4}{13} \quad (10.34)$$

$$= 0.264. \quad (10.35)$$

Let  $k := k_W + k_L = 8$  and  $n := n_W + n_L = 20$ . The variance under the null hypothesis is

$$\sigma_{\text{null}}^2 := \frac{k(n-k)}{n^2} \left( \frac{1}{n_W} + \frac{1}{n_L} \right) \quad (10.36)$$

$$= \frac{8 \cdot 12}{20^2} \left( \frac{1}{7} + \frac{1}{13} \right) \quad (10.37)$$

$$= 0.0257. \quad (10.38)$$

The p value equals

$$\text{pv}(t_{\text{data}}) = 1 - F_{\tilde{z}} \left( \frac{t_{\text{data}}}{\sigma_{\text{null}}} \right) \quad (10.39)$$

$$= 0.125. \quad (10.40)$$

- b) No, because the null hypothesis and test statistic were not decided before looking at the data.

## 10.7 (Unemployment in Spain)

- a) Our null hypothesis is that the population correlation coefficient  $\rho$  is zero. Let  $\tilde{r}$  be a random variable that represents the sample correlation coefficient computed from  $n$  samples under the null hypothesis. After applying Fisher's transformation, which we denote by  $M$ ,  $M(\tilde{r})$  is approximately Gaussian with mean  $M(\rho) = M(0) = 0$  and standard deviation  $1/\sqrt{n-3}$ .

We select the absolute value of the observed sample correlation coefficient  $|r|$  as the test statistic. Let  $\tilde{z}$  be a standard Gaussian random variable with mean 0 and unit variance, then by monotonicity of  $M$  the p-value function is

$$\text{P}(|\tilde{r}| \geq |r|) = \text{P}(|M(\tilde{r})| \geq |M(r)|) \quad (10.41)$$

$$\approx \text{P} \left( \frac{|\tilde{z}|}{\sqrt{n-3}} \geq |M(r)| \right) \quad (10.42)$$

$$= 2 \left( 1 - F_{\tilde{z}}(\sqrt{n-3}|M(r)|) \right) \quad (10.43)$$

$$= 2 \left( 1 - F_{\tilde{z}} \left( \left| \frac{\sqrt{n-3}}{2} \log \left( \frac{1+t}{1-t} \right) \right| \right) \right). \quad (10.44)$$

- b) Setting  $t := -0.21$  and  $n := 92$ , the p value equals

$$2 \left( 1 - F_{\tilde{z}} \left( \left| \frac{\sqrt{89}}{2} \log \left( \frac{0.79}{1.21} \right) \right| \right) \right) = 2(1 - F_{\tilde{z}}(2.01)) \quad (10.45)$$

$$= 0.044. \quad (10.46)$$

The correlation between unemployment and temperature in Spain is statistically significant at a significance level of 0.05.

- c) The data are not sampled i.i.d. from a population, they just include all the months over a certain period (with some exclusions due to the COVID-19 pandemic). The hypothesis test indicates that if they were i.i.d. samples, then it would be very unlikely for them to be uncorrelated. This is convincing evidence that the correlation between temperature and unemployment is not due to small sample size, but we can definitely *not* conclude that higher temperatures cause a decrease in unemployment, as explained in detail in Section 8.8.

- 10.8 (Median cholesterol) The test statistic is the median cholesterol of the ill people minus the median cholesterol of the healthy people. The value of the test statistic for each permutation is:

HHHHH: 250-240=10

HHHHH: 260-240 = 20

IIHHHI: 250-260=-10

IHHIHI: 300-250=50

HHHHH: 260-250=10.

The actual test statistic equals 300-250=50. The p value is the fraction of the time the test statistic is larger or equal to this for the different permutations. It equals 1/5.

- 10.9 (False positives of a permutation test)

- a) By Theorem 10.22, the p-value function equals

$$pv(t) = \frac{\sum_{v \in \Pi_{x_{\text{data}}}} 1(T(v) \geq t)}{|\Pi_{x_{\text{data}}}|}. \quad (10.47)$$

In our case, there are 24 permutations. For each possible value of the test statistic  $t$  (see Figure 10.7), we compute what fraction of permutations correspond to a test statistic greater or equal to  $t$ . The resulting values of the p value function are listed below, to the right of the corresponding value of  $t$ :

$t$	-4.5	-1.5	-0.5	0.5	1.5	4.5
$pv(t)$	1	$\frac{5}{6}$	$\frac{2}{3}$	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{6}$

- b) Under the null hypothesis that the entries of  $\tilde{x}_{\text{null}}$  are exchangeable, every permutation is equally likely: each occurs with probability 1/24. Since each value of the test statistic corresponds to 4 permutations (see Figure 10.7), the conditional probability that it occurs is  $4/24 = 1/6$ . Consequently, that is also the conditional probability of each of the possible values of  $\tilde{u}$ , because for any  $t$  and  $u$  such that  $pv(t) = u$ ,

$$P(\tilde{u} = u \mid \tilde{x}_{\text{null}} \in \Pi_{x_{\text{data}}}) = P(\tilde{t}_{\text{null}} = t \mid \tilde{x}_{\text{null}} \in \Pi_{x_{\text{data}}}). \quad (10.48)$$

To compute the conditional cdf of  $\tilde{u}$  given  $\tilde{x}_{\text{null}} \in \Pi_{x_{\text{data}}}$ , we sum the corresponding probabilities, obtaining:

$u$	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{2}{3}$	$\frac{5}{6}$	1
$F_{\tilde{u}}(u \mid \tilde{x}_{\text{null}} \in \Pi_{x_{\text{data}}})$	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{2}{3}$	$\frac{5}{6}$	1

Notice that this is consistent with Theorem 10.27.

- c) A false positive happens if  $\tilde{u} \leq 0.2$ , which can only happen if  $\tilde{u} = \frac{1}{6}$ . The conditional probability of a false positive is therefore equal to  $\frac{1}{6}$ , which is smaller or equal to the significance level in accordance with Theorem 10.27.

- d) Yes, as established in Corollary 10.28. Intuitively, the marginal probability of false positives is controlled by the significance level  $\alpha$ , because the conditional probability is bounded by  $\alpha$  for every possible value of the data.

#### 10.10 (Sign test)

- a) Since  $\tilde{a}$  and  $\tilde{b}$  are continuous,  $P(\tilde{a} = \tilde{b}) = 0$ , so any outcome in the probability space associated to the random variables must belong to one of the two disjoint events  $\tilde{a} > \tilde{b}$  and  $\tilde{b} > \tilde{a}$ . Consequently the sum of their probabilities must equal one

$$1 = P(\{\tilde{a} > \tilde{b}\} \cup \{\tilde{b} > \tilde{a}\}) = P(\tilde{a} > \tilde{b}) + P(\tilde{b} > \tilde{a}). \quad (10.49)$$

Under the assumptions, these two probabilities are the same. Since  $f_{\tilde{a}, \tilde{b}}(a, b) = f_{\tilde{a}}(a)f_{\tilde{a}}(b)$ :

$$P(\tilde{a} > \tilde{b}) = \int_{b=-\infty}^{\infty} \int_{a=b}^{\infty} f_{\tilde{a}, \tilde{b}}(a, b) \, da \, db \quad (10.50)$$

$$= \int_{b=-\infty}^{\infty} \int_{a=b}^{\infty} f_{\tilde{a}}(a)f_{\tilde{a}}(b) \, da \, db \quad (10.51)$$

$$= \int_{a=-\infty}^{\infty} \int_{b=a}^{\infty} f_{\tilde{a}}(b)f_{\tilde{a}}(a) \, da \, db \quad (10.52)$$

$$= \int_{a=-\infty}^{\infty} \int_{b=a}^{\infty} f_{\tilde{a}, \tilde{b}}(a, b) \, da \, db \quad (10.53)$$

$$= P(\tilde{b} > \tilde{a}). \quad (10.54)$$

We conclude that  $P(\tilde{a} > \tilde{b}) = 1/2$ .

- b) Let  $\tilde{x}_i$  and  $\tilde{y}_i$  denote the first and second entry for each data point. By the previous question,  $1(\tilde{y}_i > \tilde{x}_i)$  is a Bernoulli random variable with parameter  $1/2$ . Consequently, under the null hypothesis the test statistic is a binomial random variable with parameters  $n$  and  $\theta := 1/2$ .
- c) The test statistic equals 7. Let  $\tilde{t}_{\text{null}}$  denote its distribution under the null hypothesis. The p value equals

$$\text{pv}(7) := P(\tilde{t}_{\text{null}} \geq 7) \quad (10.55)$$

$$= \frac{1}{2^{10}} \sum_{i=7}^{10} \binom{10}{i} \quad (10.56)$$

$$= 0.172. \quad (10.57)$$

- d) The sign test is still valid because under these assumptions  $1(\tilde{y}_i > \tilde{x}_i)$  is still a Bernoulli random variable with parameter  $1/2$

#### 10.11 (Drug design)

- a) Let  $\tilde{t}_{\text{null}}$  denote the test statistic under the null hypothesis. The p value function equals

$$\text{pv}(t) := P(\tilde{t}_{\text{null}} \geq t) \quad (10.58)$$

$$= \int_{a=t}^{\infty} \exp(-a) \, da \quad (10.59)$$

$$= \exp(-t). \quad (10.60)$$

If  $\tilde{t}_{\theta}$  is the test statistic when the drug is effective, the corresponding p value equals



$\widetilde{\text{pv}}_\theta := \text{pv}(\tilde{t}_\theta) = \exp(-\tilde{t}_\theta)$ . For  $0 \leq b \leq 1$ , the cdf of this random variable is

$$F_{\widetilde{\text{pv}}_\theta}(b) = \text{P}(\widetilde{\text{pv}}_\theta \leq b) \quad (10.61)$$

$$= \text{P}(\exp(-\tilde{t}_\theta) \leq b) \quad (10.62)$$

$$= \text{P}\left(\tilde{t}_\theta \geq \log\left(\frac{1}{b}\right)\right) \quad (10.63)$$

$$= \int_{a=\log(\frac{1}{b})}^{\infty} \theta \exp(-\theta a) \, da \quad (10.64)$$

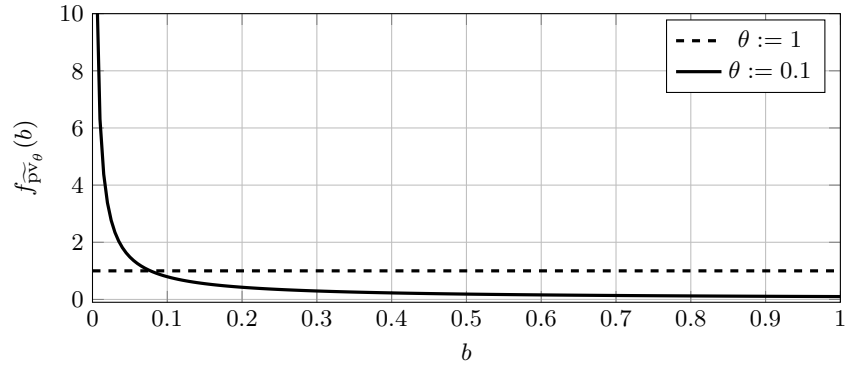
$$= \exp(-\theta \log(\frac{1}{b})) \quad (10.65)$$

$$= b^\theta. \quad (10.66)$$

The corresponding pdf equals

$$f_{\widetilde{\text{pv}}_\theta}(b) = \theta b^{\theta-1}, \quad (10.67)$$

for  $0 \leq b \leq 1$  and zero otherwise. For  $\theta := 1$  unsurprisingly, the distribution is uniform, since this corresponds to the null hypothesis. The plot for  $\theta := 0.1$  is shown below.



- b) Let  $\widetilde{\text{pv}}$  represent the p value and let  $\tilde{d}$  be a Bernoulli random variable that equals one if the drug is effective and zero otherwise. From the previous question we know that:

$$f_{\widetilde{\text{pv}}|\tilde{d}}(b|0) = f_{\widetilde{\text{pv}}_0}(b) = 1, \quad (10.68)$$

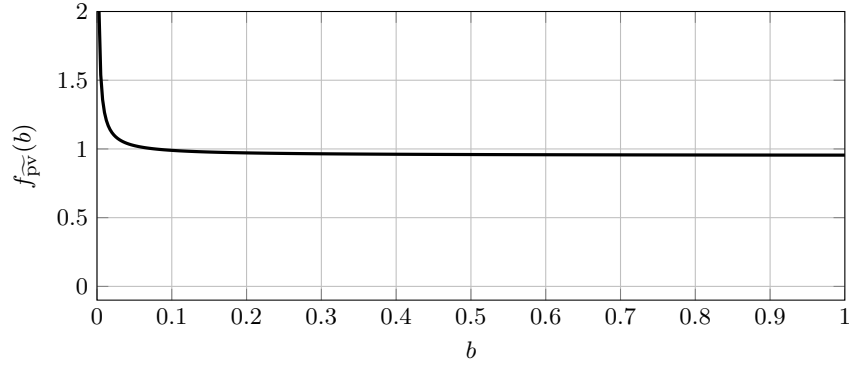
$$f_{\widetilde{\text{pv}}|\tilde{d}}(b|1) = f_{\widetilde{\text{pv}}_\theta}(b) = \theta b^{\theta-1}, \quad (10.69)$$

for  $0 \leq b \leq 1$  and zero otherwise. The pdf of the p value therefore equals

$$f_{\widetilde{\text{pv}}}(b) = \sum_{b=0}^1 p_{\tilde{d}}(b) f_{\widetilde{\text{pv}}|\tilde{d}}(b|0) \quad (10.70)$$

$$= 0.95 + 0.05\theta b^{\theta-1}, \quad (10.71)$$

for  $0 \leq b \leq 1$  and zero otherwise. Here is the pdf for  $\theta := 0.1$ :



- c) We reject the null hypothesis if the p value  $\widetilde{p}\widetilde{v} \leq \alpha$ . A false positive occurs if we reject when  $\tilde{d} = 0$ . The conditional probability of a false positive given that we reject the null hypothesis is therefore equal to

$$P(\tilde{d} = 0 | \widetilde{p}\widetilde{v} \leq \alpha) = \frac{P(\tilde{d} = 0, \widetilde{p}\widetilde{v} \leq \alpha)}{P(\widetilde{p}\widetilde{v} \leq \alpha)}. \quad (10.72)$$

We have

$$P(\tilde{d} = 0, \widetilde{p}\widetilde{v} \leq \alpha) = P(\tilde{d} = 0) P(\widetilde{p}\widetilde{v} \leq \alpha | \tilde{d} = 0) \quad (10.73)$$

$$= 0.95 F_{\widetilde{p}\widetilde{v}_0}(\alpha) \quad (10.74)$$

$$= 0.95\alpha. \quad (10.75)$$

From (b) we know that the conditional cdf of  $\widetilde{p}\widetilde{v}$  given  $\tilde{d} = 1$  is

$$F_{\widetilde{p}\widetilde{v} | \tilde{d}}(b | 1) = b^\theta, \quad (10.76)$$

so

$$P(\widetilde{p}\widetilde{v} \leq \alpha) = \sum_{d=0}^1 P(\tilde{d} = d) P(\widetilde{p}\widetilde{v} \leq \alpha | \tilde{d} = d) \quad (10.77)$$

$$= 0.95\alpha + 0.05 F_{\widetilde{p}\widetilde{v} | \tilde{d}}(\alpha | 1) \quad (10.78)$$

$$= 0.95\alpha + 0.05\alpha^\theta. \quad (10.79)$$

We conclude that

$$P(\tilde{d} = 0 | \widetilde{p}\widetilde{v} \leq \alpha) = \frac{0.95\alpha}{0.95\alpha + 0.05\alpha^\theta}. \quad (10.80)$$

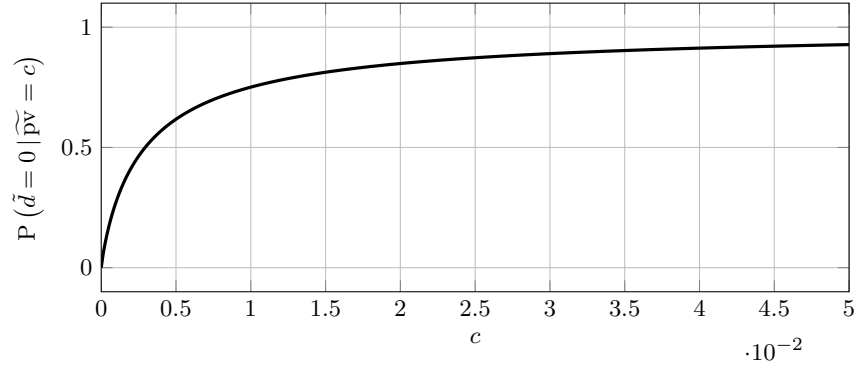
For  $\alpha = 0.05$  and  $\theta := 0.1$  the probability equals 0.562. More than half of the findings are false positives.

