

## Theory of Inference: Homework 2

1. Study the statement and proof of the Bayes Rule theorem. Now state and prove the theorem in the special case where the choice of action does not affect the client's beliefs about  $X$ .

**Answer.** Define

$$\delta^*(y) := \operatorname{argmin}_{a \in \mathcal{A}} \mathbb{E}\{L(a, X) \mid Y \doteq y\},$$

note that there is no ' $a$ ' in the expectation, because the action does not affect beliefs about  $X$ . We want to show that, for arbitrary rule  $\delta$ ,

$$\mathbb{E}\{L(\delta(Y), X)\} \geq \mathbb{E}\{L(\delta^*(Y), X)\},$$

which would prove that  $\delta^*$  was a Bayes Rule. Thus:

$$\begin{aligned} \mathbb{E}\{L(\delta(Y), X)\} &= \sum_x \sum_y L(\delta(y), x) p(x, y) && \text{FTP} \\ &= \sum_y \sum_x L(\delta(y), x) p(x \mid y) p(y) && \text{definition of } p(x \mid y) \\ &\geq \sum_y \left\{ \min_{a \in \mathcal{A}} \sum_x L(a, x) p(x \mid y) \right\} p(y) \\ &= \sum_y \left\{ \min_{a \in \mathcal{A}} \mathbb{E}\{L(a, X) \mid Y \doteq y\} \right\} p(y) && \text{CFTP} \\ &= \sum_y \mathbb{E}\{L(\delta^*(y), X) \mid Y \doteq y\} p(y) && \text{by definition} \\ &= \sum_y \left\{ \sum_x L(\delta^*(y), x) p(x \mid y) \right\} p(y) && \text{CFTP again} \\ &= \sum_y \sum_x L(\delta^*(y), x) p(x, y) && p(x \mid y) \text{ again} \\ &= \mathbb{E}\{L(\delta^*(Y), X)\} && \text{FTP again} \end{aligned}$$

as required. This is overkill on the number of steps; about four would do.

2. A volcano can be either inactive or active; it is active with probability  $\theta$ . If it is inactive its eruption rate is zero. If it is active, its eruption rate

has a Gamma distribution with shape parameter  $\alpha$  and rate parameter  $\beta$ . Let  $\lambda$  be its eruption rate. Show that

$$\Pr(\lambda \leq v) = 1 - \theta + F(v; \alpha, \beta) \cdot \theta$$

for  $v \geq 0$ , and zero otherwise, where  $F$  is the distribution function of the Gamma distribution. Hint: introduce the random quantity  $A \in \{0, 1\}$ , where  $A = 1$  exactly when the volcano is active, and use the Law of Total Probability.

**Answer.** Introduce  $A$  as in the hint, so that  $\{A \doteq 0, A \doteq 1\}$  is a belief partition. Let  $\gamma \sim \text{Gamma}(\cdot; \alpha, \beta)$ . Then  $\lambda = A\gamma$ , and

$$\begin{aligned} \Pr(\lambda \leq v) &= \Pr(\lambda \leq v \mid A \doteq 0) \Pr(A \doteq 0) + \Pr(\lambda \leq v \mid A \doteq 1) \Pr(A \doteq 1) && \text{LTP} \\ &= \Pr(A\gamma \leq v \mid A \doteq 0) \Pr(A \doteq 0) + \Pr(A\gamma \leq v \mid A \doteq 1) \Pr(A \doteq 1) \\ &= \Pr(0 \leq v \mid A \doteq 0) \Pr(A \doteq 0) + \Pr(\gamma \leq v \mid A \doteq 1) \Pr(A \doteq 1) && \text{by TOWIK} \\ &= 1 \cdot \Pr(A \doteq 0) + F(v; \alpha, \beta) \cdot \Pr(A \doteq 1) && \text{as } v \geq 0 \\ &= 1 - \theta + F(v; \alpha, \beta) \cdot \theta \end{aligned}$$

as required.

3. Here is an exam-style revision question on the first strand of lectures.

- (a) The Fundamental Theorem of Prevision (FTP) is an if-and-only-if theorem. Pick one branch (either if or only-if), state it, and prove it. [5 marks]

**Answer.** All of these answers are covered in the notes; I'll produce off-the-cuff answers here without reference to the notes—just for fun.

