# Decision Theory

Sofia Olhede

**EPFL**

October 13, 2020

## Loss functions IV

**EPFL**

- The risk is

$$
\begin{aligned}
R(\lambda, \widetilde{\lambda}) &= \mathbb{E}_\lambda \left[ \frac{n\lambda \bar{Y}}{n-1} - 1 - \log\left( \frac{n\lambda \bar{Y}}{n-1} \right) \right] \\
&= \mathbb{E}_\lambda \left[ \lambda \bar{Y} - 1 - \log(\lambda \bar{Y}) \right] + \frac{\mathbb{E}_\lambda \left[ \lambda \bar{Y} \right]}{n-1} - \log\left( \frac{n}{n-1} \right) \\
&= \mathbb{E}_\lambda \left[ \lambda \bar{Y} - 1 - \log(\lambda \bar{Y}) \right] + g(n). \qquad (1)
\end{aligned}
$$

- To derive the simplification we write $\bar{Y} = \frac{n-1}{n} \bar{Y} + \frac{1}{n} \bar{Y}$.
- Note that $\mathbb{E}_\lambda(\bar{Y}) = \lambda^{-1}$. Thus

$$
g(n) = \frac{1}{n-1} - \log\left( \frac{n}{n-1} \right).
$$

- We claim that $g(n) > 0$ once $n \geq 2$.

## Loss functions V

EPFL

- Using that $\log(x) = \int_1^x t^{-1} \, dt$ this follows if

$$\frac{1}{x} > \log(x+1) - \log(x), \quad x > 1$$

$$\Leftrightarrow \frac{1}{x} > \int_x^{x+1} t^{-1} dt, \quad x > 1. \tag{2}$$

This inequality holds by a rectangle area bound on the integral, as follows:

$$\frac{1}{x} = [(1+x) - x]\frac{1}{x} = \int_x^{x+1} \frac{1}{x} dt > \int_x^{x+1} \frac{1}{t} dt,$$

when $x > 1$.

- It therefore follows $R(\widetilde{\lambda}, \lambda) > R(\widehat{\lambda}, \lambda)$ and so $\underline{\widetilde{\lambda} \text{ dominates } \widehat{\lambda}}$.

# Decision Theory

**EPFL**

- We can push generality even further, and obtain an all encompassing framework.
- Called decision theory, it views inference as a game between nature and the statistician.
- Recall our general framework for statistical inference:

  *Reminder on how we do things*

  1. Model phenomenon by distribution $F(y_1, \ldots, y_n; \theta)$ for some $n \geq 1$.   └ *of how data was generated*
  2. Distributional form is known but $\theta \in \Theta$ is not known.
  3. Observe realisation of $(Y_1, \ldots, Y_n)$ from this distribution.
  4. Use $(Y_1, \ldots, Y_n)$ in order to make assertions concerning the true value of , and quantify the uncertainty associated with these assertions.
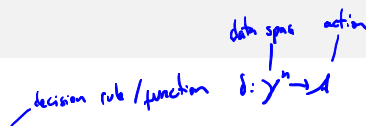
- The decision theory framework formalises step (4) to include estimation, testing, and confidence intervals.

# Decision Theory II

- In the decision theory framework we usually have these formal constructs:
  - * A family of distributions $\mathcal{F}$ usually assumed to admit densities (frequencies) ~ eg 0 a 1 in heads!
  - * A <u>parameter space</u> $\Theta$ that is used to parameterize $\mathcal{F} = \{F_\theta\}_{\theta \in \Theta}$. This models the possible realizations that may happen.
  - * The space on which observations are taken, the <u>data space</u> $\mathcal{Y}^n$.
  - * The action space which represents the space of possible actions or decisions or plays/moves available to the statistician.
  - * A <u>loss function</u> $\mathcal{L} : \Theta \times \mathcal{A} \to \mathbb{R}^+$.
  - * A <u>set $\mathcal{D}$</u> of <u>decision rules</u>. Any $\delta \in \mathcal{D}$ is a (measurable) function $\delta : \mathcal{Y}^n \mapsto \mathcal{A}$. This corresponds to possible strategies.

# Decision Theory III

*data space*    *action*

*decision rule / function*   $\delta : \mathcal{Y}^n \to \mathcal{A}$

- <u>The statistician would pick $\delta$ to limit losses</u>. <u>As the losses are random, actions are associated with risk</u>.