- d) By Bayes rule for discrete and continuous random variables,

$$P(\tilde{d} = 0 | \widetilde{p}\widetilde{v} = c) = \frac{p_{\tilde{d}}(0) f_{\widetilde{p}\widetilde{v} | \tilde{d}}(c | 0)}{f_{\widetilde{p}\widetilde{v}}(c)} \quad (10.81)$$

$$= \frac{0.95 \cdot 1}{0.95 + 0.05\theta c^{\theta-1}}. \quad (10.82)$$

Here is the plot for  $\theta := 0.1$ :



The probability increases monotonically with the p value, which suggests doing follow-up testing based on the associated p value, starting with the drug that has the smallest p value.

#### 10.12 (P-hacking)

- a) For  $u \leq \alpha$ , the conditional cdf equals

$$F_{\tilde{u} | \tilde{u} \leq \alpha}(u) = P(\tilde{u} \leq u | \tilde{u} \leq \alpha) \quad (10.83)$$

$$= \frac{P(\tilde{u} \leq u, \tilde{u} \leq \alpha)}{P(\tilde{u} \leq \alpha)} \quad (10.84)$$

$$= \frac{P(\tilde{u} \leq u)}{P(\tilde{u} \leq \alpha)} \quad (10.85)$$

$$= \frac{F_{\tilde{u}}(u)}{F_{\tilde{u}}(\alpha)} \quad (10.86)$$

$$= \frac{u}{\alpha}. \quad (10.87)$$

For  $u < 0$ ,  $F_{\tilde{u} | \tilde{u} \leq \alpha}(u) = 0$  and for  $u > \alpha$ ,  $F_{\tilde{u} | \tilde{u} \leq \alpha}(u) = 1$ . The pdf is therefore equal to  $1/\alpha$  between 0 and  $\alpha$ . The conditional distribution is uniform between 0 and  $\alpha$ .

- b) The histogram is exactly equal to the continuous pdf derived in the previous question for  $\alpha := 0.05$ . Since for continuous test statistics, the p value is uniformly distributed between 0 and 1, this is the distribution we would see if the null hypothesis always holds, but only p values smaller than  $\alpha$  are reported. This suggests that the publications could be based on meaningless results that are only statistically significant by chance.
- c) If our conjecture holds, the group is only reporting 5% of their results, which means that there would be approximately 2,000 total results, and therefore 1,900 unpublished ones.

# Principal Component Analysis And Low-Rank Models

## Exercises

### 11.1 (Random vector)

a) The covariance matrix equals

$$\Sigma_{\tilde{x}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0.25 & 0.25 \\ 0 & 0.25 & 0.25 \end{bmatrix}, \quad (11.1)$$

so  $\text{Var}[\tilde{x}_1] = 1$ ,  $\text{Var}[\tilde{x}_2] = 0.25$ ,  $\text{Var}[\tilde{x}_3] = 0.25$ .

- b) The maximum variance in any direction is given by the largest eigenvalue, which is equal to 1. There cannot be another direction with higher variance.
- c) Let  $\tilde{y} := a_1\tilde{x}_1 + a_2\tilde{x}_2 + a_3\tilde{x}_3$ . By linearity of expectation,  $\mathbb{E}[\tilde{y}] = 0$ . By Chebyshev's inequality, if  $\text{Var}[\tilde{y}] = 0$  then  $\text{P}(\tilde{y} \neq 0) = 1$ , which is exactly what we want. According to the eigendecomposition, the variance is zero in the direction of the third eigenvector, so setting  $a_1 = 0$ ,  $a_2 = 1/\sqrt{2}$ , and  $a_3 = -1/\sqrt{2}$  does the trick.

Geometrically, the random vector lies in the plane orthogonal to the vector

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix} \quad (11.2)$$

with probability one.

### 11.2 (Basketball team) Let $\tilde{x}_A$ , $\tilde{x}_B$ and $\tilde{x}_C$ denote the points scored by each player. We have

$$\text{Var}[\tilde{x}_A + \tilde{x}_B] = \text{Var}[\tilde{x}_A] + \text{Var}[\tilde{x}_B] + 2\text{Cov}[\tilde{x}_A, \tilde{x}_B] \quad (11.3)$$

$$= 100 + 81 - 160 \quad (11.4)$$

$$= 21, \quad (11.5)$$

$$\text{Var}[\tilde{x}_A + \tilde{x}_C] = \text{Var}[\tilde{x}_A] + \text{Var}[\tilde{x}_C] + 2\text{Cov}[\tilde{x}_A, \tilde{x}_C] \quad (11.6)$$

$$= 100 + 100 + 20 \quad (11.7)$$

$$= 220, \quad (11.8)$$

$$\text{Var}[\tilde{x}_B + \tilde{x}_C] = \text{Var}[\tilde{x}_B] + \text{Var}[\tilde{x}_C] + 2\text{Cov}[\tilde{x}_B, \tilde{x}_C] \quad (11.9)$$

$$= 81 + 100 + 100 \quad (11.10)$$

$$= 281. \quad (11.11)$$

If the team is winning by a lot, it is probably a good idea to minimize the variance, so we would recommend playing A and B. If they are losing by a lot, they should maximize the variance to try to win, so we would recommend playing B and C (but note that this could result in an even more humiliating loss!).

- 11.3 (Correlation and PCA) If the correlation coefficient is one, then the covariance equals the product of the standard deviations  $\text{Cov}(\tilde{x}[1], \tilde{x}[2]) = \sigma_{\tilde{x}[1]} \sigma_{\tilde{x}[2]}$ . The covariance matrix equals

$$\Sigma_{\tilde{x}} = \begin{bmatrix} \sigma_{\tilde{x}[1]}^2 & \sigma_{\tilde{x}[1]} \sigma_{\tilde{x}[2]} \\ \sigma_{\tilde{x}[1]} \sigma_{\tilde{x}[2]} & \sigma_{\tilde{x}[2]}^2 \end{bmatrix} \quad (11.12)$$

$$= \begin{bmatrix} \sigma_{\tilde{x}[1]} \\ \sigma_{\tilde{x}[2]} \end{bmatrix} \begin{bmatrix} \sigma_{\tilde{x}[1]} & \sigma_{\tilde{x}[2]} \end{bmatrix}. \quad (11.13)$$

The rank of the matrix is one. The first eigenvalue equals  $\sigma_{\tilde{x}[1]}^2 + \sigma_{\tilde{x}[2]}^2$  and the second equals zero. The variance of the second principal component is therefore zero. Intuitively, both components are linearly dependent, so the data lie along a line. The variance along the line is the variance of the first principal component. The variance orthogonal to the line, which is the variance of the second principal component, is zero because all the points are on the line.

- 11.4 (Principal components are uncorrelated) Let  $X$  be a dataset containing  $n$  vectors  $x_1, x_2, \dots, x_n \in \mathbb{R}^d$  with  $d$  features each. The sample mean of each principal component is zero, for  $1 \leq j \leq d$ ,

$$\frac{1}{n} \sum_{i=1}^n w_j[i] = \frac{1}{n} \sum_{i=1}^n u_j^T \text{ct}(x_i) \quad (11.14)$$

$$= u_j^T \left( \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{n} \sum_{i=1}^n \frac{1}{n} \sum_{l=1}^n x_l \right) = 0, \quad (11.15)$$

where  $u_1, \dots, u_d$  are the eigenvectors of the sample covariance matrix  $\Sigma_X$ . Consequently, the sample covariance between two principal components ( $j \neq k$ ) equals

$$\frac{1}{n-1} \sum_{i=1}^n w_j[i] w_k[i] = \frac{1}{n-1} \sum_{i=1}^n u_j^T \text{ct}(x_i) \text{ct}(x_i)^T u_k \quad (11.16)$$

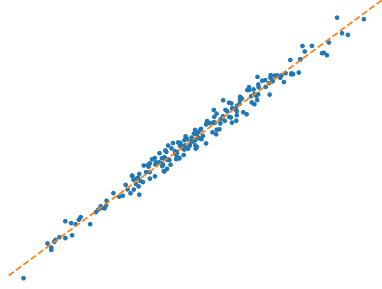
$$= u_j^T \Sigma_X u_k \quad (11.17)$$

$$= \lambda_k u_j^T u_k = 0, \quad (11.18)$$

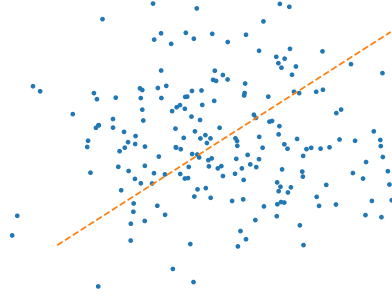
because  $u_j$  and  $u_k$  are orthogonal by the spectral theorem.

- 11.5 (Estimating a direction)

- a) To produce the samples we can first sample from  $\tilde{x}$  and multiply the value by  $v$ . This will produce a point on the line collinear with  $v$  at an average distance from the origin of  $\sigma_{\text{signal}}$ . This is then added to a sample from a Gaussian with independent entries centered at the origin with standard deviation  $\sigma_{\text{noise}}$ . If  $\sigma_{\text{signal}}$  is much larger than  $\sigma_{\text{noise}}$  then the samples will lie close to  $v$ . They look like this (the dashed line is collinear with  $v$ ):



- b) If  $\sigma_{\text{noise}}$  is much larger than  $\sigma_{\text{signal}}$  then the noise component dominates, and the samples will not lie as close to  $v$ . They look like this (the dashed line is collinear with  $v$ ):



- c) No. By linearity of expectation, the mean of the random variable that represents the data equals

$$\mathbb{E}[\tilde{y}] = \mathbb{E}[\tilde{x}v + \tilde{z}] \quad (11.19)$$

$$= \mathbb{E}[\tilde{x}]v + \mathbb{E}[\tilde{z}] \quad (11.20)$$

$$= 0, \quad (11.21)$$

so we won't be able to estimate  $v$  by averaging.

- d) By linearity of expectation and the fact that the mean of  $\tilde{x}$  and  $\tilde{z}$  is zero,

$$\mathbb{E}[\tilde{y}\tilde{y}^T] = \mathbb{E}[(\tilde{x}v + \tilde{z})(\tilde{x}v + \tilde{z})^T] \quad (11.22)$$

$$= \mathbb{E}[\tilde{x}^2]vv^T + \mathbb{E}[\tilde{z}\tilde{z}^T] \quad (11.23)$$

$$= \sigma_{\text{signal}}^2 vv^T + \sigma_{\text{noise}}^2 I. \quad (11.24)$$

- e) Setting  $U := [v \ u_2 \ \cdots \ u_d]$ , we have  $UU^T = I$  by the orthonormality assumption. We can therefore rewrite the covariance matrix as

$$\mathbb{E}[\tilde{y}\tilde{y}^T] = \sigma_{\text{signal}}^2 vv^T + \sigma_{\text{noise}}^2 UU^T \quad (11.25)$$

$$= U \begin{bmatrix} \sigma_{\text{signal}}^2 + \sigma_{\text{noise}}^2 & 0_{d-1} \\ 0_{d-1} & \sigma_{\text{noise}}^2 I_{d-1} \end{bmatrix} U^T, \quad (11.26)$$

where  $0_{d-1}$  is a vector of  $d-1$  zeros and  $I_{d-1}$  is the  $d-1 \times d-1$  identity matrix.

- f) Apply principal component analysis (i.e. compute the eigendecomposition of the sample covariance matrix) and use the first principal direction (i.e. the eigenvector with the largest eigenvalue) as an estimate of the direction of  $v$ .

#### 11.6 (Normalization)

- a) We have  $\tilde{y} = A\tilde{x}$ , where

$$A := \begin{bmatrix} \frac{1}{10} & 0 & 0 \\ 0 & \frac{1}{20} & 0 \\ 0 & 0 & \frac{1}{0.4} \end{bmatrix}, \quad (11.27)$$

so

$$\Sigma_{\tilde{y}} := \mathbb{E} [\tilde{y}\tilde{y}^T] \quad (11.28)$$

$$= \mathbb{E} [A\tilde{x}\tilde{x}^T A^T] \quad (11.29)$$

$$= A\Sigma_{\tilde{x}}A^T \quad (11.30)$$

$$= \begin{bmatrix} 1 & 0.125 & 0 \\ 0.125 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (11.31)$$

- b) No. For example, in the direction of the vector

$$v := [1/\sqrt{2} \quad 1/\sqrt{2} \quad 0] \quad (11.32)$$

the directional variance of  $\tilde{y}$  equals

$$\text{Var} [v^T \tilde{y}] = v^T \Sigma_{\tilde{y}} v \quad (11.33)$$

$$= 1.125. \quad (11.34)$$

- c) The eigendecomposition of  $\Sigma_{\tilde{x}}$  yields eigenvalues equal to 402, 97.9 and 0.16, corresponding to the eigenvectors

$$u_1(\tilde{x}) = \begin{bmatrix} -0.997 \\ 0.082 \\ 0 \end{bmatrix}, \quad u_2(\tilde{x}) = \begin{bmatrix} -0.082 \\ -0.997 \\ 0 \end{bmatrix}, \quad u_3(\tilde{x}) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad (11.35)$$

respectively. To reduce dimensionality we compute the inner product with  $u_1(\tilde{x})$  and  $u_2(\tilde{x})$ .

The eigendecomposition of  $\Sigma_{\tilde{y}}$  yields eigenvalues equal to 1.125, 1 and 0.875, corresponding to the eigenvectors

$$u_1(\tilde{y}) = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}, \quad u_2(\tilde{y}) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \quad u_3(\tilde{y}) = \begin{bmatrix} -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix} \quad (11.36)$$

respectively. To reduce dimensionality we compute the inner product with  $u_1(\tilde{y})$  and  $u_2(\tilde{y})$ .

- d) (1) Using  $\tilde{y}$  is probably better because the variances and covariances correspond to different units that are not directly comparable. Normalizing makes the analysis independent of the units. The third feature is height in meters so it is normal that its variance is smaller than one. This does not mean that it is negligible.

(2) Using  $\tilde{x}$  is probably better. The three features are dimensions of cars in the same unit. The variance of the height is very small compared to the other two, which indicates that the height of the cars does not vary at all. It therefore makes sense to ignore it while reducing the dimensionality.

