# Statistical Inference
## Lecture Three

Simon Shaw

University of Bath

APTS, 13-16 December 2022

# Overview of Lecture Three

In this lecture we will conclude the discussion of statistical principles, and move on to consider decision theory.

- Recall that two Bayesian models with the same prior distribution, $\mathcal{E}_{B,1} = \{\mathcal{X}_1, \Theta, f_{X_1}(x_1 \,|\, \theta), \pi(\theta)\}$ and $\mathcal{E}_{B,2} = \{\mathcal{X}_2, \Theta, f_{X_2}(x_2 \,|\, \theta), \pi(\theta)\}$ have the same posterior distribution when $f_{X_1}(x_1 \,|\, \theta) = c(x_1, x_2) f_{X_2}(x_2 \,|\, \theta)$. Hence, the Bayesian approach satisfies the SLP.

- Many classical procedures violate the SLP as they depend on values of the sample space $\mathcal{X}$ other than the observed value $x$.

# Overview of Lecture Three continued

- Bayesian statistical decision problem, $[\Theta, \mathcal{D}, \pi(\theta), L(\theta, d)]$.
- The risk of decision $d \in \mathcal{D}$ under the distribution $\pi(\theta)$ is $\rho(\pi(\theta), d) = \int_\theta L(\theta, d)\pi(\theta)\, d\theta$.
- The Bayes risk $\rho^*(\pi)$ minimises the expected loss,

$$\rho^*(\pi) \;\; = \;\; \inf_{d \in \mathcal{D}} \rho(\pi, d)$$

   with respect to $\pi(\theta)$.
- A decision $d^* \in \mathcal{D}$ for which $\rho(\pi, d^*) = \rho^*(\pi)$ is a Bayes rule against $\pi(\theta)$.
- A decision rule $\delta(x)$ is a function from $\mathcal{X}$ into $\mathcal{D}$,
- We view the set of decision rules, to be our possible set of inferences about $\theta$ when the sample is observed so that $\text{Ev}(\mathcal{E}, x)$ is $\delta^*(x)$
- The Bayes rule for the posterior decision respects the strong likelihood principle.

# Binomial and Negative Binomial example

- Let $\mathcal{E}_1 = \{\mathcal{X}, \Theta, f_X(x \mid \theta)\}$, where $X \mid \theta \sim Bin(n, \theta)$ so that

$$f_X(x \mid \theta) \;=\; \binom{n}{x} \theta^x (1-\theta)^{n-x}, \;\; x = 0, 1, \ldots, n.$$

- Let $\mathcal{E}_2 = \{\mathcal{Y}, \Theta, f_Y(y \mid \theta)\}$, where $Y \mid \theta \sim Nbin(r, \theta)$, so that

$$f_Y(y \mid \theta) \;=\; \binom{y-1}{r-1} \theta^r (1-\theta)^{y-r}, \;\; y = r, r+1, \ldots.$$

- Suppose we observe $x = r = 3$ and $y = n = 12$ then

$$f_X(3 \mid \theta) \;=\; \binom{12}{3} \theta^3 (1-\theta)^9, \;\; f_Y(12 \mid \theta) = \binom{11}{2} \theta^3 (1-\theta)^9$$

- Thus, $f_X(3 \mid \theta) \propto f_Y(12 \mid \theta)$.

- Consider the hypothesis test $H_0 : \theta = \frac{1}{2}$ versus $H_1 : \theta < \frac{1}{2}$ at significance level 5%.
- Let $\text{Ev}(\mathcal{E}_1, 3)$ be the result of the hypothesis test for the Binomial model where small values of $X$ support $H_1$

$$\mathbb{P}(X \leq 3 \,|\, \theta = 1/2) \;=\; \sum_{x=0}^{3} f_X(x \,|\, \theta = 1/2) \;=\; 0.0730.$$

- Thus, $\text{Ev}(\mathcal{E}_1, 3)$ is to not reject $H_0$.
- Let $\text{Ev}(\mathcal{E}_2, 12)$ be the result of the hypothesis test for the Negative Binomial model where large values of $Y$ support $H_1$

$$\mathbb{P}(Y \geq 12 \,|\, \theta = 1/2) \;=\; \sum_{y=12}^{\infty} f_Y(y \,|\, \theta = 1/2) \;=\; 0.0327.$$

- Thus, $\text{Ev}(\mathcal{E}_2, 12)$ is to reject $H_0$.
- This inference method does not respect the SLP: the choice of the model is relevant to the inference.

