

Decision-Making & Modelling Under Uncertainty (DMU)

Oliver Maclaren (oliver.maclaren@auckland.ac.nz)

[10 lectures / tutorials]

- Decision-making under uncertainty [5/10]
 - ↳ Basic concepts
 - ↳ Risk, probability, utility
 - ↳ Statistical: extended setup
 - ↳ formulation & empirical risk approx.
 - ↳ minimax & Bayes
 - ↳ Tutorial sheet
- Modelling under uncertainty {
 - models of risk & intervention}[5/10]
 - ↳ probability, graphical models, & independence
 - ↳ causal interpretations of graphical models
 - ↳ stochastic process models (esp. Markov)
 - ↳ simulation & estimation tools
 - ↳ Tutorial sheet

Lecture 4 : statistical decision theory
after Wald

- statistical decision functions
 - minimax & Bayes
-

So far we have looked at

- no data minimax
- lots of data (empirical loss)

Now we want to look at

- 'some data' case &
 - Bayesian approach & relation
to minimax
-

Setting so far

Recall that we reduced the no data problem down to eg

	θ_1	θ_2	θ_3
d_1	$\lambda(d_1, \theta_1)$	---	-
d_2	!	:	:
d_3	$\lambda(d_3, \theta_1)$...	$\lambda(d_3, \theta_3)$

etc by considering the expected loss under each distribution

indexed by θ :

$$\boxed{\lambda(d_i, \theta_j) = \mathbb{E}_{d_i, \theta_j} [\lambda(d_i, s)]}$$

note: updated
 in L3

Then we considered eg no data minimax with & without randomised (mixed) strategies - - -

statistical decision theory (wald)

The Wald model begins from tables of this form, ie:

	θ_1	θ_2	θ_3
d_1	$\lambda(d_1, \theta_1)$	--	-
d_2	!	:	:
d_3	$\lambda(d_3, \theta_1)$...	$\lambda(d_3, \theta_3)$

& typically (not always) assumes the losses are in 'regret form', ie

$$\boxed{\min_d \lambda(d, \theta) = 0 \text{ for all } \theta}$$

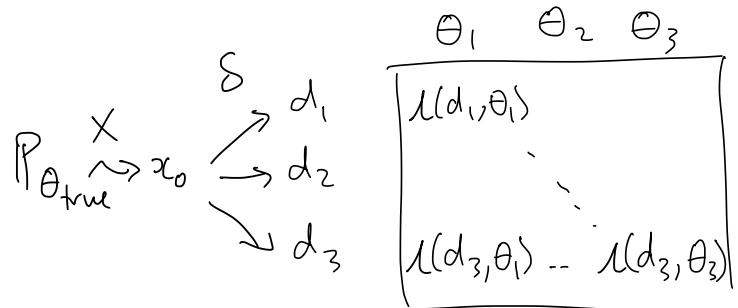
(eg $\lambda(d=\theta, \theta) = 0$)

→ The question is how to use observed data eg x_0 assumed to be a realisation of $X \sim P_{X|\theta}$

→ will use X & x here rather than s

Statistical decision functions: using data

Wald (~1950) had the idea to define 'statistical decision functions' $S(X)$ that map realisations of data to decisions:



$$\begin{aligned} \text{i.e.: } S: X &\mapsto \text{decision} \\ \text{i.e.: } \text{data} &\mapsto \text{decision.} \end{aligned}$$