- Given a <u>decision rule</u> $\delta : \mathcal{Y}^n \mapsto \mathcal{A}$ the risk is:

- $R(\delta, \theta) = \mathbb{E}\{\mathcal{L}(\delta(Y, \theta)\}$. We compare decision rules in terms of their risk functions.   *the smaller the risk, the better!*

- Risk is not uniform, but depends on $\theta$ the true state of nature.

- So comparisons can be made   *picking a decision rule which is always better*

  1. <u>Uniform</u> comparisons (hard). Seek dominance everywhere in $\Theta$.
  2. <u>Minimax</u> (relaxed). Compare <u>worst-case</u> risks over $\Theta$.
  3. <u>Bayes</u> (relaxed). Compare <u>average</u> risk over $\Theta$.

# Risk

*is look at all worst performances and pick the "best worse"*

EPFL

- Rather than look at risk at every value of $\theta$ minimax risk concentrates (focuses) on the maximum risk:

- Definition (<u>Minimax decision rule</u>): Let $\mathcal{D}$ be a class of decision rules for an experiment $(\{f_\theta\}_{\theta \in \Theta}, \mathcal{L})$. If

$$\sup_{\theta \in \Theta} R(\delta, \theta) \leq \sup_{\theta \in \Theta} R(\delta', \theta),$$

then $\delta$ is a minimax rule.

- Rather than the behaviour at all $\theta$, Bayes risk is a weighted average.

- Definition (<u>Bayes Risk</u>): Let $\pi(\theta)$ be a probability density (or *prior on $\theta$ in the very-likely values of $\theta$*) frequency) and let $\delta$ be a decision rule for the experiment $(\{f_\theta\}_{\theta \in \Theta}, \mathcal{L})$. The $\pi$-Bayes risk is defined as

*"posterior risk"*

$$r(\delta, \pi) = \int_\Theta R(\delta, \theta) \pi(\theta) \, d\theta = \int_\Theta \int_{\mathcal{Y}} \mathcal{L}(\delta(\boldsymbol{y}), \theta) f_\theta(\boldsymbol{y}) \pi(\theta) \, d\boldsymbol{y} \, d\theta.$$

If $\delta \in \mathcal{D}$ is such that $r(\delta, \pi) \leq r(\delta', \pi)$ for all $\delta' \in \mathcal{D}$ then $\delta$ is a decision rule wrt $\pi$.

# Risk II

- The choice of prior $\pi(\theta)$ places different emphasis for different values of $\theta$ considering what prior knowledge we have.

- Minimax rules are useful to establish the fundamental inferential complexity or a statistical experiment.

- Using them for more practical purposes requires caution.

- Motivated as follows: we do not know anything about $\theta$ so let us insure ourselves against the worst thing that can happen.

- Minimax is quite a conservative point-of-view.

- Bayes rules are quite attractive as they <u>can nearly never be uniformly dominated</u>.

- <u>Intuitively, if you can show your rule to be Bayes for a nice prior, you know you are doing reasonably well</u>.

*[handwritten left margin: Minimax's / not so good]*

*[handwritten left margin: Bayes is better]*

# Hypothesis Testing

EPFL

- Model a phenomenon by a distribution $F(y_1, \ldots, y_n; \theta)$ on $\mathcal{Y}^n$ for some $n \geq 1$.
- Distributional form is known but $\theta \in \Theta$ is unknown.
- Observe $\begin{pmatrix} Y_1 & \ldots & Y_n \end{pmatrix}^T \in \mathcal{Y}^n$ from this distribution.
- Use the observed data $\begin{pmatrix} Y_1 & \ldots & Y_n \end{pmatrix}^T$ to make statements about $\theta$.
- The first assertion we wish to make is <u>hypothesis testing</u>. Given two disjoint regions $\Theta_0$ and $\Theta_1$, i.e. we have $\overline{\Theta_0 \cap \Theta_1 = \emptyset}$, <u>which interval is more likely to contain the true value of $\theta$</u>?

  *eg. in a trial: guilty or not guilty regions*

- We assume we know $\theta \in \Theta_0 \cup \Theta_1$.
- We need to use $\begin{pmatrix} Y_1 & \ldots & Y_n \end{pmatrix}^T \in \mathcal{Y}^n$ to decide between the two possibilities.