One branch of the FTP states that if  $E$  is coherent then there exists a  $\mathbf{p} \in \mathbb{S}^{s-1}$  for which  $E\{g(\mathbf{X})\} = \sum_j g(\mathbf{x}^{(j)}) \cdot p_j$  for any specified  $g$ , where  $s$  is the dimension of the realm of  $\mathbf{X}$ , which is  $\mathcal{X} := \{\mathbf{x}^{(j)}\}_{j=1}^s$ . To prove this, note that

$$1 = \sum_j \mathbb{1}_{\mathbf{X} \doteq \mathbf{x}^{(j)}},$$

because  $\mathbf{X}$  must take exactly one value from its realm. And then

$$\begin{aligned}
\mathbb{E}\{g(\mathbf{X})\} &= \mathbb{E}\{g(\mathbf{X}) \cdot 1\} \\
&= \mathbb{E}\left\{g(\mathbf{X}) \sum_j \mathbf{1}_{\mathbf{X}=\mathbf{x}^{(j)}}\right\} \\
&= \mathbb{E}\left\{\sum_j g(\mathbf{X}) \mathbf{1}_{\mathbf{X}=\mathbf{x}^{(j)}}\right\} \\
&= \mathbb{E}\left\{\sum_j g(\mathbf{x}^{(j)}) \mathbf{1}_{\mathbf{X}=\mathbf{x}^{(j)}}\right\} \quad (\text{make sure you understand this step}) \\
&= \sum_j g(\mathbf{x}^{(j)}) \mathbb{E}\{\mathbf{1}_{\mathbf{X}=\mathbf{x}^{(j)}}\} \quad \text{by linearity} \\
&= \sum_j g(\mathbf{x}^{(j)}) \cdot p_j
\end{aligned}$$

say. We have  $p_j \geq 0$  by lower-boundedness, and also

$$1 = \mathbb{E}\{1\} = \mathbb{E}\left\{\sum_j \mathbf{1}_{\mathbf{X}=\mathbf{x}^{(j)}}\right\} = \sum_j \mathbb{E}\{\mathbf{1}_{\mathbf{X}=\mathbf{x}^{(j)}}\} = \sum_j p_j$$

by lower-boundedness and linearity, showing that  $\mathbf{p} := (p_1, \dots, p_s)$  is a point in  $\mathbb{S}^{s-1}$ , as required.

- (b) Outline a model for data, distinguishing between random quantities, observables, and observations. Give an example of how an observable differs from the random quantity it measures, owing to limitations in the instrument. [5 marks]

**Answer.** Suppose we have  $m$  random quantities  $\mathbf{X} := (X_1, \dots, X_m)$ . An ‘observable’  $Y_i$  is a specified function of  $\mathbf{X}$ , say  $g_i(\mathbf{X})$ , whose value will at some point become known; its value is the ‘observation’  $y_i^{\text{obs}}$ . ‘Data’ can collectively be represented as the truth of the random proposition

$$Q := q(\mathbf{X}) \quad \text{where} \quad q(\mathbf{x}) := \bigwedge_{i=1}^n (g_i(\mathbf{x}) \doteq y_i^{\text{obs}}).$$

In the simplest possible case,  $g_i(\mathbf{x}) \leftarrow x_i$ , in which case the data comprise observations on the first  $n$  random quantities: call this the *Simple Observation Model (SOM)*. More generally,  $Y_i$  might be a more complicated function of  $X_i$ , reflecting in part the measuring instrument’s limitations. For example, a meter which only reads up to the value  $v$  has  $g_i(\mathbf{x}) = \min\{x_i, v\}$ .

- (c) State and prove the Muddy Table theorem. Illustrate it with a diagram. [5 marks]

**Answer.** Let  $q(\mathbf{x})$  be a first order sentence, and define the random proposition  $Q := q(\mathbf{X})$ . The Muddle Table theorem states that

$$\Pr(\mathbf{X} \doteq \mathbf{x} \mid Q) = \frac{\mathbb{1}_{q(\mathbf{x})} p(\mathbf{x})}{\Pr(Q)} \quad \text{provided that } \Pr(Q) > 0,$$

where  $p(\cdot)$  represents the PMF of  $\mathbf{X}$  and  $\Pr(Q) = \sum_{\mathbf{x}} \mathbb{1}_{q(\mathbf{x})} p(\mathbf{x})$  by the Fundamental Theorem of Prevision (FTP).