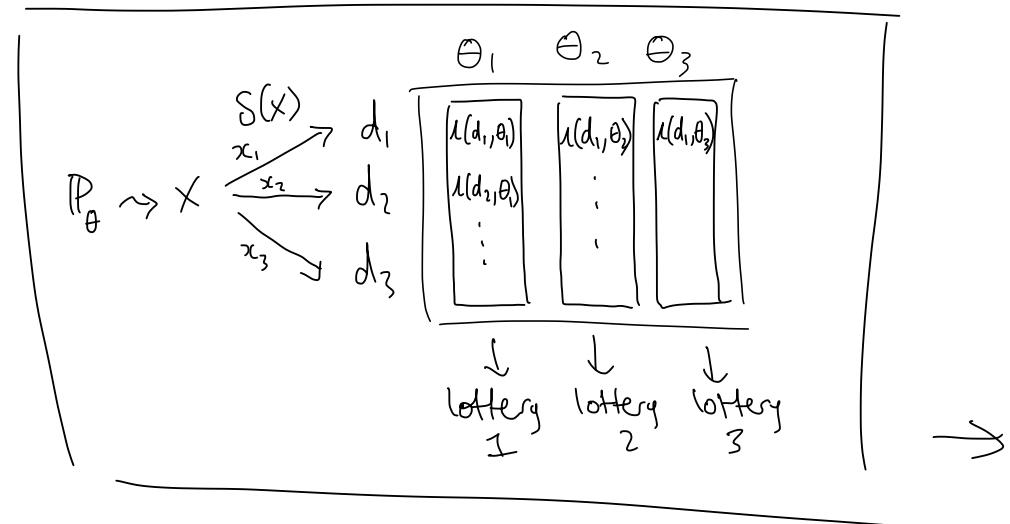
- These are also called 'decision rules' in the literature
- In statistics they are often called estimators



Statistical decision functions cont'd.

Treating the data x_0 as a realisation of a random variable X , these are like the randomised / mixed strategies considered previously, except now the random variable is related to the parameter:

→ We assume $X \sim P_{X; \theta}$, i.e. the data comes from a distribution indexed by θ . We again get lotteries:



From Barnett (1999),

'Comparative Statistical Inference'

270

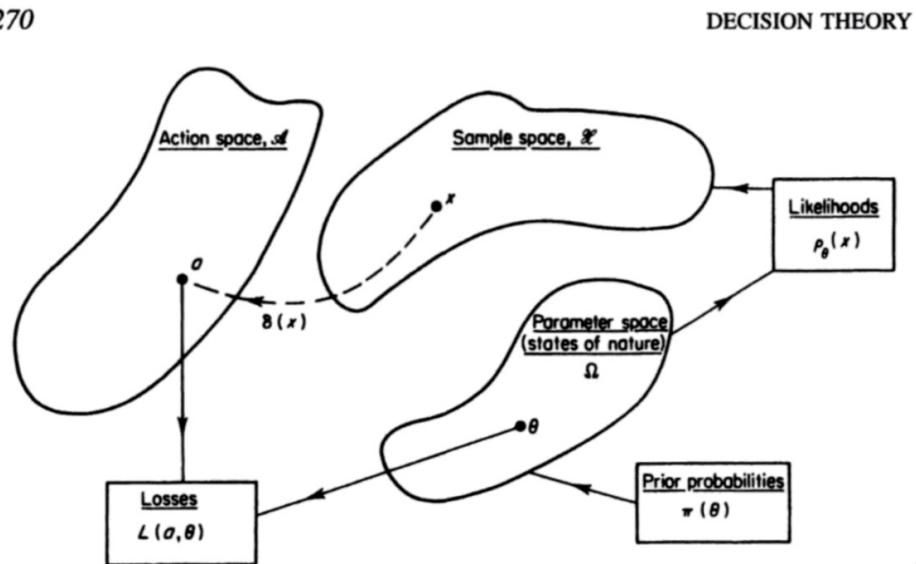


Figure 7.2.2 The superstructure of decision theory

Terminology:

- Note: in the statistical literature a decision rule / function is only called randomised if $\delta(x)$ maps to a further probability distribution over decisions (rather than to a single decision).
 - i.e. if additional randomness beyond the data distribution randomness is included.
- Even non-randomised statistical decision functions are still like 'mixed' or 'randomised' strategies in ('no data') simple decision theory / game theory as they include randomness from the data distribution
- → The question in stat. is whether additional randomness is useful (sometimes, yes!)