11.7 (PCA and sample variance) Let  $X := \{x_1, \dots, x_n\}$  denote the dataset of interest. We want to establish that the first  $k$  eigenvectors of the sample covariance matrix  $\Sigma_X$  of  $X$ , associated to the  $k$  largest eigenvalues, satisfy

$$\{u_1, \dots, u_k\} = \arg \max_{\substack{\{b_1, \dots, b_k\} \\ \|b_j\|_2=1 \ 1 \leq j \leq k \\ b_j \perp b_k \text{ for } j \neq k}} \sum_{j=1}^k v(X_{b_j}), \quad (11.37)$$

where  $v(X_{b_j})$  denotes the sample variance of the set  $X_{b_j} := \{b_j^T x_1, \dots, b_j^T x_n\}$ . We prove this mimicking the proof of the theorem that establishes the analogous result for random vectors, which leverages induction on  $k$ . The base case  $k := 1$  follows immediately from the fact that the first eigenvector maximizes the sample variance,

$$u_1 = \arg \max_{\|b\|_2=1} v(X_b). \quad (11.38)$$

To complete the proof we establish that if the induction hypothesis

$$\{u_1, \dots, u_{k-1}\} = \arg \max_{\substack{\{b_1, \dots, b_{k-1}\} \\ \|b_j\|_2=1 \ 1 \leq j \leq k-1 \\ b_j \perp b_{k-1} \text{ for } j \neq k-1}} \sum_{j=1}^{k-1} v(X_{b_j}) \quad (11.39)$$

holds, then (11.37) holds. We now set  $b_1, \dots, b_k$  to be an arbitrary fixed set of  $k$  orthonormal vectors, and show that they cannot capture more variance than  $u_1, \dots, u_k$  if the induction hypothesis holds. Consider the subspace  $\mathcal{S} := \text{span}(b_1, \dots, b_k)$  spanned by  $b_1, \dots, b_k$ . We are interested in the projection of the data onto this subspace, because

$$\sum_{j=1}^k v(X_{b_j}) = \sum_{j=1}^k \frac{1}{n-1} \sum_{i=1}^n (b_j^T x_i)^2 \quad (11.40)$$

$$= \frac{1}{n-1} \sum_{i=1}^n \left\| \sum_{j=1}^k b_j^T x_i b_j \right\|_2^2 \quad (11.41)$$

$$= \frac{1}{n-1} \sum_{i=1}^n \|\mathcal{P}_{\mathcal{S}} x_i\|_2^2. \quad (11.42)$$

We can represent this projection using an arbitrary orthonormal basis of  $\mathcal{S}$ . Indeed, any set of  $k$  orthonormal vectors  $a_1, \dots, a_k$  spanning  $\mathcal{S}$  satisfy

$$\mathcal{P}_{\mathcal{S}} x_i := \sum_{j=1}^k b_j^T x_i b_j = \sum_{j=1}^k a_j^T x_i a_j. \quad (11.43)$$

The key is to choose the basis wisely.  $\mathcal{S}$  has dimension  $k$ , so it must contain at least one vector  $a_{\perp}$  that is orthogonal to  $u_1, u_2, \dots, u_{k-1}$ . By the properties of the principal directions, the sample variance in that direction cannot be higher than in the  $k$ th principal direction

$$v(X_{u_k}) \geq v(X_{a_{\perp}}). \quad (11.44)$$



We build our wisely-chosen orthonormal basis  $a_1, a_2, \dots, a_k$  for  $\mathcal{S}$  setting  $a_k := a_\perp$  (we can construct such a basis via the Gram-Schmidt process, starting with  $a_\perp$ ). By the induction hypothesis,

$$\sum_{j=1}^{k-1} v(X_{u_j}) \geq \sum_{j=1}^{k-1} v(X_{a_j}). \quad (11.45)$$

Combining (11.45) and (11.44) yields

$$\sum_{j=1}^k v(X_{u_j}) \geq \sum_{j=1}^k v(X_{a_j}) \quad (11.46)$$

$$= \frac{1}{n-1} \sum_{i=1}^n \|\mathcal{P}_{\mathcal{S}} x_i\|_2^2 \quad (11.47)$$

$$= \sum_{j=1}^k v(X_{b_j}). \quad (11.48)$$

Since this holds for any choice of  $b_1, \dots, b_k$  the proof is complete.

- 11.8 (Whitening) If the entries of  $\tilde{y}$  are all uncorrelated and have unit variance then  $\tilde{y}$  has a covariance matrix equal to the identity. The covariance matrix equals

$$\Sigma_{\tilde{y}} := \mathbb{E}[\tilde{y}\tilde{y}^T] = \mathbb{E}[A\tilde{x}\tilde{x}^T A^T] \quad (11.49)$$

$$= A\Sigma_{\tilde{x}}A^T. \quad (11.50)$$

Now in order for this product to equal the identity we build  $A$  based on the eigendecomposition of  $\Sigma_{\tilde{x}} = U\Lambda U^T$ . Since covariance matrices are positive semidefinite, and  $\Sigma_{\tilde{x}}$  is full rank, the entries of the diagonal matrix  $\Lambda$  are positive, so we can define  $\sqrt{\Lambda}$ , which is a diagonal matrix that contains the square root of the entries of  $\Lambda$ . Setting  $A := \sqrt{\Lambda}^{-1}U^T$ , since  $U$  is an orthogonal matrix by the spectral theorem so that  $U^T U = I$ ,

$$\Sigma_{\tilde{y}} = A\Sigma_{\tilde{x}}A^T \quad (11.51)$$

$$= \sqrt{\Lambda}^{-1}U^T U \Lambda U^T U \sqrt{\Lambda}^{-1} \quad (11.52)$$

$$= \sqrt{\Lambda}^{-1} \Lambda \sqrt{\Lambda}^{-1} = I. \quad (11.53)$$

- 11.9 (Numerical rank) Let us assume that  $M$  has dimensions  $n_1 \times n_2$  and  $n_1 \leq n_2$  (this is without loss of generality, because we can always consider the transpose of  $M$  if  $n_2 < n_1$ ). Let  $s_1 > s_2 > \dots > s_{n_1}$  be the singular values of  $M$ . The numerical rank  $r$  is the smallest positive integer such that

$$\sum_{i=r+1}^n s_i^2 \leq \epsilon^2. \quad (11.54)$$

To compute it, we sum the squared singular values starting from the smallest one until the sum reaches  $\epsilon^2$ . Consider the rank- $r$  matrix

$$L_{\text{SVD}} := \sum_{l=1}^r s_l u_l v_l^T, \quad (11.55)$$

obtained by truncating the SVD of  $M$ . As shown in the notes (eq. 11.162):

$$\|M - L_{\text{SVD}}\|_F = \sqrt{\sum_{l=r+1}^{n_1} s_l^2}, \quad (11.56)$$

which is smaller or equal to  $\epsilon$  by the definition of  $r$ . Consequently, there is a rank  $r$  matrix such that the distance to  $M$  is smaller or equal to  $\epsilon$ .

Now we show that there cannot be a rank- $r'$  matrix  $L'$  with  $r' < r$ , such that  $\|L' - M\|_F \leq \epsilon$ . Consider the  $r'$  low-rank approximation to  $M$ ,

$$L_{\text{SVD}} := \sum_{l=1}^{r'} s_l u_l v_l^T. \quad (11.57)$$

This is the optimal rank  $r'$  approximation, so

$$\|L - M\|_F^2 \geq \|L_{\text{SVD}} - M\|_F^2 \quad (11.58)$$

$$= \sum_{i=r'+1}^{n_1} s_i^2. \quad (11.59)$$

If  $\|L' - M\|_F \leq \epsilon$ , then  $\sum_{i=r'+1}^{n_1} s_i^2 \leq \epsilon^2$ . However, this is impossible because  $r$  is defined to be the smallest positive integer for which  $\sum_{i=r+1}^{n_1} s_i^2 \leq \epsilon^2$ .

#### 11.10 (Double centering)

- a) The average of the entries of the column sample mean is the grand sample mean,

$$\frac{1}{n_2} \sum_{j=1}^{n_2} \mu_{\text{row}}[j] = \frac{1}{n_2} \sum_{j=1}^{n_2} \frac{1}{n_1} \sum_{i=1}^{n_1} D[i, j] \quad (11.60)$$

$$= \mu_{\text{all}}. \quad (11.61)$$

The average of the entries of the row sample mean is also the grand sample mean,

$$\frac{1}{n_1} \sum_{i=1}^{n_1} \mu_{\text{col}}[i] = \frac{1}{n_1} \sum_{i=1}^{n_1} \frac{1}{n_2} \sum_{j=1}^{n_2} D[i, j] \quad (11.62)$$

$$= \mu_{\text{all}}. \quad (11.63)$$

- b) We define a matrix  $S$  obtained by subtracting the row sample mean and the row column mean,

$$S[i, j] := D[i, j] - \mu_{\text{col}}[i] - \mu_{\text{row}}[j], \quad 1 \leq i \leq n_1, 1 \leq j \leq n_2. \quad (11.64)$$

The  $j$ th entry of the row sample mean of  $S$  equals,

$$\frac{1}{n_1} \sum_{i=1}^{n_1} S[i, j] = \frac{1}{n_1} \sum_{i=1}^{n_1} (D[i, j] - \mu_{\text{col}}[i] - \mu_{\text{row}}[j]) \quad (11.65)$$

$$= \mu_{\text{row}}[j] - \frac{1}{n_1} \sum_{i=1}^{n_1} \mu_{\text{col}}[i] - \mu_{\text{row}}[j] \quad (11.66)$$

$$= -\mu_{\text{all}}, \quad (11.67)$$

where the last equality follows from (11.63). Consequently, the row sample mean of the matrix is constant and equal to the additive inverse of the grand sample mean of the

original matrix. This is also the case for the column sample mean. The  $i$ th entry of the column sample mean equals

$$\frac{1}{n_2} \sum_{j=1}^{n_2} S[i, j] = \frac{1}{n_2} \sum_{j=1}^{n_2} (D[i, j] - \mu_{\text{col}}[i] - \mu_{\text{row}}[j]) \quad (11.68)$$

$$= \mu_{\text{col}}[i] - \mu_{\text{col}}[i] - \frac{1}{n_2} \sum_{j=1}^{n_2} \mu_{\text{row}}[j] \quad (11.69)$$

$$= -\mu_{\text{all}}, \quad (11.70)$$

where the last equality follows from (11.61).

- c) Above we see that subtracting the row and column sample means results in a matrix where the row and column sample means equal the additive inverse of the grand sample mean. This suggests adding it to each entry to cancel out both means. The entries of the resulting doubly centered matrix are defined by

$$C[i, j] := D[i, j] - \mu_{\text{col}}[i] - \mu_{\text{row}}[j] + \mu_{\text{all}}, \quad 1 \leq i \leq n_1, 1 \leq j \leq n_2. \quad (11.71)$$

We check that the matrix indeed has zero column and row sample means. The  $j$ th entry of the row sample mean of  $C$  equals,

$$\frac{1}{n_1} \sum_{i=1}^{n_1} C[i, j] = \frac{1}{n_1} \sum_{i=1}^{n_1} (D[i, j] - \mu_{\text{col}}[i] - \mu_{\text{row}}[j] + \mu_{\text{all}}) \quad (11.72)$$

$$= \mu_{\text{row}}[j] - \frac{1}{n_1} \sum_{i=1}^{n_1} \mu_{\text{col}}[i] - \mu_{\text{row}}[j] + \mu_{\text{all}} \quad (11.73)$$

$$= 0, \quad (11.74)$$

where the last equality follows from (11.63). The argument for the column sample mean is essentially the same.

#### 11.11 (Mutant mosquito)

- a) Based on the SVD analysis, which yields the rank-one model

$$L[i, j] := m(D) + s_i u_1[i] v_1[j], \quad (11.75)$$

we can set  $c := m(D) = 3$  and  $b[\text{user}]$  equal to the corresponding entry of  $v_1$ , so

$$b = \begin{pmatrix} \text{Bob} & \text{Molly} & \text{Mary} & \text{Larry} \\ 0.48 & 0.52 & -0.48 & -0.52 \end{pmatrix}. \quad (11.76)$$

- b) The residual sum of squared errors is

$$\text{error} = \sum_{\text{user} \in \{\text{Bob}, \text{Mary}, \text{Larry}\}} (\text{estimate} - \text{observed})^2 \quad (11.77)$$

$$\begin{aligned} &= (3 + 0.48a[\text{MM}] - 2)^2 + (3 - 0.48a[\text{MM}] - 4)^2 + (3 - 0.052a[\text{MM}] - 4)^2 \\ &= 3 + 2.024a[\text{MM}] + 0.464a[\text{MM}]^2. \end{aligned} \quad (11.78)$$

This is a convex quadratic (as the quadratic coefficient is positive). Setting its derivative

to zero yields the coefficient that minimizes the sum of squared errors:

$$a_{\text{ls}}[\text{MM}] = -\frac{2.024}{2 \cdot 0.464} \quad (11.79)$$

$$= -2.18. \quad (11.80)$$

c) The estimate for Molly is

$$\text{Molly} := c - a_{\text{ls}}[\text{MM}]b[\text{Molly}] \quad (11.81)$$

$$= 3 - 2.18 \cdot 0.52 = 1.87. \quad (11.82)$$

11.12 (Exact matrix completion) The matrix cannot be rank 1, because the first and second columns are clearly linearly independent. We will try to find a rank-2 matrix that is consistent with the revealed entry. If the matrix is rank 2, then the third column is a linear combination of the first and second columns with coefficients that satisfy

$$2\alpha + 2\beta = 3, \quad (11.83)$$

$$2\alpha + 3\beta = 4. \quad (11.84)$$

so  $\alpha = 0.5$  and  $\beta = 1$ . As a result,  $D[4, 3] = 5$ .

Similarly, the fourth column is a linear combination of the first and second columns with coefficients that satisfy

$$2\gamma + 3\theta = -2, \quad (11.85)$$

$$2\gamma + 4\theta = -3. \quad (11.86)$$

so  $\gamma = 0.5$  and  $\theta = -1$ . As a result,  $D[2, 4] = -1$ .

Finally, the first row must be a linear combination of the second and third rows with coefficients that satisfy

$$3\delta + 4\psi = 2, \quad (11.87)$$

$$-\delta - 2\psi = 0. \quad (11.88)$$

so  $\delta = 2$  and  $\psi = -1$ . As a result,  $D[1, 1] = 2$  and  $D[1, 2] = 1$ .

We have succeeded in completing the entries with a rank-2 matrix, which equals

$$D = \begin{bmatrix} 2 & 1 & 2 & 0 \\ 2 & 2 & 3 & -1 \\ 2 & 3 & 4 & -2 \\ 2 & 4 & 5 & -3 \end{bmatrix}. \quad (11.89)$$

11.13 (Topic modeling)

a) The singular values equal

$$D = USV^T = U \begin{bmatrix} 23.64 & 0 & 0 & 0 & 0 & 0 \\ 0 & 18.82 & 0 & 0 & 0 & 0 \\ 0 & 0 & 14.23 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3.63 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2.03 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1.36 \end{bmatrix} V^T \quad (11.90)$$

Three of the singular values are much larger than the rest, which suggests fitting a

model with three topics. The singular vectors do not have an intuitive interpretation. In particular, they do not allow to cluster the words

$$\begin{array}{rcccccc} & \text{a} & \text{b} & \text{c} & \text{d} & \text{e} & \text{f} \\ u_1 & = & (-0.24 & -0.47 & -0.24 & -0.32 & -0.58 & -0.47) \\ u_2 & = & (0.64 & -0.23 & 0.67 & -0.03 & -0.18 & -0.21) \\ u_3 & = & (-0.08 & -0.39 & -0.08 & 0.77 & 0.28 & -0.40) \end{array} \quad (11.91)$$

or the articles

$$\begin{array}{rcccccccccc} & \text{singer} & \text{GDP} & \text{senate} & \text{election} & \text{vote} & \text{stock} & \text{concert} & \text{market} & \text{band} \\ v_1 & = & (-0.18 & -0.24 & -0.51 & -0.38 & -0.46 & -0.34 & -0.2 & -0.3 & -0.22) \\ v_2 & = & (0.47 & 0.01 & -0.22 & -0.15 & -0.25 & -0.07 & 0.63 & -0.05 & 0.49) \\ v_3 & = & (-0.13 & 0.47 & -0.3 & -0.14 & -0.37 & 0.52 & -0.04 & 0.49 & -0.07) \end{array}$$

- b) The left and right singular vectors must be orthogonal, so the only way they can all have nonnegative entries is if their support is completely disjoint.  
 c) We apply NMF to obtain two matrices  $W \in \mathbb{R}^{m \times k}$  and  $H \in \mathbb{R}^{k \times n}$  such that

$$M[i, j] \approx \sum_{l=1}^3 w_l[i] h_l[j], \quad w_l[i] \geq 0, \quad 1 \leq l \leq r, \quad 1 \leq i \leq 6, \quad (11.92)$$

$$h_l[j] \geq 0, \quad 1 \leq l \leq r, \quad 1 \leq j \leq 9. \quad (11.93)$$

The vectors  $h_1$ ,  $h_2$  and  $h_3$  cluster the words into three topics:

$$\begin{array}{rcccccccccc} & \text{singer} & \text{GDP} & \text{senate} & \text{election} & \text{vote} & \text{stock} & \text{concert} & \text{market} & \text{band} \\ h_1 & = & (0.34 & 0 & 3.73 & 2.54 & 3.67 & 0.52 & 0 & 0.35 & 0.35) \\ h_2 & = & (0 & 2.21 & 0.21 & 0.45 & 0 & 2.64 & 0.21 & 2.43 & 0.22) \\ h_3 & = & (3.22 & 0.37 & 0.19 & 0.2 & 0 & 0.12 & 4.13 & 0.13 & 3.43) \end{array}$$

The first topic corresponds to the entries that are not zero (or very small) in  $h_1$ : senate, election and vote. The second corresponds to  $h_2$ : GDP, stock and market. The third corresponds to  $h_3$ : singer, concert and band.