- Suppose that $\text{Ev}(\mathcal{E}, x)$ depends on the value of $f_X(x' \,|\, \theta)$ for some $x' \neq x$. Then, typically, Ev does not respect the SLP.
- We could create an alternate experiment $\mathcal{E}_1 = \{\mathcal{X}, \Theta, f_1(x \,|\, \theta)\}$ where:
  - $f_1(x \,|\, \theta) = f_X(x \,|\, \theta)$ for the observed $x$.
  - $f_1(x \,|\, \theta) \neq f_X(x \,|\, \theta)$ for all $x \in \mathcal{X}$.
- In particular, that $f_1(x' \,|\, \theta) \neq f_X(x' \,|\, \theta)$.
  - Let $\tilde{x} \neq x, x'$ and set

$$
\begin{array}{rcl}
f_1(x' \,|\, \theta) & = & \alpha f_X(x' \,|\, \theta) + \beta f_X(\tilde{x} \,|\, \theta) \\
f_1(\tilde{x} \,|\, \theta) & = & (1 - \alpha) f_X(x' \,|\, \theta) + (1 - \beta) f_X(\tilde{x} \,|\, \theta)
\end{array}
$$

  - By suitable choice of $\alpha$, $\beta$ we can redistribute the mass to ensure $f_1(x' \,|\, \theta) \neq f_X(x' \,|\, \theta)$. We then let $f_1 = f_X$ elsewhere.
- Consequently, whilst $f_1(x \,|\, \theta) = f_X(x \,|\, \theta)$ we will not have that $\text{Ev}(\mathcal{E}, x) = \text{Ev}(\mathcal{E}_1, x)$ and so will violate the SLP.

The two main difficulties with violating the SLP are:

1. To reject the SLP is to reject at least one of the WIP and the WCP. Yet both of these principles seem self-evident. Therefore violating the SLP is either illogical or obtuse.

2. In their everyday practice, statisticians use the SRP (ignoring the intentions of the experimenter) which is not self-evident, but is implied by the SLP. If the SLP is violated, it needs an alternative justification which has not yet been forthcoming.

# Reflections

- This chapter does not explain how to choose Ev but instead describes desirable properties of Ev.

- What is evaluated is the algorithm, the method by which $(\mathcal{E}, x)$ is turned into an inference about the parameter $\theta$.

- It is quite possible that statisticians of quite different persuasions will produce effectively identical inferences from different algorithms.

- A Bayesian statistician might produce a 95% High Density Region, and a classical statistician a 95% confidence set, but they might be effectively the same set.

- Primary concern for the auditor is why the particular inference method was chosen and they might also ask if the statistician is worried about the SLP.

- Classical statistician might argue a long-run frequency property but the client might wonder about their interval.

# Introduction

- Statistical Decision Theory allows us to consider ways to construct the Ev function that reflects our needs, which will vary from application to application, and which assesses the consequences of making a good or bad inference.
- The set of possible inferences, or decisions, is termed the decision space, denoted $\mathcal{D}$.
- For each $d \in \mathcal{D}$, we want a way to assess the consequence of how good or bad the choice of decision $d$ was under the event $\theta$.

## Definition (Loss function)

A loss function is any function $L$ from $\Theta \times \mathcal{D}$ to $[0, \infty)$.

- The loss function measures the penalty or error, $L(\theta, d)$ of the decision $d$ when the parameter takes the value $\theta$.
- Thus, larger values indicate worse consequences.

The three main types of inference about $\theta$ are

1. point estimation,
2. set estimation,
3. hypothesis testing.

It is a great conceptual and practical simplification that Statistical Decision Theory distinguishes between these three types simply according to their decision spaces.

| Type of inference | Decision space $\mathcal{D}$ |
| --- | --- |
| Point estimation | The parameter space, $\Theta$. |
| Set estimation | A set of subsets of $\Theta$. |
| Hypothesis testing | A specified partition of $\Theta$, denoted $\mathcal{H}$. |

# Bayesian statistical decision theory

In a Bayesian approach, a statistical decision problem $[\Theta, \mathcal{D}, \pi(\theta), L(\theta, d)]$ has the following ingredients.