For the proof, take  $\Pr(Q) > 0$  as stated, and then

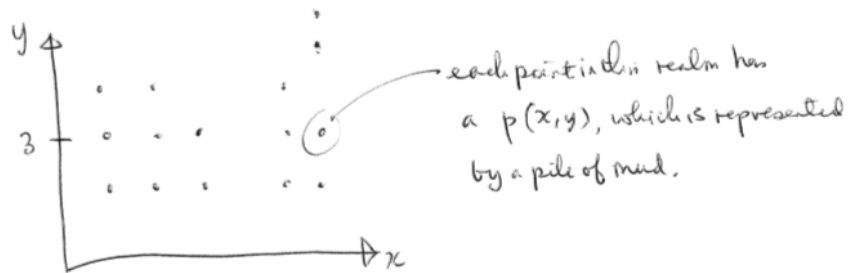
$$\Pr(\mathbf{X} \doteq \mathbf{x} \mid Q) = \frac{\Pr(\mathbf{X} \doteq \mathbf{x}, Q)}{\Pr(Q)}$$

according to the definition of hypothetical/conditional probability. Evaluate the numerator using the FTP to get

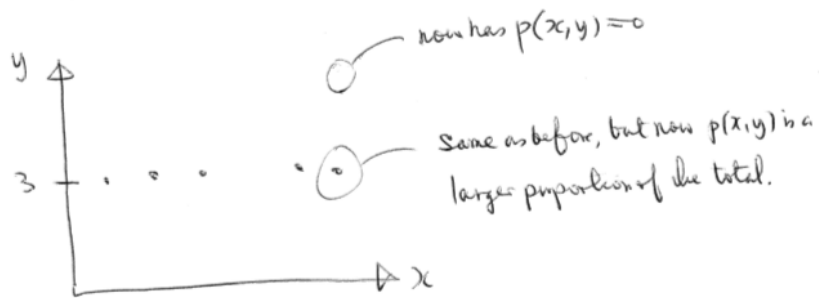
$$\Pr(\mathbf{X} \doteq \mathbf{x}, Q) = \sum_{\mathbf{x}'} \mathbb{1}_{\mathbf{x}' \doteq \mathbf{x} \wedge q(\mathbf{x}')} p(\mathbf{x}') = \mathbb{1}_{q(\mathbf{x})} p(\mathbf{x})$$

as required. Here is my diagram:

Let  $\underline{X} = (X, Y)$  with realm shown by the dots below:



Suppose that  $q(x) = (y = 3)$ . All the points incompatible with  $Y=3$  are swept clear of mud, leaving



- (d) 'Bayesian conditionalisation' is the name given to a model for learning in which new information is incorporated into beliefs through conditioning. What are the attractive features of this model? Why does it *not* describe the typical practice of statistical inference? [10 marks]

**Answer.** Here are some bullet points for this answer:

- Bayesian conditionalisation asserts that beliefs are represented as probability statements and updated by conditioning. A set of propositions  $\Psi$  is already known to be true, and another proposition  $Q$  is newly learnt to be true, and is to be added to  $\Psi$ . Before knowing  $Q$ , beliefs have the form  $\text{Bel}_\Psi(\cdot) = \Pr(\cdot | \Psi)$ . After knowing  $Q$ , updated beliefs have the form  $\text{Bel}_{Q \wedge \Psi}(\cdot) = \Pr(\cdot | Q, \Psi)$ ,

provided that  $\Pr(Q | \Psi) > 0$ .

- One attractive feature of this approach is that it has a very simple form: it is easy to compute the updated beliefs by applying sequential Bayes's theorem. Other attractive features include:
  - i. It is a theorem that conditional probabilities represented as  $\Pr(Q | R)$  in the defining relationship

$$\Pr(Q, R) = \Pr(Q | R) \Pr(R)$$

are indeed probabilities (i.e. satisfy the axioms of expectation), so this type of update is 'coherence preserving'.

- ii. As a special case of the above, the update satisfies the minimal property of being logically consistent, so that if  $Q$  implies  $R$  then  $\Pr(R | Q) = 1$ .
  - iii. The update is order-invariant, so that if  $Q = Q_1 \wedge Q_2$ , then the same beliefs arise from conditioning first on  $Q_1$  and then on  $Q_2$ , or the other way around: in other words, only  $\Psi$  matters, not the order in which its elements were acquired.
- The learning model has a simple sequential form in which we start with some beliefs  $\Psi$ , acquire new knowledge  $Q$ , and then update to  $Q \wedge \Psi$ . Practically speaking, it can only work if we anticipate  $Q$ , in order that  $\Pr(Q | \Psi) > 0$ . It is very hard to anticipate everything that one might learn.

More generally, in a statistical inference the nature of  $Q$  is itself shaped through dialogue with the client, and analysis of the datapool. Often the  $Q$  that is selected is only a small subset of the datapool, reflecting the statistician's and the client's limitations in specifying defensible beliefs about many of the elements of the datapool. This 'backwards and forwards' activity to choose  $Q$  is quite different from the simple sequential form represented in Bayesian conditionalisation.

Please hand in your answers for marking next Tue (3 Mar), at the lecture or by 5pm in the box outside my office door. I will return them in Thu's

lecture.

Jonathan Rougier

Feb 2015