Risk function

Lotteries \Rightarrow can value any choice of decision function s (like how we valued coin flip randomised strategies):

$$\boxed{1(s, \theta) = E_{x; \theta} [l(\tilde{s}(x), \theta)]}$$

note:

decision function

random decision resulting from decision function applied to random variable

The above is often called the 'risk function' of a decision rule s in statistics: expected (regret) loss for decision function:

$$\boxed{R(s, \theta) := 1(s, \theta)}$$

Risk function?

\rightarrow This can be justified in terms of expected utility theory for one-off decisions

\rightarrow You can also think of it as providing a measure of the 'expected performance' of the rule $s(x)$ if you applied it over & over again for repeated realisations of X

\hookrightarrow '(frequentist) performance under repeated sampling'

Recall: expectation notation & calc.

→ given a distribution over a random variable X write:

$$\begin{aligned} E[x] &\equiv E_{P(x)}[x] \equiv E_x[x] \\ &= \begin{cases} \sum x \cdot p(x), & \text{discrete} \\ \int x \cdot p(x) dx, & \text{continuous} \end{cases} \end{aligned}$$

where $p(x) = \begin{cases} \text{o prob. mass function (pmf)} \\ \text{on } x \text{ (discrete)} \\ \text{o prob. density function (pdf)} \\ \text{on } x \text{ (continuous)} \end{cases}$
↳ ' $p(x)dx$ ' is pmf for 'dx' at x .

(more later)

→ indexed by θ : $P_{x|\theta} \equiv P(x; \theta)$

$$\begin{aligned} E[x] &\equiv E_{P(x; \theta)}[x] \equiv E_{x|\theta}[x] \\ &\text{etc.} \end{aligned}$$

Using risk functions

A risk function further transforms our table to the form:

	θ_1	θ_2	θ_3
s_1	$R(s_1, \theta_1)$	---	
s_2	:	⋮	
s_3			⋮

$(R(s_1, \theta_1) = 1(s_1, \theta_1))$
etc

(like reducing coin flip problem to choice of P)

In general, the risk of one rule will not be uniformly better than another

Example →

Interlude: Properties of squared error loss

suppose $\mathcal{L}(\delta(x), \theta) = (\delta(x) - \theta)^2$

define $\bar{\delta} = \mathbb{E}_{X;\theta}[\delta(X)]$ & $\delta = \delta(X)$

(Note θ not mean of $\delta(X)$ dist in general)

$$\begin{aligned}\Rightarrow R(\delta, \theta) &= \mathbb{E}_{X;\theta}[\mathcal{L}(\delta(x), \theta)] = \mathbb{E}_{X;\theta}[\mathcal{L}(\delta, \theta)] \\ &= \mathbb{E}_{X;\theta}[(\delta - \theta)^2] = \text{"mean squared error"} \\ &= \mathbb{E}_{X;\theta}[((\delta - \bar{\delta}) + (\bar{\delta} - \theta))^2] \\ &= \mathbb{E}_{X;\theta}[(\delta - \bar{\delta})^2 + 2(\delta - \bar{\delta})(\bar{\delta} - \theta) + (\bar{\delta} - \theta)^2]\end{aligned}$$

Now: $2(\delta - \bar{\delta})(\bar{\delta} - \theta)$
 $\underset{\text{const.}}{\sim} \underset{\text{const.}}{\sim} \underset{\text{const.}}{\sim} \underset{\text{constant}}{\text{constant}}$
 $\delta(x)$

$$\Rightarrow \mathbb{E}_{X;\theta}[2(\delta - \bar{\delta})(\bar{\delta} - \theta)]$$

$$= 2(\bar{\delta} - \theta) \mathbb{E}_{X;\theta}[\delta(x) - \bar{\delta}]$$

$$= 0 \text{ since } \mathbb{E}_{X;\theta}[\delta(x)] = \bar{\delta}$$

& $\mathbb{E}_{X;\theta}[(\bar{\delta} - \theta)^2] = (\bar{\delta} - \theta)^2 \Rightarrow$
 $\underset{\text{const.}}{\sim}$