1. The possible values of the parameter: $\Theta$, the parameter space.
2. The set of possible decisions: $\mathcal{D}$, the decision space.
3. The probability distribution on $\Theta$, $\pi(\theta)$. For example,
   1. this could be a prior distribution, $\pi(\theta) = f(\theta)$.
   2. this could be a posterior distribution, $\pi(\theta) = f(\theta \mid x)$ following the receipt of some data $x$.
   3. this could be a posterior distribution $\pi(\theta) = f(\theta \mid x, y)$ following the receipt of some data $x, y$.
4. The loss function $L(\theta, d)$.

In this setting, only $\theta$ is random and we can calculate the expected loss, or risk.

## Definition (Risk)

The risk of decision $d \in \mathcal{D}$ under the distribution $\pi(\theta)$ is

$$\rho(\pi(\theta), d) = \int_\theta L(\theta, d)\pi(\theta)\, d\theta.$$

We choose $d$ to minimise this risk.

## Definition (Bayes rule and Bayes risk)

The Bayes risk $\rho^*(\pi)$ minimises the expected loss,

$$\rho^*(\pi) = \inf_{d \in \mathcal{D}} \rho(\pi, d)$$

with respect to $\pi(\theta)$. A decision $d^* \in \mathcal{D}$ for which $\rho(\pi, d^*) = \rho^*(\pi)$ is a Bayes rule against $\pi(\theta)$.

The Bayes rule may not be unique, and in weird cases it might not exist. We solve $[\Theta, \mathcal{D}, \pi(\theta), L(\theta, d)]$ by finding $\rho^*(\pi)$ and (at least one) $d^*$.

## Example - quadratic loss

Suppose that $\Theta \subset \mathbb{R}$ and we wish to find a point estimate for $\theta$. We consider the loss function $L(\theta, d) = (\theta - d)^2$.

- The risk of decision $d$ is

$$
\begin{aligned}
\rho(\pi, d) = \mathbb{E}\{L(\theta, d) \,|\, \theta \sim \pi(\theta)\} &= \mathbb{E}_{(\pi)}\{(\theta - d)^2\} \\
&= \mathbb{E}_{(\pi)}(\theta^2) - 2d\mathbb{E}_{(\pi)}(\theta) + d^2,
\end{aligned}
$$

where $\mathbb{E}_{(\pi)}(\cdot)$ denotes the expectation with respect to $\pi(\theta)$.

- Differentiating with respect to $d$ we have

$$
\frac{\partial}{\partial d}\rho(\pi, d) = -2\mathbb{E}_{(\pi)}(\theta) + 2d.
$$

- So, the Bayes rule is $d^* = \mathbb{E}_{(\pi)}(\theta)$.

## Example - quadratic loss (continued)

- The corresponding Bayes risk is

$$
\begin{aligned}
\rho^*(\pi) \;=\; \rho(\pi, d^*) \;&=\; \mathbb{E}_{(\pi)}(\theta^2) - 2d^*\mathbb{E}_{(\pi)}(\theta) + (d^*)^2 \\
&=\; Var_{(\pi)}(\theta) + (d^* - \mathbb{E}_{(\pi)}(\theta))^2 \\
&=\; Var_{(\pi)}(\theta)
\end{aligned}
$$

where $Var_{(\pi)}(\theta)$ is the variance of $\theta$ computed with respect to $\pi(\theta)$.

1. If $\pi(\theta) = f(\theta)$, a prior for $\theta$, then the Bayes rule of an immediate decision is $d^* = \mathbb{E}(\theta)$ with corresponding Bayes risk $\rho^* = Var(\theta)$.

2. If we observe sample data $x$ then the Bayes rule given this sample information is $d^* = \mathbb{E}(\theta \,|\, X)$ with corresponding Bayes risk $\rho^* = Var(\theta \,|\, X)$ as $\pi(\theta) = f(\theta \,|\, x)$.

- Typically we solve:
  1. $[\Theta, \mathcal{D}, f(\theta), L(\theta, d)]$, the immediate decision problem,
  2. $[\Theta, \mathcal{D}, f(\theta \,|\, x), L(\theta, d)]$, the decision problem after sample information.

- We may also want to consider the risk of the sampling procedure, before observing the sample, to decide whether or not to sample.