The entries of  $w_1$ ,  $w_2$  and  $w_3$  allow us to assign the topics to articles:

$$\begin{array}{rcccccc} & \text{a} & \text{b} & \text{c} & \text{d} & \text{e} & \text{f} \\ w_1 & = & (0.03 & 2.23 & 0 & 0 & 1.59 & 2.24) \\ w_2 & = & (0.1 & 0 & 0.08 & 3.13 & 2.32 & 0) \\ w_3 & = & (2.13 & 0 & 2.22 & 0 & 0 & 0.03) \end{array} \quad (11.94)$$

Articles  $b$ ,  $e$  and  $f$  are about politics (topic 1),  $d$  and  $e$  about economics (topic 3) and  $a$  and  $c$  about music (topic 3).

- d) The Frobenius-norm error of the SVD-based approximation is smaller (4.38 compared to 4.44). This makes sense because it is optimal, and nonnegative matrix factorization has an additional constraint that can only make the fit worse. Nevertheless, the NMF model is clearly better for the purpose of topic modeling, because it yields a reasonable choice of models that fit the data well and allow us to cluster the words and articles.

# 12

## Regression And Classification

### Exercises

12.1 (Heart beat)

a) The simple-regression linear MMSE estimator of  $\tilde{b}$  given  $\tilde{x}[1]$  is

$$\ell_{\text{MMSE}}(\tilde{x}[1]) = \frac{\text{Cov}[\tilde{x}[1], \tilde{b}]\tilde{x}[1]}{\text{Var}[\tilde{x}[1]]}. \quad (12.1)$$

By the independence assumptions,

$$\text{Var}[\tilde{x}[1]] = \text{Var}[\tilde{b} + \tilde{m} + \tilde{z}_1] \quad (12.2)$$

$$= \text{Var}[\tilde{b}] + \text{Var}[\tilde{m}] + \text{Var}[\tilde{z}_1] \quad (12.3)$$

$$= 12, \quad (12.4)$$

$$\text{Cov}[\tilde{b}, \tilde{x}[1]] = \mathbb{E}[\tilde{b}(\tilde{b} + \tilde{m} + \tilde{z}_1)] \quad (12.5)$$

$$= \mathbb{E}[\tilde{b}^2] + \mathbb{E}[\tilde{b}\tilde{m}] + \mathbb{E}[\tilde{b}\tilde{z}_1] \quad (12.6)$$

$$= \mathbb{E}[\tilde{b}^2] + \mathbb{E}[\tilde{b}]\mathbb{E}[\tilde{m}] + \mathbb{E}[\tilde{b}]\mathbb{E}[\tilde{z}_1] \quad (12.7)$$

$$= 1. \quad (12.8)$$

Therefore,

$$\ell_{\text{MMSE}}(\tilde{x}[1]) = \frac{\tilde{x}[1]}{12}. \quad (12.9)$$

The estimator just shrinks the measurement.

The coefficient of determination equals

$$R^2 := \frac{\text{Var}[\ell_{\text{MMSE}}(\tilde{x}[1])]}{\text{Var}[\tilde{b}]} \quad (12.10)$$

$$= \frac{12}{12^2} = \frac{1}{12}, \quad (12.11)$$

which indicates that the estimate is not great.

b) The linear MMSE estimator equals

$$\ell_{\text{MMSE}}(\tilde{x}) = \Sigma_{\tilde{x}b}^T \Sigma_{\tilde{x}}^{-1} \tilde{x}. \quad (12.12)$$

By the independence assumptions,  $\text{Var} [\tilde{x}[1]] = 12$ ,

$$\text{Var} [\tilde{x}[2]] = \text{Var} [\tilde{m}] + \text{Var} [\tilde{z}_2] \quad (12.13)$$

$$= 11, \quad (12.14)$$

$$\text{Cov} [\tilde{x}[1]\tilde{x}[2]] = \mathbb{E}[(\tilde{b} + \tilde{m} + \tilde{z}_1)(\tilde{m} + \tilde{z}_2)] \quad (12.15)$$

$$= \mathbb{E}[\tilde{m}^2] \quad (12.16)$$

$$= 10, \quad (12.17)$$

so

$$\Sigma_{\tilde{x}} = \begin{bmatrix} 12 & 10 \\ 10 & 11 \end{bmatrix} \quad (12.18)$$

Again by the independence assumptions,  $\text{Cov} [\tilde{b}, \tilde{x}[1]] = 1$ , and

$$\text{Cov} [\tilde{b}\tilde{x}[2]] = \mathbb{E}[\tilde{b}(\tilde{m} + \tilde{z}_2)] \quad (12.19)$$

$$= 0, \quad (12.20)$$

so

$$\Sigma_{\tilde{x}\tilde{b}} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (12.21)$$

Therefore

$$\ell_{\text{MMSE}}(\tilde{x}) = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 12 & 10 \\ 10 & 11 \end{bmatrix}^{-1} \tilde{x} \quad (12.22)$$

$$= \frac{11\tilde{x}[1] - 10\tilde{x}[2]}{32}. \quad (12.23)$$

The estimate approximately cancels out the signal from the mother.

The linear MMSE estimator can be expressed as

$$\ell_{\text{MMSE}}(\tilde{x}) = \frac{11(\tilde{b} + \tilde{m} + \tilde{z}_1) - 10(\tilde{m} + \tilde{z}_2)}{32} \quad (12.24)$$

$$= \frac{11\tilde{b} + \tilde{m} + 11\tilde{z}_1 - 10\tilde{z}_2}{32}. \quad (12.25)$$

By the independence assumption,

$$\text{Var} [\ell_{\text{MMSE}}(\tilde{x})] = \frac{11^2 \text{Var} [\tilde{b}] + \text{Var} [\tilde{m}] + 11^2 \text{Var} [\tilde{z}_1] + 10^2 \text{Var} [\tilde{z}_2]}{32^2} \quad (12.26)$$

$$= \frac{11^2 + 10 + 11^2 + 10^2}{32^2} = 0.344. \quad (12.27)$$

The coefficient of determination equals

$$R^2 := \frac{\text{Var} [\ell_{\text{MMSE}}(\tilde{x})]}{\text{Var} [\tilde{b}]} \quad (12.28)$$

$$= 0.344, \quad (12.29)$$

which is much higher than  $1/12 = 0.083$ .

12.2 (PCA and OLS) For PCA, we have

$$u_1 = \arg \max_{\|v\|_2=1} v^T \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) v \quad (12.30)$$

$$= \arg \max_{\|v\|_2=1} \sum_{i=1}^n (v^T x_i)^2. \quad (12.31)$$

Let us define the  $\ell_2$ -norm distance between  $x_i$  and the line collinear with a unit  $\ell_2$ -norm vector  $v$ , as the distance between  $x_i$  and its nearest point on the line:

$$d_v(x_i) := \arg \min_{\gamma \in \mathbb{R}} \|x_i - \gamma v\|_2^2. \quad (12.32)$$

The nearest point to  $x_i$  on the line is the projection  $(x_i^T v)v$  (you can check this by computing the derivative with respect to  $\gamma$ ), which results in a residual that is orthogonal to the projection. Consequently, by Pythagoras' theorem,

$$d_v(x_i) = \|x_i - (x_i^T v)v\|_2^2 \quad (12.33)$$

$$= \|x_i\|_2^2 - \|(x_i^T v)v\|_2^2 \quad (12.34)$$

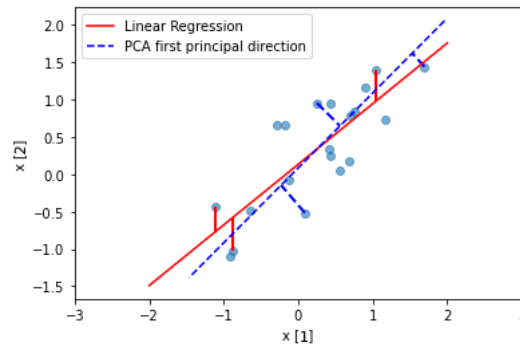
$$= \|x_i\|_2^2 - (x_i^T v)^2. \quad (12.35)$$

We conclude that, by Eq. (12.31), the line collinear with  $u_1$  minimizes the sum of squared  $\ell_2$ -norm distances to the data.

The line corresponding to the OLS estimator has slope  $\beta_{\text{OLS}}$ , satisfying

$$\beta_{\text{OLS}} := \arg \min_{\beta} \sum_{i=1}^n (x_i[2] - \beta x_i[1])^2. \quad (12.36)$$

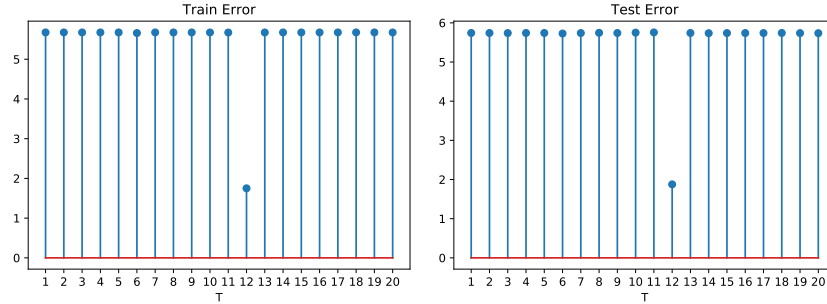
The OLS cost function can be interpreted as the sum of squared differences between the second component of the data point  $x_i$  and the second component of the point  $(x_i[1], \beta x_i[1])$ , which is on the line with slope  $\beta$ . On a 2D plane, if the second component is on the vertical axis, this exactly equal to the sum of squared vertical distances between the data and the line. The line obtained from linear regression minimizes this sum. In contrast PCA minimizes the sum of square Euclidean distances as explained above, so they are different.



12.3 (Global warming)

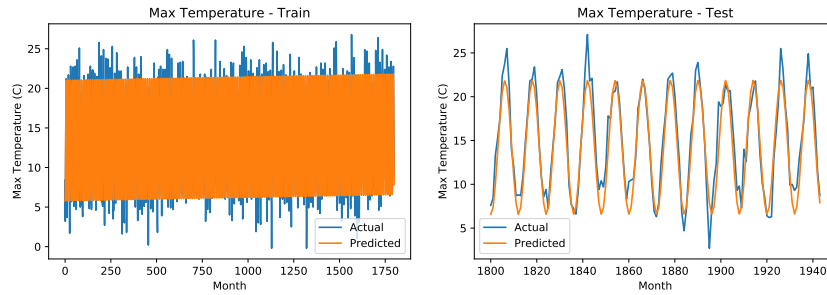


- a) We have four parameters in the model and 1800 data points to fit the model. The model is very unlikely to overfit.
- b) Below we show a plot of the test and train error as a function of the period  $T$ :



The best value is  $T^* = 12$ , which is not surprising, since the weather is cyclical with a 12 month period.

- c) Here is the fit of the model to the training and test data



- d) The OLS coefficients are:

Term	OLS coefficients
1	$1.34 \cdot 10^1$
$t$	$4.38 \cdot 10^{-4}$
$\cos(2\pi t/12)$	$-7.61 \cdot 10^0$
$\sin(2\pi t/12)$	$-4.83 \cdot 10^{-1}$

The constant term provides an offset for all temperatures. The linear term linear trends in temperature. The cosine and sine terms account for the yearly periodic pattern.

- e) The coefficient of the linear term indicates the increase or decrease rate in temperatures. According to the model, the temperature in Oxford is increasing at a rate of  $4.38 \cdot 10^{-4}$  per month, which translates to 0.526 degrees per century.

- 12.4 (OLS estimator) We denote the residual sum of squares, as a function of the vector of coefficients  $\beta$  and the additive constant  $\alpha$ , by

$$\text{RSS}(\beta, \alpha) := \sum_{i=1}^n \left( y_i - \beta^T x_i - \alpha \right)^2. \quad (12.37)$$

We denote by  $\alpha^*(\beta)$  the optimal value of  $\alpha$  for a fixed value of  $\beta$ . Equivalently,  $\alpha^*(\beta)$  is the best constant estimate of  $y_i - \beta x_i$ ,  $1 \leq i \leq n$ , in terms of the residual sum of squares.

By Theorem 7.31,

$$\alpha_*(\beta) := \arg \min_{\alpha} \sum_{i=1}^n \left( y_i - \beta^T x_i - \alpha \right)^2 \quad (12.38)$$

$$= \frac{1}{n} \sum_{i=1}^n \left( y_i - \beta^T x_i \right) \quad (12.39)$$

$$= m(Y) - \beta^T m(X). \quad (12.40)$$

Consequently, for any  $\beta$  and any  $\alpha$   $\text{RSS}(\beta, \alpha) \geq \text{RSS}(\beta, \alpha^*(\beta))$ . To obtain  $\beta_{\text{OLS}}$  we fix  $\alpha$  to  $\alpha^*(\beta)$  and minimize the sum of squared errors,

$$\beta_{\text{OLS}} = \arg \min_{\beta} \text{RSS}(\beta, \alpha^*(\beta)). \quad (12.41)$$

Let  $\text{ct}(x_i) := x_i - m(X)$  and  $y_i - m(Y)$  be the centered features and response. We have

$$\text{RSS}(\beta, \alpha^*(\beta)) = \sum_{i=1}^n \left( y_i - \beta^T x_i - m(Y) + \beta^T m(X) \right)^2 \quad (12.42)$$

$$= \sum_{i=1}^n \text{ct}(y_i)^2 + \beta^T \left( \sum_{i=1}^n \text{ct}(x_i) \right) \beta - 2\beta^T \sum_{i=1}^n \text{ct}(x_i) \text{ct}(y_i) \quad (12.43)$$

$$= (n-1) \left( v(Y) + \beta^T \Sigma_X \beta - 2\beta^T \Sigma_{XY} \right). \quad (12.44)$$

As a function of  $\beta$ , the residual sum of squares is a quadratic form. Its gradient and Hessian with respect to  $\beta$  equal

$$\nabla_{\beta} \text{RSS}(\beta, \alpha^*(\beta)) = 2(n-1)\Sigma_X \beta - 2\Sigma_{XY}, \quad (12.45)$$

$$\nabla_{\beta}^2 \text{RSS}(\beta, \alpha^*(\beta)) = 2(n-1)\Sigma_X. \quad (12.46)$$

Covariance matrices are positive semidefinite, because by the properties of sample covariance matrices, for any vector  $a \in \mathbb{R}^d$

$$a^T \Sigma_X a = a^T \left( \frac{1}{n-1} \sum_{i=1}^n \text{ct}(x_i) \text{ct}(x_i)^T \right) a \quad (12.47)$$

$$= \frac{1}{n-1} \sum_{i=1}^n \left( a^T \text{ct}(x_i) \right)^2 \geq 0. \quad (12.48)$$

Since  $\Sigma_X$  is invertible, there cannot be a nonzero vector such that  $\Sigma_X a$  equals the zero vector, so the inequality is strict as long as  $a$  is not the zero vector. This means that the quadratic function is strictly convex and we can set its gradient to zero to find the value of  $\beta$  that achieves the unique minimum:

$$\beta_{\text{OLS}} = \Sigma_X^{-1} \Sigma_{XY}. \quad (12.49)$$

The corresponding value of  $\alpha^*(\beta)$  is

$$\alpha_{\text{OLS}} = \alpha^*(\beta_{\text{OLS}}) \quad (12.50)$$

$$= m(Y) - \Sigma_{XY}^T \Sigma_X^{-1} m(X). \quad (12.51)$$

- 12.5 (Ice cream sales) To adjust for the effect of temperature, we fit a long regression model where the response is the sales and the features are advertising and temperature. Let  $X$  denote the set of features (i.e. advertising revenue and temperature pairs), and  $Y$  the

response. We can extract the sample covariance matrix of  $X$  and the cross-covariance between  $X$  and  $Y$  directly from the sample covariance matrix of the data. The OLS coefficients equal

$$\beta_{\text{OLS}} = \Sigma_X^{-1} \Sigma_{X,Y} \quad (12.52)$$

$$= \begin{bmatrix} 900 & 540 \\ 540 & 400 \end{bmatrix}^{-1} \begin{bmatrix} 960 \\ 720 \end{bmatrix} \quad (12.53)$$

$$= \begin{bmatrix} -0.07 \\ 1.89 \end{bmatrix}. \quad (12.54)$$

Consequently, according to the model, for a fixed temperature the sales are on average *inversely proportional* to advertising! The correlation between advertising and sales is due to the confounding effect of temperature. We definitely don't agree with the marketing department; there is no evidence that advertising is beneficial.