Interlude cont'd

$$\begin{aligned}\Rightarrow R(\delta, \theta) &= \mathbb{E}_{X;\theta}[(\delta(x) - \bar{\delta})^2] + \mathbb{E}_{X;\theta}[(\bar{\delta} - \theta)^2] \\ &= \mathbb{E}_{X;\theta}[(\delta(x) - \bar{\delta})^2] + (\bar{\delta} - \theta)^2\end{aligned}$$

$$\Rightarrow R(\delta, \theta) = \underbrace{\text{Var}(\delta(x))}_{\text{risk from summary of } \delta(X) \text{ dist.}} + \underbrace{[\text{bias}(\delta(x))]^2}_{\text{risk from summary not matching } \theta}$$

Implies:

$$\text{"Bias - Variance trade-off"}$$

If bias ≈ 0 then minimise variance

But trading some bias increase for further variance reduction
can be useful.



Example (Wasserman 2004)

Suppose $X \sim N(\theta, 1)^*$, take one sample

& we use squared error loss

Consider $\delta_1(x) = x$

$$\delta_2(x) = 3$$

which is better?

Can show

$$R(\delta_1, \theta) = \text{var}(\delta_1) = 1 \quad (\text{unbiased})$$

$$R(\delta_2, \theta) = (\text{bias}(\delta_2))^2 = (3 - \theta)^2 \quad (\text{zero variance})$$



Wasserman (2004) 'All of Statistics':

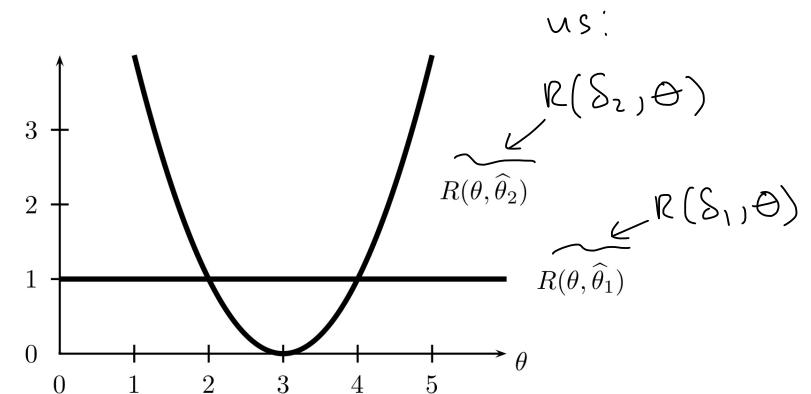


FIGURE 12.1. Comparing two risk functions. Neither risk function dominates the other at all values of θ .

'Silly' estimator δ_2 (here $\hat{\theta}_2$)

⇒ Can get lucky if $\theta = 3$ etc.

* X is normally distributed with mean = θ , variance = 1

Minimax again!

- One approach is to apply minimax again to choose a decision rule (decision function)

Example:

$$\max R(S_1, \theta) \leq \max R(S_2, \theta)$$

| |
| eg 4 if $\theta = 1$

\Rightarrow minimax says choose
 $S_1(x) = x$ over S_2

(S_1 is minimax wrt any other as well!)

Bayes

- Minimax is one way to reduce a risk function to a single number to compare rules
- Another way is to use a Bayesian approach.
 - Here the idea is to assume a distribution over parameters (states of nature), i.e. a distribution over distributions (prior), & consider the associated expected loss

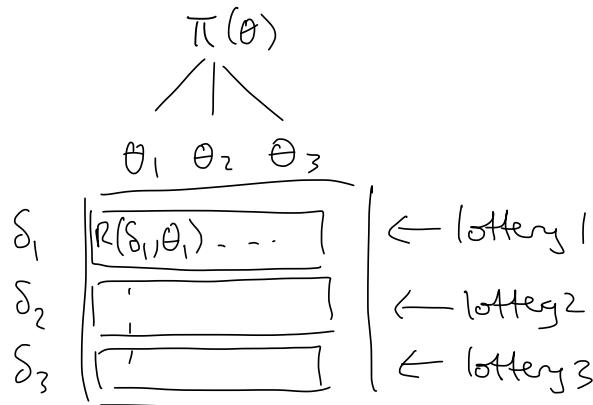
(Prior distribution: your belief about true parameter)