  *eg in hospital. sick vs healthy*

# Hypothesis Testing II

- Often in science two concurrent theories need to be confronted with the empirical evidence.

  The <u>null hypothesis</u> $H_0$ which states that $\theta \in \Theta_0$

  $$H_0 \; : \; \theta \in \Theta_0.$$

  The <u>alternative hypothesis</u> that postulates $\theta \in \Theta_1$

  $$H_1 \; : \; \theta \in \Theta_1.$$

Definition (Test Function): A test function is a map $\delta : \mathcal{Y}^n \mapsto \{0,1\}$.

- Obtaining 0 or 1 must be decided on whether or not the sample satisfies a certain condition:

  → critical region

  $$\delta(\mathsf{Y}) = \left\{ \begin{array}{lll} 1, & \text{if} & T(Y_1 \ldots Y_n) \in C \\ 0, & \text{if} & T(Y_1 \ldots Y_n) \notin C \end{array} \right. .$$

# Hypothesis Testing III

- $T$ is a statistic called a test statistic and;
- $C$ is a subset of the range of $T$ and is called the critical region.
- We can write

$$\delta(Y) = I(T(Y_1 \ldots Y_n) \in C).$$

- To choose good test functions we need to quantify the performance of a test function.
- Remark that, obviously, $\delta$ has a Bernoulli distribution:

$$\delta = \left\{ \begin{array}{ll} 1, & \text{if} \quad T(Y_1 \ldots Y_n) \in C \\ 0, & \text{if} \quad T(Y_1 \ldots Y_n) \notin C \end{array} \right. .$$

- So a good test function must have a sampling distribution concentrated around the right decision.
- The difference from point estimation is that our action space is discrete.
- But how do we compare $\delta$?

# Hypothesis Testing IV

- What possible errors are there?
- Take action 0 when $H_1$ is true–this is a type II error. Take action 1 when $H_0$ is true–this is a type I error. ?
- If we abused terminology we could define

$$\mathrm{MSE}\{\delta, H_i\} = \mathbb{E}_\theta(\delta - i)^2.$$

but result ⤷ the index of the right hypothesis

- We can then deduce that

$$\mathrm{MSE}\{\delta, H_i\} = \begin{cases} \mathbb{E}_\theta(\delta) & \text{if } \theta \in \Theta_0 \\ \mathbb{E}_\theta(1-\delta) & \text{if } \theta \in \Theta_1 \end{cases} \quad ?$$

$$= \begin{cases} \Pr_\theta(\delta = 1) & \text{if } \theta \in \Theta_0 \\ \Pr_\theta(\delta = 0) & \text{if } \theta \in \Theta_1 \end{cases}$$

$$= \begin{cases} \Pr_\theta(\delta = 1) & \text{if } \theta \in \Theta_0 \\ 1 - \Pr_\theta(\delta = 1) & \text{if } \theta \in \Theta_1 \end{cases}.$$

Asymmetry of false positive versus false negative.

# Hypothesis Testing V

EPFL

- In decision theory terms, the action space is $\mathcal{A} = \{0, 1\}$ and the loss function is the so-called "0–1" loss,
- The loss is then

$$\mathcal{L}(a, \theta) = \begin{cases} 1 & \text{if} & \theta \in \Theta_0 \;\&\; a = 1 \;(\text{Type } I \text{ error}) \\ 1 & \text{if} & \theta \in \Theta_1 \;\&\; a = 0 \;(\text{Type } II \text{ error}) \\ 0 & \text{o/w} & (\text{no error}) \end{cases}.$$

*ie wrong guess*

- Thus, we lose 1 unit whenever committing a type I or type II error.
- The risk function is

$$\mathcal{R}(\delta, \theta) = \begin{cases} \Pr\{\delta = 1\} & \text{if} & \theta \in \Theta_0 \;(\text{Prob type } I \text{ error}) \\ \Pr\{\delta = 0\} & \text{if} & \theta \in \Theta_1 \;(\text{Prob type } II \text{ error}) \end{cases}.$$

*probability of $\delta$ being 0.*

- Thus

$$\mathcal{R}(\delta, \theta) = \Pr\{\delta = 1\}\mathrm{I}(\theta \in \Theta_0) + \Pr\{\delta = 0\}\mathrm{I}(\theta \in \Theta_1).$$

Can we simultaneously control both errors??? NO :(.