12.6 (Linear regression and maximum likelihood estimation) To alleviate notation, we define the design matrix

$$X_{\text{train}} := \begin{bmatrix} x_1^T \\ x_2^T \\ \dots \\ x_n^T \end{bmatrix}, \quad (12.55)$$

the observed response vector

$$y_{\text{train}} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} \quad (12.56)$$

and the random vector

$$\tilde{y} = \begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \\ \dots \\ \tilde{y}_n \end{bmatrix}. \quad (12.57)$$

Under the modeling assumptions, the entries of  $\tilde{y}$  are independent Gaussians with mean  $x_i^T \beta$  and variance  $\sigma^2$ , so its joint pdf evaluated at the observed data equals

$$f_{\tilde{y}}(y_{\text{train}}) = \prod_{i=1}^n \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{1}{2\sigma^2} (y_i - x_i^T \beta)^2\right) \quad (12.58)$$

$$= \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{1}{2\sigma^2} \|y_{\text{train}} - X_{\text{train}}\beta\|_2^2\right), \quad (12.59)$$

which is equal to the likelihood, if we interpret it as a function of  $\beta$ . Consequently the log likelihood equals

$$\log \mathcal{L}_y(\beta) = -\frac{1}{2\sigma^2} \|y_{\text{train}} - X_{\text{train}}\beta\|_2^2 - \frac{n}{2} \log(2\pi\sigma^2), \quad (12.60)$$

so the maximum-likelihood estimator is

$$\beta_{\text{ML}} = \arg \max_{\beta} \log \mathcal{L}_y(\beta) \quad (12.61)$$

$$= \arg \max_{\beta} -\frac{1}{2\sigma^2} \|y_{\text{train}} - X_{\text{train}}\beta\|_2^2 - \frac{n}{2} \log(2\pi\sigma^2) \quad (12.62)$$

$$= \arg \min_{\beta} \|y_{\text{train}} - X_{\text{train}}\beta\|_2^2, \quad (12.63)$$

which is *the same as the OLS estimator*.

12.7 (Best linear unbiased estimator) To alleviate notation, we set  $\tilde{y} := \tilde{y}_{\text{train}}$ ,  $X := X_{\text{train}}$ ,  $\tilde{z} := \tilde{z}_{\text{train}}$  and  $\beta := \beta_{\text{true}}$ .

a) By linearity of expectation

$$\mathbb{E}[C\tilde{y}] = \mathbb{E}[C(X\beta + \tilde{z})] = CX\beta + C\mathbb{E}[\tilde{z}] = CX\beta. \quad (12.64)$$

b) Let  $\mu := \mathbb{E}[C\tilde{y}] = CX\beta$ . The covariance matrix equals

$$\Sigma_{C\tilde{y}} = \mathbb{E}[(C\tilde{y})(C\tilde{y})^T] - \mu\mu^T = C\mathbb{E}[\tilde{y}\tilde{y}^T]C^T - \mu\mu^T \quad (12.65)$$

$$= C\mathbb{E}[X\beta\beta^T X^T + \tilde{z}\tilde{z}^T]C^T - \mu\mu^T \quad (12.66)$$

$$= CX\beta\beta^T X^T C^T + C(\sigma^2 I)C^T - \mu\mu^T \quad (12.67)$$

$$= \mu\mu^T + \sigma^2 CC^T - \mu\mu^T \quad (12.68)$$

$$= \sigma^2 CC^T. \quad (12.69)$$

c) Since  $\mathbb{E}[C\tilde{y}] = CX\beta$ ,

$$\mathbb{E}[C\tilde{y}] = CX\beta = \left((X^T X)^{-1} X^T + D\right) X\beta = (I + DX)\beta = \beta + DX\beta. \quad (12.70)$$

Thus we require  $DX = 0$ .

d) We have

$$\text{Var} \left[ v^T C\tilde{y} \right] = v^T \Sigma_{C\tilde{y}} v \geq v^T \Sigma_{\tilde{\beta}_{\text{OLS}}} v = \text{Var} \left[ v^T \tilde{\beta}_{\text{OLS}} \right] \quad (12.71)$$

for any  $v \in \mathbb{R}^d$ , if and only if the matrix  $\Sigma_{C\tilde{y}} - \Sigma_{\tilde{\beta}_{\text{OLS}}}$  is positive semidefinite. By (12.69), since  $D := C - (X^T X)^{-1} X^T$  and  $DX = 0$  (which must hold if  $C\tilde{y}$  is unbiased),

$$\Sigma_{C\tilde{y}} = \sigma^2 CC^T \quad (12.72)$$

$$= \sigma^2 \left( (X^T X)^{-1} X^T + D \right) \left( (X^T X)^{-1} X^T + D \right)^T \quad (12.73)$$

$$= \sigma^2 \left( (X^T X)^{-1} X^T X (X^T X)^{-1} + DD^T + (X^T X)^{-1} X^T D^T + DX (X^T X)^{-1} \right)$$

$$= \sigma^2 ((X^T X)^{-1} + DD^T) \quad (12.74)$$

$$= \Sigma_{\tilde{\beta}_{\text{OLS}}} + \sigma^2 DD^T, \quad (12.75)$$

where we have used the fact that by (12.69)

$$\Sigma_{\tilde{\beta}_{\text{OLS}}} = \Sigma_{(X^T X)^{-1} X^T \tilde{y}} \quad (12.76)$$

$$= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \quad (12.77)$$

$$= \sigma^2 (X^T X)^{-1}. \quad (12.78)$$

$DD^T$  is positive semidefinite since for any  $v$   $v^T DD^T v = \|D^T v\|_2^2 \geq 0$ , so  $\Sigma_{C\tilde{y}} - \Sigma_{\tilde{\beta}_{\text{OLS}}} = \sigma^2 DD^T$  is indeed positive semidefinite and the proof is complete.

- 12.8 (Augmented dataset) The ridge-regression cost function can be reformulated as the following least-squares problem:

$$\beta_{\text{RR}} := \arg \min_{\beta} \left\| \begin{bmatrix} y \\ 0 \end{bmatrix} - \begin{bmatrix} X^T \\ \sqrt{\lambda} I \end{bmatrix} \beta \right\|_2^2. \quad (12.79)$$

There are  $d$  additional examples. For the  $i$ th example the feature vector equals a one-hot vector where the  $i$ th entry equals  $\lambda$  and the rest equal zero. All the response values for the additional examples equal zero. The additional examples force the linear coefficients to fit zeros, preventing them from becoming too large.

- 12.9 (Correlated features)

a) Let us define  $\bar{\alpha} := \sqrt{1 - \alpha^2}$  to alleviate notation. We have

$$\beta_{\text{OLS}} = (X_{\text{train}}^T X_{\text{train}})^{-1} X_{\text{train}}^T y \quad (12.80)$$

$$= \begin{bmatrix} 1 & \alpha \\ \alpha & 1 \end{bmatrix}^{-1} \begin{bmatrix} x_1^T \\ \alpha x_1^T + \bar{\alpha} x_{\perp}^T \end{bmatrix} (\beta_{\text{true}} x_1 + z) \quad (12.81)$$

$$= \frac{1}{1 - \alpha^2} \begin{bmatrix} 1 & -\alpha \\ -\alpha & 1 \end{bmatrix} \begin{bmatrix} \beta_{\text{true}} + 0.1 \\ \alpha(\beta_{\text{true}} + 0.1) + 0.1\bar{\alpha} \end{bmatrix} \quad (12.82)$$

$$= \frac{1}{1 - \alpha^2} \begin{bmatrix} (1 - \alpha^2)(\beta_{\text{true}} + 0.1) - 0.1\alpha\bar{\alpha} \\ 0.1\bar{\alpha} \end{bmatrix} \quad (12.83)$$

$$= \begin{bmatrix} \beta_{\text{true}} + 0.1 - \frac{0.1\alpha}{\bar{\alpha}} \\ \frac{0.1}{\bar{\alpha}} \end{bmatrix}. \quad (12.84)$$

The estimated coefficients tend to

$$\lim_{\alpha \rightarrow 1} \beta_{\text{OLS}} = \begin{bmatrix} -\infty \\ \infty \end{bmatrix}. \quad (12.85)$$

The OLS estimator overfits the response vector, which causes the entries of the coefficient estimate to explode.

- b) The estimate of the response equals

$$y_{\text{OLS}} := X_{\text{train}} \beta_{\text{OLS}} \quad (12.86)$$

$$= \begin{bmatrix} x_1 & \alpha x_1 + \bar{\alpha} x_{\perp} \end{bmatrix} \begin{bmatrix} \beta_{\text{true}} + 0.1 - \frac{0.1\alpha}{\bar{\alpha}} \\ \frac{0.1}{\bar{\alpha}} \end{bmatrix} \quad (12.87)$$

$$= \left( \beta_{\text{true}} + 0.1 - \frac{0.1\alpha}{\bar{\alpha}} \right) x_1 + \frac{0.1}{\bar{\alpha}} (\alpha x_1 + \bar{\alpha} x_{\perp}) \quad (12.88)$$

$$= (\beta_{\text{true}} + 0.1)x_1 + 0.1x_{\perp}, \quad (12.89)$$

which does not change as  $\alpha \rightarrow 1$ . The estimate of the response is not collinear with  $x_1$  even when  $\alpha \rightarrow 1$ , although the second feature tends to be collinear with the first one in this limit! This is because the noise in the training response vector has a component in that direction, and the linear regression estimate fits it, even if the corresponding component in the second feature vector is tiny.

c)

$$\beta_{\text{RR}} = (X_{\text{train}}^T X_{\text{train}} + \lambda I)^{-1} X_{\text{train}}^T y \quad (12.90)$$

$$= \begin{bmatrix} 1 + \lambda & \alpha \\ \alpha & 1 + \lambda \end{bmatrix}^{-1} \begin{bmatrix} x_1^T \\ \alpha x_1^T + \bar{\alpha} x_{\perp}^T \end{bmatrix} (\beta_{\text{true}} x_1 + z) \quad (12.91)$$

$$= \frac{1}{(1 + \lambda)^2 - \alpha^2} \begin{bmatrix} 1 + \lambda & -\alpha \\ -\alpha & 1 + \lambda \end{bmatrix} \begin{bmatrix} \beta_{\text{true}} + 0.1 \\ \alpha(\beta_{\text{true}} + 0.1) + 0.1\bar{\alpha} \end{bmatrix} \quad (12.92)$$

$$= \frac{1}{(1 + \lambda)^2 - \alpha^2} \begin{bmatrix} (1 + \lambda - \alpha^2)(\beta_{\text{true}} + 0.1) - 0.1\alpha\bar{\alpha} \\ \lambda\alpha(\beta_{\text{true}} + 0.1) + 0.1\bar{\alpha}(1 + \lambda) \end{bmatrix}. \quad (12.93)$$

The estimated coefficients tend to

$$\lim_{\alpha \rightarrow 1} \beta_{\text{RR}} = \frac{1}{2\lambda + \lambda^2} \begin{bmatrix} \lambda(\beta_{\text{true}} + 0.1) \\ \lambda(\beta_{\text{true}} + 0.1) \end{bmatrix} \quad (12.94)$$

$$= \frac{1}{2 + \lambda} \begin{bmatrix} \beta_{\text{true}} + 0.1 \\ \beta_{\text{true}} + 0.1 \end{bmatrix}. \quad (12.95)$$

The coefficients no longer explode as in OLS. The two coefficients are the same, which makes sense because when  $\alpha \rightarrow 1$  both features are equal to  $x_1$ .

d)

$$\lim_{\alpha \rightarrow 1} y_{\text{RR}} := (\lim_{\alpha \rightarrow 1} X_{\text{train}}) \lim_{\alpha \rightarrow 1} \beta_{\text{RR}} \quad (12.96)$$

$$= \begin{bmatrix} x_1 & x_1 \end{bmatrix} \frac{1}{2 + \lambda} \begin{bmatrix} \beta_{\text{true}} + 0.1 \\ \beta_{\text{true}} + 0.1 \end{bmatrix} \quad (12.97)$$

$$= \frac{2(\beta_{\text{true}} + 0.1)x_1}{2 + \lambda}. \quad (12.98)$$

In contrast to the OLS response estimate, the ridge-regression response estimate is collinear with the true feature vector in the limit where  $\alpha \rightarrow 1$ .

#### 12.10 (Prior knowledge)

- a) A natural way to modify the ridge-regression cost function is to incorporate  $\beta_{\text{prior}}$  in the regularization term, to promote solutions that are close to it:

$$\beta_{\text{RRP}} := \arg \min_{\beta} \|y_{\text{train}} - X_{\text{train}}\beta\|_2^2 + \lambda \|\beta_{\text{prior}} - \beta\|_2^2. \quad (12.99)$$

- b) The cost function can be reformulated to equal a modified least-squares problem

$$\beta_{\text{RRP}} := \arg \min_{\beta} \left\| \begin{bmatrix} \tilde{y}_{\text{train}} \\ \sqrt{\lambda}\beta_{\text{prior}} \end{bmatrix} - \begin{bmatrix} X_{\text{train}} \\ \sqrt{\lambda}I \end{bmatrix} \beta \right\|_2^2. \quad (12.100)$$

By Lemma 12.8, the solution to the modified least-squares problem equals

$$\tilde{\beta}_{\text{RRP}} = \left( \begin{bmatrix} X_{\text{train}}^T & \sqrt{\lambda}I \end{bmatrix} \begin{bmatrix} X_{\text{train}} \\ \sqrt{\lambda}I \end{bmatrix} \right)^{-1} \begin{bmatrix} X_{\text{train}}^T & \sqrt{\lambda}I \end{bmatrix} \begin{bmatrix} \tilde{y}_{\text{train}} \\ \sqrt{\lambda}\beta_{\text{prior}} \end{bmatrix} \quad (12.101)$$

$$= \left( X_{\text{train}}^T X_{\text{train}} + \lambda I \right)^{-1} (X_{\text{train}}^T \tilde{y}_{\text{train}} + \lambda \beta_{\text{prior}}) \quad (12.102)$$

$$= \tilde{\beta}_{\text{RR}} + \left( \Sigma_X + \frac{\lambda}{n-1} I \right)^{-1} \lambda \beta_{\text{prior}}, \quad (12.103)$$

because by (12.254) in the main text, the ridge-regression coefficients equal

$$\tilde{\beta}_{\text{RR}} = \left( X_{\text{train}}^T X_{\text{train}} + \lambda I \right)^{-1} X_{\text{train}}^T \tilde{y}_{\text{train}}. \quad (12.104)$$

Consequently, the estimator is the same as the ridge-regression estimator except for a deterministic constant. As a result, the covariance matrix is unchanged and the mean equals  $\mathbb{E} [\tilde{\beta}_{\text{RR}}] + \lambda \left( \Sigma_X + \frac{\lambda}{n-1} I \right)^{-1} \beta_{\text{prior}}$ , which approaches  $\beta_{\text{prior}}$  as  $\lambda$  increases (since  $\mathbb{E} [\tilde{\beta}_{\text{RR}}]$  approaches zero by Theorem 12.33).