Bayes loss / Bayes risk

distribution over Θ

- Given a prior $\pi(\theta)$ over Θ we can compute the expected risk under that prior of any decision function
- Common to use the symbol $r(s, \theta)$ for Bayes risk & $\pi(\theta)$ for prior prob. distribution

Let:



(follow principles of expected utility given prior.

Gives Bayes risk:

$$\Rightarrow r(s, \pi) = \underbrace{\mathbb{E}_{\theta|\pi}[R(s, \theta)]}_{\text{expectation under } \pi(\theta)}$$

"Bayes risk" of s under $\pi(\theta)$

How to use?

Note:

$$r(s, \bar{\pi}) = \mathbb{E}_{\theta|\pi} \left[\mathbb{E}_{x|\theta} [l(s(x), \theta)] \right]$$

Averages over both data & parameter distributions

also $\{ \mathbb{E}_x; \theta = \mathbb{E}_{x|\theta} \text{ since } \theta \text{ has prob. dist}$

$$\{ \mathbb{E}_{\theta|\pi} \equiv \mathbb{E}_{\pi(\theta)} \equiv \text{expectation under } \pi(\theta) \}$$



Note: Conditional expectation notation:

" $E[X|Y]$ " another name for ...

$$E_X[X|Y] \equiv E_{X|Y}[X] \equiv E_{P(X|Y)}[X] \text{ etc}$$

$$= \begin{cases} \sum x \cdot P(x|y) \\ \int x \cdot P(x|y) dy \end{cases}$$

i.e. take expectation of X but
with $X \sim P(X|Y)$ rather
than $P(X)$.

Same for $f(x)$ in general.

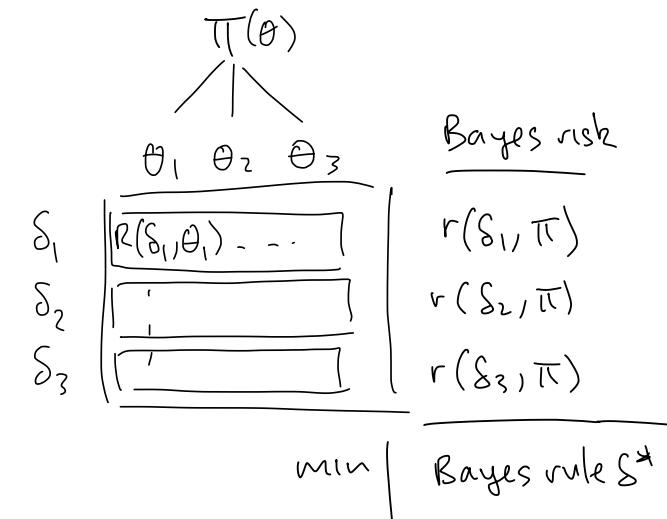
Should really write:

$$\left\{ E[X|Y=y] = E[X|y] \right. \quad \text{will also} \\ \text{for fixed } y. \quad \left. \right\} \text{use.}$$

Bayes decision rule

- Given the Bayes risk for a chosen prior & set of decision rules, we can then find the decision rule that minimises the Bayes risk for this prior

→ This is called the Bayes rule with respect to the prior:



Posterior expected loss

So far the expectation has been with respect to the (pre-data!) prior (& data dist).

