

These are exercises for self-study and exam preparation. All material is examinable unless otherwise mentioned.

Exercise 1. Optimal actions for different loss functions

Assume that $h \in \{1, ..., 5\}$ with probabilities f(h) equal to

$$f(h=1) = 0.2$$
, $f(h=2) = 0.3$, $f(h=3) = 0$, $f(h=4) = 0.1$, $f(h=5) = 0.4$ (1)

Compute the action a that minimises the expected loss $\mathbb{E}_{f(h)}[\ell(h,a)]$ when:

(a)
$$\ell(h, a) = (h - a)^2$$

Solution. The optimal action $a^*(f)$ for the quadratic loss is given by the mean $\mathbb{E}_{f(h)}[h]$, which equals

$$\mathbb{E}_{f(h)}[h] = \sum_{h} hp(h) \tag{S.1}$$

$$= 1 \cdot 0.2 + 2 \cdot 0.3 + 3 \cdot 0 + 4 \cdot 0.1 + 5 \cdot 0.4 \tag{S.2}$$

$$=3.2\tag{S.3}$$

$$(b) \ \ell(h,a) = |h - a|$$

Solution. The optimal action $a^*(f)$ for the absolute error loss is given by the median of f. To compute the median we compute the cumulative distribution function (cdf) $F(\alpha) = \mathbb{P}(h \leq \alpha)$. Even for discrete random variables, the cdf is defined for all values of α . For discrete h, the cdf thus will be piecewise constant with jump points at

$$F(1) = 0.2, \quad F(2) = 0.5, \quad F(4) = 0.6, \quad F(5) = 1$$
 (S.4)

The median is any number m that satisfies $F(m^-) \leq 0.5 \leq F(m)$ where m^- is the lower limit.

Since $F(\alpha)$ is equal to 0.5 on the interval [2,4] with a jump point at 4, the median is here not unique, any number in the interval [2,4] satisfies the definition to be a median. The endpoint 4 is included since $F(4^-) = 0.5 \le 0.5$ and $F(4) = 0.6 \ge 0.5$. There are several approaches to deal with the non-uniqueness: One is to report the whole interval [2,4]; if a single number is required, a popular approach is to report the smallest median, i.e. 2 here.

(c)
$$\ell(h, a) = 1 - \mathbb{1}(h = a)$$

Solution. The optimal action $a^*(f)$ for the zero-one loss is given by the mode of f, i.e. $\operatorname{argmax}_h f(h)$, which is here h = 5.

(d) $\ell(h,a) = -\log a(h)$ where choosing an action a means choosing a distribution over h.

Solution. The optimal action for the log-loss when the expectation is taken over f is f itself.

Exercise 2. Causal decision theory applied to the kidney stone example

The table summarizes outcomes for two types of surgery (T = a and T = b) to remove kidney stones. Success rates are presented both overall and by stone size. For each treatment and size category, the table also reports the number of successes (R = 1) and the number of patients treated. For example, 192/263 means that 263 patients received treatment a and for 192 of them the surgery was successful (R = 1).

	Overall success rate	$Small\ stones$	Large stones
Treatment a Treatment b	78% (273/350)	93% (81/87)	73% (192/263)
	83% (289/350)	87% (234/270)	69% (55/80)

To decide which treatment to choose when the size of the kidney stone is unknown, we specify the following loss function:

$$\ell_0(R) = \begin{cases} c & \text{if } R = 0\\ 0 & \text{if } R = 1 \end{cases}$$
 (2)

where c > 0. The loss ignores recovery time and other considerations, and simply assigns a penalty c > 0 to an unsuccessful surgery.

(a) Determine the treatment T that minimises

$$\mathcal{R}(T) = \mathbb{E}_{p(R:do(T))} \left[\ell_o(R) \right]. \tag{3}$$

You may use that

$$p(R = 1; do(T) = a) = 0.833$$
 $p(R = 1; do(T) = b) = 0.779$ (4)

Solution. The loss equals $\mathcal{R}(T) = c \cdot p(R=0; do(T))$. Hence we choose the treatment for which the probability of failure is smallest, or, in other words, the treatment for which the probability p(R=1; do(T)) is largest. This is the case for T=a.

(b) Assume now that the success rates are considered uncertain (e.g because they are computed from a small number of patients only). How can this uncertainty be incorporated in the decision making?

Solution. We can consider the success rates in the table to be random variables \mathbf{h} . Their distribution is not affected by the action that we are taking now, hence $p(\mathbf{h}; do(T)) = p(\mathbf{h})$. As seen in the lecture, we can then express the risk as

$$\mathcal{R}(T) = \mathbb{E}_{p(\mathbf{h})} \mathbb{E}_{p(R|\mathbf{h};do(T))} \left[\ell_o(R) \right]. \tag{S.5}$$

$$= c\mathbb{E}_{p(\mathbf{h})} \left[p(R = 0 | \mathbf{h}; do(T)) \right]. \tag{S.6}$$

This has the form of an expected loss with $\ell(\mathbf{h}, T) = p(R = 0 | \mathbf{h}; do(T))$, the probability of failure when taking action T for a given \mathbf{h} . Hence, for each value of \mathbf{h} , we compute $p(R = 0 | \mathbf{h}; do(T))$ and then choose the treatment T for which the expected probability of failure is smallest.