12.11 (Lasso estimation with two features)

a) If  $c_{\text{true}, \text{noise}} < -a_{\text{true}}$

$$\frac{x_{\text{true}}^T y_{\text{train}}}{n-1} = \frac{a_{\text{true}} x_{\text{true}}^T x_{\text{true}} + x_{\text{true}}^T z}{n-1} \quad (12.105)$$

$$= a_{\text{true}} + c_{\text{true}, \text{noise}} < 0. \quad (12.106)$$

b) The derivative of  $\mathcal{L}_{\text{lasso}}(a, 0)$  under our assumptions is

$$\frac{d\mathcal{L}_{\text{lasso}}(a, 0)}{da} = \frac{d(y_{\text{train}} - x_{\text{true}}a)^T (y_{\text{train}} - x_{\text{true}}a) + \lambda a}{da} \quad (12.107)$$

$$= -2x_{\text{true}}^T (y_{\text{train}} - x_{\text{true}}a) + \lambda \quad (12.108)$$

$$= -2x_{\text{true}}^T (x_{\text{true}}(a_{\text{true}} - a) + z) + \lambda \quad (12.109)$$

$$= -2x_{\text{true}}^T (x_{\text{true}}(a_{\text{true}} - a) + z) + \lambda \quad (12.110)$$

$$= 2(n-1)(a - a_{\text{true}} - c_{\text{true}, \text{noise}}) + \lambda. \quad (12.111)$$

The second derivative equals

$$\frac{d^2 \mathcal{L}_{\text{lasso}}(a, 0)}{da^2} = 2(n-1), \quad (12.112)$$

which is positive. The function is therefore convex and we can find the minimum by setting the first derivative to zero. Dividing the resulting equation by  $n-1$ , yields

$$2(a - a_{\text{true}}) - 2c_{\text{true}, \text{noise}} + \lambda_n = 0. \quad (12.113)$$

The derivative is zero at

$$a_{\text{lasso}} := a_{\text{true}} + c_{\text{true}, \text{noise}} - \frac{\lambda_n}{2}. \quad (12.114)$$

c) Since  $a_{\text{true}} \geq 0$  and  $c_{\text{true}, \text{noise}} \geq -a_{\text{true}}$ , we have  $a_{\text{lasso}} \geq 0$  as long as

$$\lambda_n \leq 2(a_{\text{true}} + c_{\text{true}, \text{noise}}) := \lambda_{\text{max}}. \quad (12.115)$$

d) The derivative equals

$$\frac{d\mathcal{L}_{\text{lasso}}(a_{\text{lasso}}, b)}{db} \quad (12.116)$$

$$= \frac{d(y_{\text{train}} - x_{\text{true}}a_{\text{lasso}} - x_{\text{other}}b)^T (y_{\text{train}} - x_{\text{true}}a_{\text{lasso}} - x_{\text{other}}b) + \lambda a_{\text{lasso}} + \lambda b}{db}$$

$$= -2x_{\text{other}}^T (y_{\text{train}} - x_{\text{true}}a_{\text{lasso}} - x_{\text{other}}b) + \lambda \quad (12.117)$$

$$= -2x_{\text{other}}^T (x_{\text{true}}(a_{\text{true}} - a_{\text{lasso}}) + z - x_{\text{other}}b) + \lambda. \quad (12.118)$$

Dividing this expression by  $n - 1$  and using the fact that  $x_{\text{other}}^T x_{\text{other}} = n - 1$  yields

$$2\rho \left( c_{\text{true},\text{noise}} - \frac{\lambda_n}{2} \right) - 2c_{\text{other},\text{noise}} + v(x_{\text{other}})b + \lambda_n \quad (12.119)$$

$$= (1 - \rho)\lambda_n - 2c_{\text{other},\text{noise}} + 2\rho c_{\text{true},\text{noise}} + b, \quad (12.120)$$

which is positive for small enough  $b$  if

$$\frac{\lambda_n}{2} > \frac{c_{\text{other},\text{noise}} - \rho c_{\text{true},\text{noise}}}{1 - \rho}. \quad (12.121)$$

e) By the same argument used to derive (12.120), the derivative  $\frac{d\mathcal{L}_{\text{lasso}}(a_{\text{lasso}}, b)}{db}$  for  $b < 0$  is negative if

$$-(1 - \rho)\lambda_n - 2c_{\text{other},\text{noise}} + 2\rho c_{\text{true},\text{noise}} + b < 0, \quad (12.122)$$

which is the case when the magnitude of  $b$  is small enough, as long as

$$\frac{\lambda_n}{2} > \frac{\rho c_{\text{true},\text{noise}} - c_{\text{other},\text{noise}}}{1 - \rho}. \quad (12.123)$$

12.12 (Ridge regression and sparsity) Under the assumptions, the sample covariance matrix of the features equals

$$\Sigma_X = \begin{bmatrix} 1 & 0.2 \\ 0.2 & 1 \end{bmatrix}. \quad (12.124)$$

Since everything has zero sample mean, the sample covariance between the response and the true feature is

$$c_{\text{true},\text{response}} = \frac{x_{\text{true}}^T y}{n - 1} \quad (12.125)$$

$$= 1 + c_{\text{true},\text{noise}} = 0.85, \quad (12.126)$$

whereas the sample covariance between the response and the spurious feature is

$$c_{\text{other},\text{response}} = \frac{x_{\text{other}}^T y}{n - 1} \quad (12.127)$$

$$= \rho + c_{\text{other},\text{noise}} = 0.355. \quad (12.128)$$

Consequently the OLS coefficient estimate equals

$$\beta_{\text{OLS}} = \Sigma_X^{-1} \Sigma_{XY} \quad (12.129)$$

$$= \begin{bmatrix} 1 & 0.2 \\ 0.2 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 0.85 \\ 0.355 \end{bmatrix} \quad (12.130)$$

$$= \begin{bmatrix} 0.81 \\ 0.19 \end{bmatrix}. \quad (12.131)$$

The eigenvectors of the sample covariance matrix (normalized to have unit  $\ell_2$  norm) equal

$$u_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad u_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \quad (12.132)$$

and the corresponding eigenvalues equal  $\xi_1 = 1.2$  and  $\xi_2 = 0.8$ . In the basis of eigenvectors, the OLS coefficients equal

$$c_{\text{OLS}} := \begin{bmatrix} u_1^T \\ u_2^T \end{bmatrix} \beta_{\text{OLS}} = \begin{bmatrix} 0.71 \\ -0.44 \end{bmatrix}, \quad (12.133)$$



so the ridge regression coefficients equal

$$\beta_{\text{RR}} = \frac{c_{\text{OLS}}[1]}{1 + \lambda_n/\xi_1} u_1 + \frac{c_{\text{OLS}}[2]}{1 + \lambda_n/\xi_2} u_2 \quad (12.134)$$

$$= \frac{0.71}{1 + \lambda_n/1.2} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \frac{0.44}{1 + \lambda_n/0.8} \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \quad (12.135)$$

where  $\lambda_n := \frac{\lambda}{n-1}$ . The first entry only cancels out for  $\lambda_n = 0.518$ , so the ridge-regression coefficient estimate is never sparse.

### 12.13 (Derivative analysis of logistic regression)

a) By basic calculus

$$\frac{d \lg f(\ell)}{d\ell} = \frac{d}{d\ell} \frac{1}{1 + \exp(-\ell)} \quad (12.136)$$

$$= \frac{\exp(-\ell)}{(1 + \exp(-\ell))^2} \quad (12.137)$$

$$= \frac{\lg f(\ell) \exp(-\ell)}{1 + \exp(-\ell)} \quad (12.138)$$

$$= \lg f(\ell) \left( 1 - \frac{1}{1 + \exp(-\ell)} \right) \quad (12.139)$$

$$= \lg f(\ell)(1 - \lg f(\ell)). \quad (12.140)$$

b) By the chain rule for derivatives and (12.140), since  $p_{\alpha,\beta}(x) = \lg f(\beta^T x + \alpha)$ ,

$$\frac{\partial p_{\alpha,\beta}(x)}{\partial x[j]} = \beta[j] p_{\alpha,\beta}(x)(1 - p_{\alpha,\beta}(x)), \quad 1 \leq j \leq d. \quad (12.141)$$

As the probability approaches 0 or 1, the derivative tends to zero, which indicates that the model becomes less sensitive to changes in the features.

c) By the chain rule for derivatives and (12.140), since  $p_{\alpha,\beta}(x) = \lg f(\beta^T x + \alpha)$ ,

$$\frac{\partial p_{\alpha,\beta}(x)}{\partial \beta[j]}(\alpha, \beta) = x[j] p_{\alpha,\beta}(x)(1 - p_{\alpha,\beta}(x)). \quad (12.142)$$

Consequently, by basic calculus

$$\frac{\partial \log \mathcal{L}_{XY}}{\partial \beta[j]}(\alpha, \beta) \quad (12.143)$$

$$= \frac{\partial}{\partial \beta[j]} \left( \sum_{\{i: y_i=0\}} \log(1 - p_{\alpha,\beta}(x_i)) + \sum_{\{l: y_l=1\}} \log p_{\alpha,\beta}(x_l) \right) \quad (12.144)$$

$$= \sum_{\{i: y_i=0\}} x_i \frac{-p_{\alpha,\beta}(x_i)(1 - p_{\alpha,\beta}(x_i))}{(1 - p_{\alpha,\beta}(x_i))} + \sum_{\{l: y_l=1\}} x_l \frac{p_{\alpha,\beta}(x_l)(1 - p_{\alpha,\beta}(x_l))}{p_{\alpha,\beta}(x_l)} \quad (12.145)$$

$$= \sum_{\{i: y_i=0\}} -x_i[j] p_{\alpha,\beta}(x_i) + \sum_{\{l: y_l=1\}} x_l[j](1 - p_{\alpha,\beta}(x_l)) \quad (12.146)$$

where  $y_m$  is the  $m$ th label. The contribution of the  $m$ th data point to the gradient-ascent update is therefore  $\eta x_m[j](y_m - p_{\alpha^{(t-1)}, \beta^{(t-1)}}(x_m))$ . It consists of the  $j$ th feature scaled by a small constant and by the error in the probability estimate for that

data point. When the label equals one, adding the contribution increases the logit  $(\beta^{(t-1)})^T x_j$  for the data point, bringing the estimated probability closer to one. Conversely, when the label equals zero, it decreases the logit, bringing it closer to zero.

#### 12.14 (Building regression trees)

- a) We use induction to prove that the leaves of a regression tree form a partition of the feature space. The induction is applied on the number of leaves  $L$ . The base of induction is a tree with a single leaf,  $L := 1$ , which represents the whole feature space and is clearly a partition (it's a single set that is equal to the whole space).

For the induction step, we assume that for a tree with  $L$  leaves, the corresponding regions  $R_1, \dots, R_L$  form a partition (their intersections are all empty and their union equals their whole space). To add a new leaf, we choose one of the regions in the partition and split it by thresholding one of the features. Let  $R_l$  denote the chosen region,  $j$  the feature used for thresholding and  $\tau$  the threshold. The two new regions are  $R_{l1} := R_l \cap \{x : x_j \leq \tau\}$  and  $R_{l2} := R_l \cap \{x : x_j > \tau\}$ . If we replace  $R_l$  by these two regions, then the resulting set of regions is still a partition. All the regions are disjoint because  $R_{l1}$  and  $R_{l2}$  are disjoint, and they cannot have any points in common with any of the other regions, because the other regions are disjoint with  $R_l$ . In addition, the union of  $R_{l1}$  and  $R_{l2}$  is  $R_l$ , so the union of all the regions still equals the whole feature space. We conclude that the  $L + 1$  leaves in the new tree correspond to regions of the feature space that form a partition, so the proof is complete.

- b) Since  $\mathcal{A}_r(s, t) \cup \mathcal{B}_r(s, t) = R_r$ , we can separate the first term in the expression, as follows

$$\begin{aligned} \Delta \text{RSS}(r, s, t) &= \sum_{\{i: x_i \in R_r\}}^n (y_i - \alpha_r^{\text{tree}})^2 - \sum_{\{i: x_i \in \mathcal{A}_r(s, t)\}}^n (y_i - \alpha_{\mathcal{A}_r(s, t)})^2 \\ &\quad - \sum_{\{i: x_i \in \mathcal{B}_r(s, t)\}}^n (y_i - \alpha_{\mathcal{B}_r(s, t)})^2 \quad (12.147) \\ &= \sum_{\{i: x_i \in \mathcal{A}_r(s, t)\}}^n (y_i - \alpha_r^{\text{tree}})^2 + \sum_{\{i: x_i \in \mathcal{B}_r(s, t)\}}^n (y_i - \alpha_r^{\text{tree}})^2 \\ &\quad - \sum_{\{i: x_i \in \mathcal{A}_r(s, t)\}}^n (y_i - \alpha_{\mathcal{A}_r(s, t)})^2 - \sum_{\{i: x_i \in \mathcal{B}_r(s, t)\}}^n (y_i - \alpha_{\mathcal{B}_r(s, t)})^2. \quad (12.148) \end{aligned}$$

By 7.31

$$\alpha_{\mathcal{A}_r(s, t)} = \arg \min_c \sum_{x_i \in \mathcal{A}_r(s, t)} (x_i - c)^2, \quad (12.149)$$

so in particular

$$\sum_{\{i: x_i \in \mathcal{A}_r(s, t)\}}^n (y_i - \alpha_{\mathcal{A}_r(s, t)})^2 \leq \sum_{\{i: x_i \in \mathcal{A}_r(s, t)\}}^n (y_i - \alpha_r^{\text{tree}})^2. \quad (12.150)$$

For the same reason,

$$\sum_{\{i: x_i \in \mathcal{B}_r(s, t)\}}^n (y_i - \alpha_{\mathcal{B}_r(s, t)})^2 \leq \sum_{\{i: x_i \in \mathcal{B}_r(s, t)\}}^n (y_i - \alpha_r^{\text{tree}})^2. \quad (12.151)$$

Consequently,  $\Delta \text{RSS}(r, s, t) \geq 0$ .

## 12.15 (Tree with categorical features)

Leaf 1:

The mean duration is

$$\alpha_{\text{all}} := \frac{30 + 20 + 50 + 15 + 40 + 40 + 40 + 30}{8} = 33.125. \quad (12.152)$$

The RSS of the single-leaf tree is

$$\triangle \text{RSS} := \sum_{i=1}^8 (y_i - \alpha_{\text{all}})^2 \quad (12.153)$$

$$= (30 - 33.125)^2 + (20 - 33.125)^2 + (50 - 33.125)^2 + (15 - 33.125)^2 \quad (12.154)$$

$$+ (40 - 33.125)^2 + (40 - 33.125)^2 + (40 - 33.125)^2 + (30 - 33.125)^2 \quad (12.155)$$

$$= 946.9. \quad (12.156)$$

Leaf 2:

The mean of the durations for humans is

$$\alpha_h := \frac{20 + 15 + 40 + 30}{4} = 26.25. \quad (12.157)$$

The mean of the durations for animals is

$$\alpha_a := \frac{30 + 50 + 40 + 40}{4} = 40. \quad (12.158)$$

If we threshold according to the categorical variable, then the reduction in RSS is

$$\triangle \text{RSS} = \sum_{i=1}^8 (y_i - \alpha_{\text{all}})^2 - \sum_{\{i:\text{human}\}} (y_i - \alpha_h)^2 - \sum_{\{i:\text{animal}\}} (y_i - \alpha_a)^2 \quad (12.159)$$

$$= 946.9 - \left( (20 - 26.25)^2 + (15 - 26.25)^2 + (40 - 26.25)^2 + (30 - 26.25)^2 \right) \\ - \left( (30 - 40)^2 + (50 - 40)^2 + (40 - 40)^2 + (40 - 40)^2 \right) \\ = 378.1. \quad (12.160)$$

Now we consider thresholding the number of scans using a threshold  $\tau$ .