→ we then choose the Bayes rule as the best rule under this prior & data distribution

It turns out that Bayes rules are equivalent to those that minimise the posterior loss:

$$r(s, \pi(\theta|x)) = \underbrace{E_{\pi(\theta|x)}[l(s, \theta)]}$$

expectation
under posterior

Bayes' theorem:

$$\pi(\theta|x) = \frac{\pi(x|\theta)\pi(\theta)}{\pi(x)}$$

(*proof: eg Wasserman (2004) & appendix →)
 $E_{\pi(\theta|x)} = (E_{\theta|x, \pi} \}$ diff. notation for same thing

Bayes rule: find posterior!

This is actually much simpler

→ for any given realisation x of X we find the posterior $\pi(\theta|x)$ use Bayes' rule (see later for more on computing)

→ we compute a posterior summary by taking the minimum expected value of $l(s, \theta)$ for given $\pi(\theta|x)$ & 'free choice' of s

→ this is our decision rule for any x

i.e.

$$s^*(x) = \min_s [E_{\pi(\theta|x)}[l(s, \theta)]]$$

$$E_{\pi(\theta|x)}[l(s, \theta)] = \int l(s(x), \theta) \pi(\theta|x) d\theta$$

Bayes summary:

- Given prior $\pi(\theta)$
- Given x (observed value)
- Find posterior $\pi(\theta|x)$
- Compute 'best summary' s^* of posterior $\pi(\theta|x)$ under loss $\lambda(s, \theta)$ & $\theta \sim \pi(\theta|x)$
- Defines $s(x) = s^*$ as Bayes rule value for given x
- Only need to consider x that occurs!

Example : squared error

$$\lambda(s, \theta) = (s - \theta)^2$$

$$\min_s \mathbb{E}_{\pi(\theta|x)} [(s - \theta)^2]$$

$$|\mathbb{E}_\theta|_{\pi, x} =$$

i.e. best summary of distribution over θ in terms of squared error

$$\Rightarrow \boxed{\text{Posterior mean!}} \quad (\text{same as L3})$$

$$s^* = \mathbb{E}_{\pi(\theta|x)} [\theta] \approx \frac{1}{n} \sum_i \theta_i \text{ for } \theta_i \text{ samples from posterior}$$

If $\lambda(s, \theta) = |s - \theta|$

\Rightarrow posterior median etc]

Note We will consider computing posteriors etc in more detail in 'modelling under uncertainty' part.

Bayes & minimax

A limitation of Bayes cf minimax is that the prior is fairly arbitrary ... & a different (personal) interpretation of probability is required!

Compromise:

There are a number of theorems connecting Bayes rules using a least favourable prior to minimax

→ use Bayes as computational tool to find minimax solutions (Wald!)

Beyond scope to go into detail but:

Bayes & minimax

Basic idea: instead of 'belief' think of prior as a 'randomised' or 'mixed' strategy of nature!

→ Nature's best choice = your worst
=
least favourable prior distribution

Eg Key results

(1) $\begin{cases} \text{If } S_0 \text{ is a Bayes rule wrt } \pi_0 \\ \& R(S_0, \theta) \leq r(S_0, \pi_0) \end{cases}$
then π_0 is a least favourable prior & S_0 is minimax

(2) $\begin{cases} \text{An 'equaliser' decision rule, ie} \\ |R(S_0, \theta)| = \text{constant} | \text{in } \theta, \text{ that} \\ \text{is also Bayes (or 'extended Bayes')} \\ \text{for some prior is minimax} \& \\ \text{the prior is least favourable} \end{cases}$

Further reading (selected!)

Wasserman (2004) All of statistics

Barnett (1999) Comparative statistical inference

Parmigiani & Inoue (2009) Decision Theory

Chernoff & Moses (1959) Elementary Decision Theory

Ferguson (1967) Mathematical statistics

Wald (1950) Statistical decision functions

+ most 'standard' books on mathematical statistics

Wasserman (2004) : good reference!
(see library)

12

Statistical Decision Theory

12.1 Preliminaries

We have considered several point estimators such as the maximum likelihood estimator, the method of moments estimator, and the posterior mean. In fact, there are many other ways to generate estimators. How do we choose among them? The answer is found in **decision theory** which is a formal theory for comparing statistical procedures.