Exercise 3. Classification with an asymmetric loss

We here derive the optimal policy for binary classification when the costs of false positives and false negatives are unequal. Denote by $h \in \{0,1\}$ the true class label and the predicted label by $a \in \{0,1\}$. We consider the following loss

_		
h	a	$\ell(h,a)$
0	0	0
1	0	c_{fn}
0	1	c_{fp}
1	1	0

where $c_{fn} > 0$ indicates the cost of a false negative and $c_{fp} > 0$ the cost of a false positive. We assume that we were given data \mathbf{x} and that can compute the posterior $p(h|\mathbf{x})$. Derive the policy $a^*(\mathbf{x})$ that minimises the posterior expected loss, i.e.

$$a^*(\mathbf{x}) = \underset{a}{\operatorname{argmin}} \mathbb{E}_{p(h|\mathbf{x})} \left[\ell(h, a) \right]$$
 (5)

Solution. We compute the risk $\mathbb{E}_{p(h|\mathbf{x})}[\ell(h,a)]$ for the two actions a=0 and a=1:

$$\mathbb{E}_{p(h|\mathbf{x})}\left[\ell(h, a=0)\right] = p(h=1|\mathbf{x})c_{fn} \qquad \mathbb{E}_{p(h|\mathbf{x})}\left[\ell(h, a=1)\right] = p(h=0|\mathbf{x})c_{fp} \tag{S.7}$$

The decision rule thus becomes $a^*(\mathbf{x}) = 1$ if $p(h = 0|\mathbf{x})c_{fp} < p(h = 1|\mathbf{x})c_{fn}$, and $a^*(\mathbf{x}) = 0$ if the inequality is reversed (if we have a tie, we may take a random decision). That is, we choose label a = 1 if

$$p(h=0|\mathbf{x})c_{fp} < p(h=1|\mathbf{x})c_{fn} \tag{S.8}$$

$$\frac{c_{fp}}{c_{fn}} < \frac{p(h=1|\mathbf{x})}{p(h=0|\mathbf{x})} \tag{S.9}$$

Denoting c_{fp}/c_{fn} by λ and using that $p(h=0|\mathbf{x})=1-p(h=1|\mathbf{x})$ we can manipulate the above equation as follows:

$$\lambda < \frac{p(h=1|\mathbf{x})}{1 - p(h=1|\mathbf{x})} \tag{S.10}$$

$$\lambda - \lambda p(h = 1|\mathbf{x}) < p(h = 1|\mathbf{x}) \tag{S.11}$$

$$\lambda < p(h = 1|\mathbf{x})(1 + \lambda) \tag{S.12}$$

$$\frac{\lambda}{1+\lambda} < p(h=1|\mathbf{x}) \tag{S.13}$$

We thus see that as the relative cost λ of a false positive increases, the evidence in favour of h = 1 must become stronger, i.e. the posterior $p(h = 1|\mathbf{x})$ must become closer to one for the optimal decision to be a = 1.

Exercise 4. Penalised squared loss

We consider the squared loss subject to a squared penalty for the action to deviate from zero,

$$\ell(h,a) = (h-a)^2 + \lambda a^2,\tag{6}$$

where $\lambda > 0$ is a weighting factor.

Derive the action that minimise the expected loss $\mathbb{E}_{f(h)}[\ell(h,a)]$.

Solution. Let $m = \mathbb{E}_{f(h)}[h]$. Since

$$\mathbb{E}_{f(h)} \left[(h - a)^2 \right] = \mathbb{E}_{f(h)} \left[(h - m + m - a)^2 \right]$$
 (S.14)

$$= \mathbb{E}_{f(h)} \left[(h-m)^2 + 2(h-m)(m-a) + \right.$$

$$+ (m-a)^2$$
 (S.15)

$$= \mathbb{E}_{f(h)} \left[(h-m)^2 \right] + 0 + (m-a)^2$$
 (S.16)

$$= \mathbb{V}(h) + (m-a)^2 \tag{S.17}$$

we have

$$\mathbb{E}_{f(h)} [\ell(h, a)] = \mathbb{V}(h) + (m - a)^2 + \lambda a^2$$
 (S.18)

The derivative with respect to a is

$$-2(m-a) + 2\lambda a = 2a(1+\lambda) - 2m.$$
 (S.19)

Setting the derivative to zero and solving for a gives

$$a = \frac{m}{1+\lambda}. (S.20)$$

Since the second derivative of the expected loss is positive, the above is a minimum. Hence the optimal action is

$$a^* = \frac{1}{1+\lambda} \mathbb{E}_{f(h)}[h].$$
 (S.21)

The role of the penalty term is to shrink the mean $\mathbb{E}_{f(h)}[h]$ to zero.