If  $\tau := 1$ , the mean of the durations when the number of scans is 1 or less is

$$\alpha_{\leq 1} := \frac{30 + 15 + 40}{3} = 28.33. \quad (12.161)$$

The mean of the durations when the number of scans is more than 1 is

$$\alpha_{> 1} := \frac{20 + 50 + 40 + 40 + 30}{5} = 36. \quad (12.162)$$

The reduction in RSS is

$$\Delta \text{RSS} = \sum_{i=1}^8 (y_i - \alpha_{\text{all}})^2 - \sum_{\{i:\text{duration} \leq 1\}} (y_i - \alpha_{\leq 1})^2 - \sum_{\{i:\text{duration} > 1\}} (y_i - \alpha_{> 1})^2 \quad (12.163)$$

$$\begin{aligned} &= 946.9 - \left( (30 - 28.33)^2 + (15 - 28.33)^2 + (40 - 28.33)^2 \right) \\ &\quad - \left( (20 - 36)^2 + (50 - 36)^2 + (40 - 40)^2 + (40 - 36)^2 + (30 - 36)^2 \right) \\ &= 126.2. \end{aligned} \quad (12.164)$$

If  $\tau := 2$ , the mean of the durations when the number of scans is 2 or less is

$$\alpha_{\leq 2} := \frac{30 + 20 + 15 + 40 + 40 + 30}{6} = 29.17. \quad (12.165)$$

The mean of the durations when the number of scans is more than 2 is

$$\alpha_{> 2} := \frac{50 + 40}{2} = 45. \quad (12.166)$$

The reduction in RSS is

$$\Delta \text{RSS} = \sum_{i=1}^8 (y_i - \alpha_{\text{all}})^2 - \sum_{\{i:\text{duration} \leq 1\}} (y_i - \alpha_{\leq 2})^2 - \sum_{\{i:\text{duration} > 1\}} (y_i - \alpha_{> 2})^2 \quad (12.167)$$

$$\begin{aligned} &= 946.9 - \left( (30 - 29.17)^2 + (20 - 29.17)^2 + (15 - 29.17)^2 \right. \\ &\quad \left. + (40 - 29.17)^2 + (40 - 29.17)^2 + (30 - 29.17)^2 \right) \\ &\quad - \left( (50 - 45)^2 + (40 - 45)^2 \right) \end{aligned} \quad (12.168)$$

$$= 376.0. \quad (12.169)$$

The highest reduction in RSS is achieved by thresholding according to the categorical variable. The first leaf corresponds to humans, and the second to animals.

### Leaf 3:

Now we compare the reduction in RSS achieved by thresholding the examples associated with each leaf according to the number of scans using different values of the threshold  $\tau$ . The RSS for the first leaf is:

$$\text{RSS}_h = \sum_{\{i:\text{human}\}} (y_i - \alpha_h)^2 \quad (12.170)$$

$$= (20 - 26.25)^2 + (15 - 26.25)^2 + (40 - 26.25)^2 + (30 - 26.25)^2 \quad (12.171)$$

$$= 368.8 \quad (12.172)$$

If  $\tau := 1$ , the mean of the durations when the number of scans is 1 or less is

$$\alpha_{\leq 1} := 15. \quad (12.173)$$

The mean of the durations when the number of scans is more than 1 is

$$\alpha_{> 1} := \frac{20 + 40 + 30}{3} = 30. \quad (12.174)$$

The reduction in RSS is

$$\Delta \text{RSS} = \text{RSS}_h - \sum_{\{i:\text{duration} \leq 1\}} (y_i - \alpha_{\leq 1})^2 - \sum_{\{i:\text{duration} > 1\}} (y_i - \alpha_{> 1})^2 \quad (12.175)$$

$$\begin{aligned} &= 368.8 - (15 - 15)^2 - \left( (20 - 30)^2 + (40 - 30)^2 + (30 - 30)^2 \right) \\ &= 168.8. \end{aligned} \quad (12.176)$$

If  $\tau := 2$ , the mean of the durations when the number of scans is 2 or less is

$$\alpha_{\leq 2} := \frac{15 + 20 + 30}{3} = 21.67. \quad (12.177)$$

The mean of the durations when the number of scans is more than 2 is

$$\alpha_{> 2} := 40. \quad (12.178)$$

The reduction in RSS is

$$\Delta \text{RSS} = \text{RSS}_h - \sum_{\{i:\text{duration} \leq 2\}} (y_i - \alpha_{\leq 2})^2 - \sum_{\{i:\text{duration} > 2\}} (y_i - \alpha_{> 2})^2 \quad (12.179)$$

$$\begin{aligned} &= 368.8 - \left( (15 - 21.67)^2 + (20 - 21.67)^2 + (30 - 21.67)^2 \right) - (40 - 40)^2 \\ &= 252.1. \end{aligned} \quad (12.180)$$

The RSS for the second leaf is:

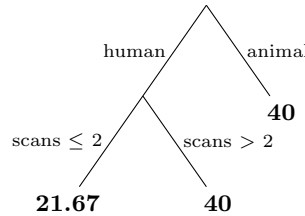
$$\text{RSS}_a = \sum_{\{i:\text{animal}\}} (y_i - \alpha_h)^2 \quad (12.181)$$

$$= (30 - 40)^2 + (50 - 40)^2 + (40 - 40)^2 + (40 - 40)^2 \quad (12.182)$$

$$= 200. \quad (12.183)$$

We don't need to compute the reductions in RSS for the second leaf, as they cannot be greater than 252.1. Therefore, the split has to be associated with the first leaf, and the optimal threshold is  $\tau := 2$ .

We end up with the following regression tree, which estimates a duration of 40 minutes unless the subject is a human and the number of scans is 2 or less, in which case the estimate is 21.67 minutes.



a)

$$\nabla \text{RSS}(\beta) = \frac{1}{2(n-1)} \sum_{i=1}^n \nabla \left( y_i - \beta^T x_i \right)^2 \quad (12.184)$$

$$= \frac{1}{2(n-1)} \sum_{i=1}^n \nabla \left( y_i^2 + \beta^T x_i x_i^T \beta - 2y_i x_i^T \beta \right) \quad (12.185)$$

$$= \frac{1}{2(n-1)} \sum_{i=1}^n \left( 2x_i x_i^T \beta - 2y_i x_i \right) \quad (12.186)$$

$$= \Sigma_X \beta - \Sigma_{XY}. \quad (12.187)$$

Consequently

$$\beta^{(t)} := \beta^{(t-1)} - \eta \Sigma_X \beta^{(t-1)} + \eta \Sigma_{XY}. \quad (12.188)$$

b) If we set  $\beta^{(0)} := 0$ , then

$$\beta^{(1)} = \eta \Sigma_{XY}, \quad (12.189)$$

$$\beta^{(2)} = (I - \eta \Sigma_X) \eta \Sigma_{XY} + \eta \Sigma_{XY}, \quad (12.190)$$

and after  $t$  iterations

$$\beta^{(t)} = \sum_{l=0}^{t-1} (I - \eta \Sigma_X)^l \eta \Sigma_{XY}. \quad (12.191)$$

To alleviate notation we denote diagonal matrices as  $\text{diag}_{j=1}^d (a_i)$ , where  $a_1, \dots, a_d$  are the entries of the diagonal. Setting  $\Sigma_X = U \Lambda U^T$ , where  $U$  contains the eigenvectors and  $\Lambda$  the eigenvalues of  $\Sigma_X$ , and  $I = U U^T$  (which holds since  $U$  is an orthogonal matrix):

$$\beta^{(t)} = \eta \sum_{l=0}^{t-1} \left( U U^T - \eta U \Lambda U^T \right)^l \Sigma_{XY} \quad (12.192)$$

$$= \eta \sum_{l=0}^{t-1} U (I - \eta \Lambda)^l U^T \Sigma_{XY} \quad (12.193)$$

$$= \eta \sum_{l=0}^{t-1} U \text{diag}_{j=1}^d (1 - \eta \xi_j)^l U^T \Sigma_{XY} \quad (12.194)$$

$$= \eta U \text{diag}_{j=1}^d \left( \sum_{l=0}^{t-1} (1 - \eta \xi_j)^l \right) U^T \Sigma_{XY} \quad (12.195)$$

$$= \eta U \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\eta \xi_j} \right) U^T \Sigma_{XY}, \quad (12.196)$$

where  $\xi_j$  is the  $j$ th eigenvalue of  $\Sigma_X$  and the last step follows by the geometric sum

formula. Since  $\beta_{\text{OLS}} = \Sigma_X^{-1} \Sigma_{XY}$  and  $\Sigma_X^{-1} = U \Lambda^{-1} U^T$ , we conclude that

$$\beta^{(t)} = U \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) U^T \Sigma_{XY} \quad (12.197)$$

$$= U \Lambda^{-1} U^T \Sigma_{XY} + U \text{diag}_{j=1}^d \left( \frac{(1 - \eta \xi_j)^t}{\xi_j} \right) U^T \Sigma_{XY} \quad (12.198)$$

$$= \Sigma_X^{-1} \Sigma_{XY} + U \text{diag}_{j=1}^d \left( \frac{(1 - \eta \xi_j)^t}{\xi_j} \right) U^T \Sigma_{XY} \quad (12.199)$$

$$= \beta_{\text{OLS}} + U \text{diag}_{j=1}^d \left( \frac{(1 - \eta \xi_j)^t}{\xi_j} \right) U^T \Sigma_{XY}. \quad (12.200)$$

Assuming  $\eta < 1/\xi_{\min} < 1/\xi_j$ ,  $\eta \xi_j < 1$  so that  $(1 - \eta \xi_j)^t \rightarrow 0$  as  $t \rightarrow \infty$  for all  $j$ . Consequently  $\beta^{(t)} \rightarrow \beta_{\text{OLS}}$  as  $t \rightarrow \infty$ .

- c) Under the assumptions  $\tilde{\Sigma}_{XY} = \Sigma_X \beta_{\text{true}} + \tilde{\Sigma}_{XZ}$  (see (12.154) in the main text), so by (12.197)

$$\tilde{\beta}^{(t)} = U \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) U^T \left( \Sigma_X \beta_{\text{true}} + \tilde{\Sigma}_{XZ} \right), \quad (12.201)$$

which implies

$$\tilde{c}_{\text{GD}} = \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) U^T \left( \Sigma_X \beta_{\text{true}} + \tilde{\Sigma}_{XZ} \right). \quad (12.202)$$

Since  $\mathbb{E} [\tilde{\Sigma}_{XZ}] = 0$  (see (12.167) in the main text), by linearity of expectation,

$$\mathbb{E} [\tilde{c}_{\text{GD}}] = \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) U^T \Sigma_X \beta_{\text{true}} \quad (12.203)$$

$$= \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) U^T U \Lambda U^T \beta_{\text{true}} \quad (12.204)$$

$$= \text{diag}_{j=1}^d \left( 1 - (1 - \eta \xi_j)^t \right) c_{\text{true}}, \quad (12.205)$$

so that

$$\mathbb{E} [\tilde{c}_{\text{GD}}[j]] = \left( 1 - (1 - \eta \xi_j)^t \right) c_{\text{true}}[j]. \quad (12.206)$$

- d) By (12.202),  $\mathbb{E} [\tilde{\Sigma}_{XZ}] = 0$  and  $\mathbb{E} [\tilde{\Sigma}_{XZ} \tilde{\Sigma}_{XZ}^T] = \frac{\sigma^2}{n-1} \Sigma_X = \frac{\sigma^2}{n-1} U \Lambda U^T$  (see (12.177))

in the main paper),

$$\Sigma_{\tilde{c}_{\text{GD}}} = \mathbb{E} \left[ \text{ct}(\tilde{c}_{\text{GD}}) \text{ct}(\tilde{c}_{\text{GD}})^T \right] \quad (12.207)$$

$$\begin{aligned} &= \mathbb{E} \left[ \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) U^T \tilde{\Sigma}_{XZ} \tilde{\Sigma}_{XZ}^T U \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) \right] \\ &= \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) U^T \mathbb{E} \left[ \tilde{\Sigma}_{XZ} \tilde{\Sigma}_{XZ}^T \right] U \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) \\ &= \frac{\sigma^2}{n-1} \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) U^T U \Lambda U^T U \text{diag}_{j=1}^d \left( \frac{1 - (1 - \eta \xi_j)^t}{\xi_j} \right) \\ &= \frac{\sigma^2}{n-1} \text{diag}_{j=1}^d \left( \frac{\left( 1 - (1 - \eta \xi_j)^t \right)^2}{\xi_j} \right). \end{aligned} \quad (12.208)$$

We conclude that

$$\text{Var}[\tilde{c}_{\text{GD}}[j]] = \frac{\sigma^2}{n-1} \left( \frac{\left( 1 - (1 - \eta \xi_j)^t \right)^2}{\xi_j} \right). \quad (12.209)$$

e) We pick  $t$  small enough to ensure that for the small eigenvalue

$$(1 - \eta \xi_{\min})^t \approx 1, \quad (12.210)$$

but for the large eigenvalues  $\xi_{\text{large}}$

$$(1 - \eta \xi_{\text{large}})^t \ll 1, \quad (12.211)$$

which is possible because  $\xi_{\text{large}} \gg \xi_{\min}$ . For the index  $j_{\min}$  corresponding to the small eigenvalue, by (12.206) and (12.209),

$$\mathbb{E}[\tilde{c}_{\text{GD}}[j_{\min}]] \approx 0, \quad (12.212)$$

$$\text{Var}[\tilde{c}_{\text{GD}}[j_{\min}]] \approx 0, \quad (12.213)$$

which introduces some bias, but avoids noise amplification due to the small eigenvalue. For any of the remaining indices, which we denote by  $j_{\text{large}}$ ,

$$\mathbb{E}[\tilde{c}_{\text{GD}}[j_{\text{large}}]] \approx c_{\text{true}}[j], \quad (12.214)$$

$$\text{Var}[\tilde{c}_{\text{GD}}[j_{\text{large}}]] \approx \frac{\sigma^2}{(n-1)\xi_{\text{large}}}, \quad (12.215)$$

which does not produce noise amplification, because  $\xi_{\text{large}}$  is large.