Consider a parameter θ which lives in a parameter space Θ . Let $\hat{\theta}$ be an estimator of θ . In the language of decision theory, an estimator is sometimes called a **decision rule** and the possible values of the decision rule are called **actions**.

We shall measure the discrepancy between θ and $\hat{\theta}$ using a **loss function** $L(\theta, \hat{\theta})$. Formally, L maps $\Theta \times \Theta$ into \mathbb{R} . Here are some examples of loss functions:

$L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$	squared error loss,
$L(\theta, \hat{\theta}) = \theta - \hat{\theta} $	absolute error loss,
$L(\theta, \hat{\theta}) = \theta - \hat{\theta} ^p$	L_p loss,
$L(\theta, \hat{\theta}) = 0$ if $\theta = \hat{\theta}$ or 1 if $\theta \neq \hat{\theta}$	zero-one loss,
$L(\theta, \hat{\theta}) = \int \log \left(\frac{f(x; \theta)}{f(x; \hat{\theta})} \right) f(x; \theta) dx$	Kullback-Leibler loss.

Appendix : 'proof' of

Bayes decision rule \rightarrow min over (x, θ)

\rightarrow

min posterior expected loss = min over $(\theta|x)$

Write: $|E_\theta| = |E_{\pi(\theta)}|$ (expectation wrt $\pi(\theta)$)

Bayes risk:

$$\begin{aligned}
 r(\delta, \pi) &= |E_\theta| \left[|E_{x|\theta}| \left[\mathbb{1}(\delta(x), \theta) \right] \right] \\
 &= \int \left[\int \mathbb{1}(\delta(x), \theta) p(x|\theta) dx \right] P(\theta) d\theta \\
 &= \int \int \mathbb{1}(\delta(x), \theta) p(x|\theta) P(\theta) dx d\theta \\
 &= \int \int \mathbb{1}(\delta(x), \theta) P(\theta|x) P(x) dx d\theta \\
 &= \int \left[\int \mathbb{1}(\delta(x), \theta) P(\theta|x) dx \right] P(x) d\theta \\
 &= |E_x| \left[|E_{\theta|x}| \left[\mathbb{1}(\delta(x), \theta) \right] \right]
 \end{aligned}$$

\Rightarrow

Appendix cont'd.

$$\Rightarrow r(\delta, \pi) = |E_x| \left[|E_{\theta|x}| \left[\mathbb{1}(\delta(x), \theta) \right] \right]$$

$$\min_{\delta} r(\delta, \pi) = \min_{\delta} |E_x| \left[|E_{\theta|x}| \left[\mathbb{1}(\delta(x), \theta) \right] \right]$$

\Rightarrow minimise this
for each x
ie solve via

$$\Rightarrow \underset{\delta}{\text{minimise}} \quad |E_{\theta|x}| \left[\mathbb{1}(\delta, \theta) \right]$$

for each fixed value of x .

$$|E_{\theta|x}| \left[\mathbb{1}(\delta, \theta) \right] = \int \mathbb{1}(\delta, \theta) p(\theta|x) d\theta$$

= posterior expected loss.

Make best Bayes decision for each
posterior induced by each realisation
 x of X

Uses: Rule of iterated expectations

General form:

$$\begin{aligned} E_{(X,Y)}[f(X,Y)] &= E_Y \left[E_{X|Y}[f(X,Y)] \right] \\ &= E_X \left[E_{Y|X}[f(X,Y)] \right] \end{aligned}$$

Note:

if $f(X,Y) = Y$ then reduces to

$$\begin{aligned} E_X \left[E_{Y|X}[Y] \right] &= E_{(X,Y)}[Y] \\ &= E_Y[Y] \end{aligned}$$

i.e. $E_X \left[E_{Y|X}[Y] \right] = E_Y[Y]$

also written $E \left[E[Y|X] \right] = E[Y]$

which is the standard but (imo)
unnecessarily confusing form often used