# Hypothesis Testing VI

EPFL

- Let us understand why.
- Let $\delta(Y_1 \ldots Y_n) = I(T(Y_1 \ldots Y_n) \in C)$.
- <u>Suppose we wish to reduce the type I error</u> probability $\Pr_\theta(\delta = 1)$ if $\theta \in \Theta_0$.

  *eg we take a lower threshold for being guilty in trial*
- To do this, <u>we would replace C by a subset $C^* \subset C$</u>. This gives
  $\delta_*(Y_1 \ldots Y_n) = I(T(Y_1 \ldots Y_n) \in C^*)$.

  *↳ eg instead of $C = 8$ crimes or more = $[8, \infty)$ we will have $C^* = 9$ crimes or more = $[9, \infty)$*
- Observe that for $\theta \in \Theta_0$ we have

$$\Pr_\theta\{\delta_* = 1\} = \Pr\{T \in C^*\} \leq \Pr\{T \in C\} = \Pr_\theta\{\delta = 1\}.$$

  *↳ checks if count is in the new region*
- Whilst for $\theta \in \Theta_1$ we have

$$\Pr_\theta\{\delta_* = 0\} = \Pr\{T \notin C^*\} \geq \Pr\{T \notin C\} = \Pr_\theta\{\delta = 0\}.$$

- By reducing the type I error probability we increased the type II error probability! *→ obvious but a good reminder*

# Hypothesis Testing VII

- We have to make a philosophical choice regarding the importance of different errors.
- In applications, one type of error (false positive or negative) is typically more severe.
- Say this is the type I error, and exploit the asymmetry: fix a tolerance ceiling for the probability of this error.
- Given this ceiling, consider only test functions that respect it, and <u>focus on minimising type II error (i.e. maximising power).</u>

# Hypothesis Testing VIII

EPFL

- The Neyman-Pearson Framework

- We declare that we only consider test functions $\delta : \mathcal{X} \mapsto \{0, 1\}$ such that

*observation space*

$$\delta \in \mathcal{D}(\Theta_0, \alpha) = \{\delta : \sup_{\theta \in \Theta_0} \Pr_\theta \{\delta = 1\} \le \alpha\}.$$

*ie the probability of the success hypothesis is kept limited (bounded)*

- i.e. rules for which prob of type I error is bounded above by $\alpha$.

- Jargon: we fix a significance level for our test. *here we need good data but also anomalies*

- Within this restricted class of rules, choose $\delta$ to minimize prob of type II error:

$$\Pr\{\delta(\boldsymbol{X}) = 0\} = 1 - \Pr\{\delta(\boldsymbol{X}) = 1\}.$$

- Equivalently, maximize the power

$\Rightarrow$ $\beta(\theta, \delta) = \Pr\{\delta(\boldsymbol{X}) = 1\} = \mathbb{E}\,\mathrm{I}\{\delta(\boldsymbol{X}) = 1\} = \mathbb{E}\{\delta(\boldsymbol{X}))\}, \quad \theta \in \Theta_1.$

# Hypothesis Testing IX

**EPFL**

- Neyman-Pearson setup naturally exploits any asymmetric structure.
- But, if natural asymmetry absent, need judicious choice of $H_0$.
- Consider simplest situation: $\Theta_0 = \{\theta_0\}$ and $\Theta_1 = \{\theta_1\}$.
- The <u>Neyman Pearson Lemma</u>: Let $\boldsymbol{Y}$ have joint density/frequency $f$ where $f \in \{f_0, f_1\}$. We wish to test

$$H_0: \quad f = f_0 \quad \text{and} \quad f = f_1.$$

$$\text{or?}$$

- If $\Lambda(Y) = f_1(Y)/f_0(Y)$ is a continuous random variable, then there exists a $k > 0$ such that $\Pr\{\Lambda(Y) \geq k | H_0\} = \alpha$ and the test whose test function is given by $\delta(Y) = I(\Lambda(Y) \geq k)$ is a <u>most powerful (MP) test</u> of $H_0$ versus $H_1$ at significance level $\alpha$.