#### 12.17 (Product)

a) The mean of the features is equal to zero, because they equal 1 or -1 with the same probability. By the independence assumption, the mean of the response equals

$$\mathbb{E}[\tilde{y}] = \sum_{x[1] \in \{-1, 1\}} \sum_{x[2] \in \{-1, 1\}} x[1]x[2]p_{\tilde{x}}(x[1], x[2]) \quad (12.216)$$

$$= p_{\tilde{x}}(-1, -1) - p_{\tilde{x}}(-1, 1) - p_{\tilde{x}}(1, -1) + p_{\tilde{x}}(1, 1) \quad (12.217)$$

$$\begin{aligned} &= p_{\tilde{x}[1]}(-1)p_{\tilde{x}[2]}(-1) - p_{\tilde{x}[1]}(-1)p_{\tilde{x}[2]}(1) - p_{\tilde{x}[1]}(1)p_{\tilde{x}[2]}(-1) + p_{\tilde{x}[1]}(1)p_{\tilde{x}[2]}(1) \\ &= 0. \end{aligned} \quad (12.218)$$



The covariance between the response and each of the features equals

$$\text{Cov}[\tilde{x}[1], \tilde{y}] = \mathbb{E}[\tilde{x}[1]\tilde{y}] \quad (12.219)$$

$$= \mathbb{E}[\tilde{x}[1]^2 \tilde{x}[2]] \quad (12.220)$$

$$= \mathbb{E}[\tilde{x}[2]] \quad (12.221)$$

$$= 0, \quad (12.222)$$

because  $\tilde{x}[1]^2 = 1$  with probability one. By the same reasoning,  $\tilde{x}[2]$  and  $\tilde{y}$  are also uncorrelated. Consequently, the linear MMSE estimator is

$$\ell_{\text{MMSE}}(\tilde{x}) = \Sigma_{\tilde{x}\tilde{y}}^T \Sigma_{\tilde{x}}^{-1} \tilde{x} = 0, \quad (12.223)$$

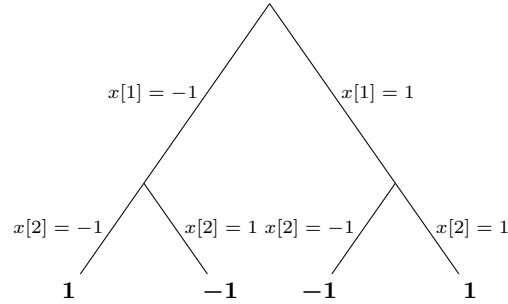
for all values of the features. The corresponding MSE is

$$\mathbb{E}[(\tilde{y} - \ell_{\text{MMSE}}(\tilde{x}))^2] = \mathbb{E}[\tilde{y}^2] \quad (12.224)$$

$$= \mathbb{E}[\tilde{x}[1]^2 \tilde{x}[2]^2] \quad (12.225)$$

$$= 1. \quad (12.226)$$

- b) The minimum MSE estimator of  $\tilde{y}$  given  $\tilde{x}[1] = x_1$  and  $\tilde{x}[2] = x_2$ , is the conditional mean function which equals  $x_1 x_2$  and is exactly equal to  $\tilde{y}$ . Therefore the corresponding MSE is zero. The estimator is perfect!
- c) The regression tree is shown below.



- d) There are many different neural networks that implement the minimum MSE estimator. A possibility is to have a first linear layer with four hidden variables that equal

$$h[1] := \frac{x[1] + x[2]}{2}, \quad (12.227)$$

$$h[2] := -\frac{x[1] + x[2]}{2}, \quad (12.228)$$

$$h[3] := \frac{x[1] - x[2]}{2}, \quad (12.229)$$

$$h[4] := \frac{-x[1] + x[2]}{2}, \quad (12.230)$$

followed by a ReLU to produce the rectified hidden variables  $h_+$ . Due to the ReLU, if  $x[1] = 1$  and  $x[2] = 1$ , then  $h_+[1] = 1$  and the rest of rectified hidden variables are all zero. If  $x[1] = -1$  and  $x[2] = -1$ ,  $h_+[2] = 1$  and the rest are zero. If  $x[1] = 1$  and  $x[2] = -1$ ,  $h_+[3] = 1$  and the rest are zero. If  $x[1] = -1$  and  $x[2] = 1$ ,  $h_+[4] = 1$  and

the rest are zero. Consequently, we can implement the minimum MSE estimator by weighting the rectified hidden variables as follows in the last layer:

$$n\ell_{\text{nnet}}(x) := h_+[1] + h_+[2] - h_+[3] - h_+[4]. \quad (12.231)$$

#### 12.18 (Diagnosing cancer)

- a) The model is probably learning features that enable it to identify each patient and memorize its corresponding label, but have nothing to do with cancer recurrence. This explains the perfect accuracy on the test set, which includes images from the same patients as the training set, and its bad performance on data from new patients.
- b) Laura can apply the model setting the additional input to *treatment*, and also to *no treatment*. If *treatment* leads to a lower predicted probability of recurrence than *no treatment*, then she concludes that the patient would benefit from the treatment. However, this only works if the treatment in the training set is randomized. Otherwise, the conditional probabilities estimated by the model cannot be interpreted causally, because the treatment and the potential outcomes associated with recurrence are not necessarily independent (there could be confounders).

#### 12.19 (ROC curve of random classifier)

- a) The probability that we select each of the  $P$  positive examples is  $p_{\tilde{x}_+}(x_i) = 1/P$  for any positive example  $x_i$ , so

$$P(\tilde{s}_+ > \tau) = P(p_{\text{pred}}(\tilde{x}_+) > \tau) \quad (12.232)$$

$$= \sum_{\{i: y_i = +, p_{\text{pred}}(x_i) > \tau\}} p_{\tilde{x}_+}(x_i) \quad (12.233)$$

$$= \frac{\text{TP}}{P} = \text{TPR}(\tau), \quad (12.234)$$

since the positive examples classified as positive are exactly the true positives (TP). Similarly, since  $p_{\tilde{x}_-}(x_i) = 1/N$  for any negative example  $x_i$ ,

$$P(\tilde{s}_- > \tau) = P(p_{\text{pred}}(\tilde{x}_-) > \tau) \quad (12.235)$$

$$= \sum_{\{i: y_i = -, p_{\text{pred}}(x_i) > \tau\}} p_{\tilde{x}_-}(x_i) \quad (12.236)$$

$$= \frac{\text{FP}}{N} = \text{FPR}(\tau). \quad (12.237)$$

- b) By definition of the random classifier, the predicted probabilities are independent samples from a uniform distribution between 0 and 1. Consequently, for a large number of examples, the histograms of the probabilities associated with positive and negative labels will both approximate a uniform distribution between 0 and 1. In that case,  $\tilde{s}_+ := p_{\text{pred}}(\tilde{x}_+)$  and  $\tilde{s}_- := p_{\text{pred}}(\tilde{x}_-)$  are both approximately uniformly distributed between 0 and 1. As a result, by (12.234)

$$\text{TPR}(\tau) = P(\tilde{s}_+ > \tau) \quad (12.238)$$

$$\approx \int_{\tau}^1 du \quad (12.239)$$

$$= 1 - \tau, \quad (12.240)$$

and by (12.237)  $\text{FPR}(\tau) \approx 1 - \tau$ . We conclude that  $\text{TPR}(\tau) = \text{FPR}(\tau)$  for all  $\tau$ , so that the ROC curve of the classifier approximates a diagonal line from the origin ( $\tau := 1$ ) to  $(1, 1)$  ( $\tau := 0$ ).

## 12.20 (Probabilistic interpretation of AUC)

- a) The area under the ROC curve can be computed by integrating the TPR as a function of the FPR. To compute the integral, we express both quantities as a function of the threshold  $\tau$  that we apply to the predicted probability. For  $\tau = 1$ ,  $\text{FPR} = 0$  and  $\text{TPR} = 0$ . Then as  $\tau$  decreases, the two quantities never decrease and eventually reach  $\text{FPR} = 1$  and  $\text{TPR} = 1$ . Consequently, we can express the area under the ROC curve as

$$\text{AUC} = \int_{\tau=1}^0 \text{TPR}(\tau) \, d\text{FPR}(\tau). \quad (12.241)$$

By (12.237) and basic calculus, the differential of FPR equals

$$d\text{FPR}(\tau) = \frac{d\text{FPR}(\tau)}{d\tau} d\tau \quad (12.242)$$

$$= \frac{d(1 - F_{\tilde{s}_-}(\tau))}{d\tau} d\tau \quad (12.243)$$

$$= -f_{\tilde{s}_-}(\tau) d\tau. \quad (12.244)$$

Combining this with (12.234) and performing the change of variable  $\tau = -t$ , we obtain

$$\text{AUC} = \int_{t=0}^1 (1 - F_{\tilde{s}_+}(t)) f_{\tilde{s}_-}(t) dt. \quad (12.245)$$

- b) Since  $\tilde{x}_+$  and  $\tilde{x}_-$  are independent, then so are  $\tilde{s}_+$  and  $\tilde{s}_-$ , so

$$P(\tilde{s}_+ > \tilde{s}_-) = \int_{t=0}^1 \int_{a=t}^1 f_{\tilde{s}_+, \tilde{s}_-}(a, t) dt da \quad (12.246)$$

$$= \int_{t=0}^1 \int_{a=t}^1 f_{\tilde{s}_+}(a) f_{\tilde{s}_-}(t) dt da \quad (12.247)$$

$$= \int_{t=0}^1 (1 - F_{\tilde{s}_+}(t)) f_{\tilde{s}_-}(t) dt \quad (12.248)$$

$$= \text{AUC}, \quad (12.249)$$

where the last step follows from (12.245).

- c) For the random classifier, the ROC is a diagonal line from (0,0) to (1,1), so the area under the curve is  $1/2$ . Since  $\tilde{s}_+ := p_{\text{pred}}(\tilde{x}_+)$  and  $\tilde{s}_- := p_{\text{pred}}(\tilde{x}_-)$  are both approximately uniformly distributed between 0 and 1, the probability that the score for the positive example is larger than that of the negative example, is  $1/2$ ,

$$P(\tilde{s}_+ > \tilde{s}_-) = \int_{t=0}^1 \int_{a=t}^1 dt da \quad (12.250)$$

$$= \int_{t=0}^1 (1 - t) dt \quad (12.251)$$

$$= \frac{1}{2}. \quad (12.252)$$

- d) We can approximate  $P(\tilde{s}_+ > \tilde{s}_-)$  via the Monte Carlo method, by randomly sampling a large number of pairs, each consisting of a positive and a negative example, from the data and then computing the fraction of them for which the score associated with the positive example is larger. For the example considered in Section 12.9, using  $10^5$  pairs produces an estimate equal to 0.8470, which is very close to the value of the AUC obtained by computing the area under the ROC curve (0.8472).

## 12.21 (F1 score)

$$F_1 = \frac{2 \cdot \text{TPR} \cdot \text{Precision}}{\text{TPR} + \text{Precision}} \quad (12.253)$$

$$= \frac{1}{\frac{1}{2} \cdot \frac{1}{\text{TPR}} + \frac{1}{2} \cdot \frac{1}{\text{Precision}}} \quad (12.254)$$

$$= \frac{1}{\frac{1}{2} \cdot \frac{\text{TP} + \text{FN}}{\text{TP}} + \frac{1}{2} \cdot \frac{\text{TP} + \text{FP}}{\text{TP}}} \quad (12.255)$$

$$= \frac{\text{TP}}{\text{TP} + \frac{\text{FN}}{2} + \frac{\text{FP}}{2}} \quad (12.256)$$

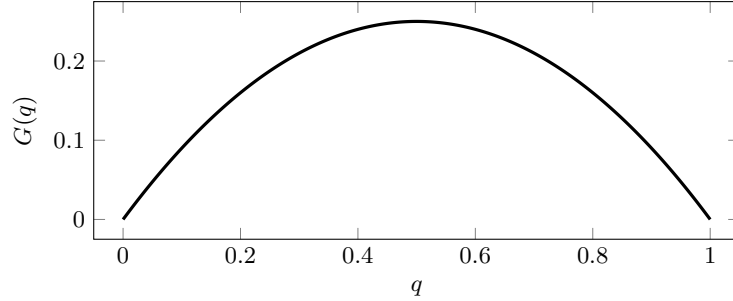
$$= \frac{\text{TP}}{\text{TP} + \text{AFE}}. \quad (12.257)$$

## 12.22 (Decomposition of the Brier score)

- a) As explained in Section 12.9.2, we evaluate calibration by dividing the unit interval in bins, assigning the probability estimates to each bin, and then comparing the estimated probabilities in each bin to the empirical probability (fraction of data with label 1) in each bin. This is exactly what the calibration component does, it sums the squared differences between the estimated probabilities ( $p_b$ ) and the empirical probabilities ( $q_b$ ) weighted by the fraction of the total data that belongs to each bin ( $n_b/n$ ). Consequently, a smaller calibration component corresponds to a better calibrated classifier.
- b) The graph below shows the Gini index

$$G(q) := q(1 - q) \quad (12.258)$$

for  $q \in [0, 1]$ .



Similar to the entropy, the Gini index is highest when  $q$  and  $1 - q$  are the same, and lowest when one equals zero and the other equals one. The discrimination component in the Brier score is equal to the Gini index of the empirical probability of one in each bin. The component is small if the empirical probability in each bin is either close to zero or close to one, which means that the classifier is able to separate the data according to the labels. Consequently, a smaller discrimination component corresponds to a more discriminative classifier.

- c) We begin by reformulating the Brier score in terms of the bins:

$$\text{Brier score} := \frac{1}{n} \sum_{i=1}^n (y_i - p_i)^2 \quad (12.259)$$

$$= \frac{1}{n} \sum_{b=1}^m \sum_{\{p_i \in \mathcal{B}_b\}} (y_i - p_i)^2. \quad (12.260)$$

By the assumption that the probability estimate in the  $b$ th bin is approximately equal to  $p_b$ ,

$$\sum_{\{p_i \in \mathcal{B}_b\}} (y_i - p_i)^2 \quad (12.261)$$

$$\approx \sum_{\{p_i \in \mathcal{B}_b\}} (y_i - p_b)^2 \quad (12.262)$$

$$= \sum_{\{p_i \in \mathcal{B}_b\}} (y_i - p_b + q_b - q_b)^2 \quad (12.263)$$

$$= \sum_{\{p_i \in \mathcal{B}_b\}} (q_b - p_b)^2 + \sum_{\{p_i \in \mathcal{B}_b\}} (y_i - q_b)^2 + 2(q_b - p_b) \sum_{\{p_i \in \mathcal{B}_b\}} (y_i - q_b). \quad (12.264)$$

The first term in (12.264) is equal to the calibration component,

$$\sum_{\{p_i \in \mathcal{B}_b\}} (q_b - p_b)^2 = n_b (q_b - p_b)^2. \quad (12.265)$$

Since the labels are equal to zero or one,  $y_i^2 = y_i$  for  $1 \leq i \leq n$ , so

$$\frac{1}{n_b} \sum_{\{p_i \in \mathcal{B}_b\}} y_i^2 = \frac{1}{n_b} \sum_{\{p_i \in \mathcal{B}_b\}} y_i \quad (12.266)$$

$$:= q_b. \quad (12.267)$$

Consequently, the second term in (12.264) is equal to the discrimination component,

$$\sum_{\{p_i \in \mathcal{B}_b\}} (y_i - q_b)^2 = \sum_{\{p_i \in \mathcal{B}_b\}} y_i^2 - 2q_b \sum_{\{p_i \in \mathcal{B}_b\}} y_i + n_b q_b^2 \quad (12.268)$$

$$= \sum_{\{p_i \in \mathcal{B}_b\}} (y_i^2 - 2q_b y_i + q_b^2) \quad (12.269)$$

$$= n_b q_b - 2n_b q_b^2 + n_b q_b^2 \quad (12.270)$$

$$= n_b q_b (1 - q_b). \quad (12.271)$$

Finally, the third term in (12.264) is zero by the definition of  $q_b$ ,

$$\sum_{\{p_i \in \mathcal{B}_b\}} (y_i - q_b) = n_b q_b - n_b q_b = 0. \quad (12.272)$$

- d) The perfectly discriminative classifier only produces two different probability estimates, so we can divide them into two bins ( $m = 2$ ). Bin 1 contains the probability estimates equal to 0.8, all of which have label zero, and bin 2 the estimates equal to 0.9, all of which have label one. Consequently, the empirical probabilities in each bin are zero and one, respectively:

$$q_1 := \frac{1}{n_1} \sum_{\{p_i \in \mathcal{B}_1\}} y_i = 0, \quad (12.273)$$

$$q_2 := \frac{1}{n_2} \sum_{\{p_i \in \mathcal{B}_2\}} y_i = \frac{n_2}{n_2} = 1. \quad (12.274)$$

As a result, the discrimination term is equal to zero. This means that the calibration term must equal the Brier score (0.504).

- e) The perfectly calibrated classifier assigns every data point to the same estimate, so we only need a single bin ( $m = 1$ ). The empirical probability  $q_1$  of label 1 in the bin is equal to the fraction of ones by definition, which is equal to the probability estimate by assumption. Consequently,  $p_1 = q_1$  and the calibration component is zero. This means that the discrimination component must equal 0.169.