- We now consider both $\theta$ and $X$ as random.

- For each possible sample, we need to specify which decision to make.

### Definition (Decision rule)

A decision rule $\delta(x)$ is a function from $\mathcal{X}$ into $\mathcal{D}$,

$$\delta : \mathcal{X} \to \mathcal{D}.$$

If $X = x$ is the observed value of the sample information then $\delta(x)$ is the decision that will be taken. The collection of all decision rules is denoted by $\Delta$ so that $\delta \in \Delta \Rightarrow \delta(x) \in \mathcal{D} \;\forall x \in X$.

- We wish to solve the problem $[\Theta, \Delta, f(\theta, x), L(\theta, \delta(x))]$.

> **Definition (Bayes (decision) rule and risk of the sampling procedure)**
>
> The decision rule $\delta^*$ is a Bayes (decision) rule exactly when
>
> $$\mathbb{E}\{L(\theta, \delta^*(X))\} \quad \leq \quad \mathbb{E}\{L(\theta, \delta(X))\}$$
>
> for all $\delta(x) \in \mathcal{D}$. The corresponding risk $\rho^* = \mathbb{E}\{L(\theta, \delta^*(X))\}$ is termed the risk of the sampling procedure.

- If the sample information consists of $X = (X_1, \ldots, X_n)$ then $\rho^*$ will be a function of $n$ and so can be used to help determine sample size choice.

## Bayes rule theorem, BRT

Suppose that a Bayes rule exists for $[\Theta, \mathcal{D}, f(\theta \,|\, x), L(\theta, d)]$. Then

$$\delta^*(x) = \underset{d \in \mathcal{D}}{\arg\min} \, \mathbb{E}(L(\theta, d) \,|\, X = x).$$

## Proof

Let $\delta$ be arbitrary. Then

$$
\begin{aligned}
\mathbb{E}\{L(\theta, \delta(X))\} &= \int_x \int_\theta L(\theta, \delta(x)) f(\theta, x) \, d\theta dx \\
&= \int_x \int_\theta L(\theta, \delta(x)) f(\theta \,|\, x) f(x) \, d\theta dx \\
&= \int_x \left\{ \int_\theta L(\theta, \delta(x)) f(\theta \,|\, x) \, d\theta \right\} f(x) \, dx \\
&= \int_x \mathbb{E}\{L(\theta, \delta(x)) \,|\, X\} f(x) \, dx
\end{aligned}
$$

## Proof continued

Now, as $f(x) > 0$, the $\delta^* \in \Delta$ which minimises $\mathbb{E}\{L(\theta, \delta(X))\}$ may equivalently be found as the $\delta^*$ which satisfies

$$\rho(f(\theta), \delta^*) \quad = \quad \inf_{\delta(x) \in \mathcal{D}} \mathbb{E}\{L(\theta, \delta(x)) \,|\, X\},$$

giving the result. $\qquad\square$

- The minimisation of expected loss over the space of all functions from $\mathcal{X}$ to $\mathcal{D}$ can be achieved by the pointwise minimisation over $\mathcal{D}$ of the expected loss conditional on $X = x$.
- The risk of the sampling procedure is $\rho^* = \mathbb{E}[\mathbb{E}\{L(\theta, \delta^*(x)) \,|\, X\}]$.

## Example - quadratic loss

We have $\delta^* = \mathbb{E}(\theta \,|\, X)$ and $\rho^* = \mathbb{E}\{Var(\theta \,|\, X)\}$.

We could consider $\Delta$, the set of decision rules, to be our possible set of inferences about $\theta$ when the sample is observed so that $\mathrm{Ev}(\mathcal{E}, x)$ is $\delta^*(x)$. We thus have the following result.

## Theorem

The Bayes rule for the posterior decision respects the strong likelihood principle.

## Proof

If we have two Bayesian models with the same prior distribution then if $f_{X_1}(x_1 \,|\, \theta) = c(x_1, x_2) f_{X_2}(x_2 \,|\, \theta)$ the corresponding posterior distributions are the same and so the corresponding Bayes rule (and risk) is the same. $\quad\square$

# Admissible rules

- Bayes rules rely upon a prior distribution for $\theta$: the risk is a function of $d$ only.

- In classical statistics, there is no distribution for $\theta$ and so another approach is needed.

### Definition (The classical risk)

For a decision rule $\delta(x)$, the classical risk for the model $\mathcal{E} = \{\mathcal{X}, \Theta, f_X(x \mid \theta)\}$ is

$$R(\theta, \delta) \;\; = \;\; \int_X L(\theta, \delta(x)) f_X(x \mid \theta) \, dx.$$

- The classical risk is thus, for each $\delta$, a function of $\theta$.

## Example

Let $X = (X_1, \ldots, X_n)$ where $X_i \sim N(\theta, \sigma^2)$ and $\sigma^2$ is known. Suppose that $L(\theta, d) = (\theta - d)^2$ and consider a conjugate prior $\theta \sim N(\mu_0, \sigma_0^2)$. Possible decision functions include:

1. $\delta_1(x) = \overline{x}$, the sample mean.

2. $\delta_2(x) = \text{med}\{x_1, \ldots, x_n\} = \tilde{x}$, the sample median.

3. $\delta_3(x) = \mu_0$, the prior mean.

4. $\delta_4(x) = \mu_n$, the posterior mean where

$$\mu_n = \left( \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1} \left( \frac{\mu_0}{\sigma_0^2} + \frac{n\overline{x}}{\sigma^2} \right),$$

the weighted average of the prior and sample mean accorded to their respective precisions.

## Example - continued

The respective classical risks are

1. $R(\theta, \delta_1) = \frac{\sigma^2}{n}$, a constant for $\theta$, since $\overline{X} \sim N(\theta, \sigma^2/n)$.

2. $R(\theta, \delta_2) = \frac{\pi\sigma^2}{2n}$, a constant for $\theta$, since $\tilde{X} \sim N(\theta, \pi\sigma^2/2n)$ (approximately).

3. $R(\theta, \delta_3) = (\theta - \mu_0)^2 = \sigma_0^2 \left(\frac{\theta - \mu_0}{\sigma_0}\right)^2$.

4. $R(\theta, \delta_4) = \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)^{-2} \left\{\frac{1}{\sigma_0^2}\left(\frac{\theta - \mu_0}{\sigma_0}\right)^2 + \frac{n}{\sigma^2}\right\}$.

Which decision do we choose? We observe that $R(\theta, \delta_1) < R(\theta, \delta_2)$ for all $\theta \in \Theta$ but other comparisons depend upon $\theta$.

- The accepted approach for classical statisticians is to narrow the set of possible decision rules by ruling out those that are obviously bad.

## Definition (Admissible decision rule)

A decision rule $\delta_0$ is inadmissible if there exists a decision rule $\delta_1$ which dominates it, that is

$$R(\theta, \delta_1) \ \leq \ R(\theta, \delta_0)$$

for all $\theta \in \Theta$ with $R(\theta, \delta_1) < R(\theta, \delta_0)$ for at least one value $\theta_0 \in \Theta$. If no such $\delta_1$ exists then $\delta_0$ is admissible.

- If $\delta_0$ is dominated by $\delta_1$ then the classical risk of $\delta_0$ is never smaller than that of $\delta_1$ and $\delta_1$ has a smaller risk for $\theta_0$.
- Thus, you would never want to use $\delta_0$.[1]
- The accepted approach is to reduce the set of possible decision rules under consideration by only using admissible rules.

---

[1]Here I am assuming that all other considerations are the same in the two cases: e.g. for all $x \in \mathcal{X}$, $\delta_1(x)$ and $\delta_0(x)$ take about the same amount of resource to compute.

- We now show that admissible rules can be related to a Bayes rule $\delta^*$ for a prior distribution $\pi(\theta)$.

### Theorem

If a prior distribution $\pi(\theta)$ is strictly positive for all $\Theta$ with finite Bayes risk and the classical risk, $R(\theta, \delta)$, is a continuous function of $\theta$ for all $\delta$, then the Bayes rule $\delta^*$ is admissible.

### Proof (Robert, 2007)

Letting $f(\theta, x) = f_X(x \mid \theta)\pi(\theta)$ we have

$$
\begin{aligned}
\mathbb{E}\{L(\theta, \delta(X))\} &= \int_x \int_\theta L(\theta, \delta(x)) f(\theta, x) \, d\theta dx \\
&= \int_\theta \left\{ \int_x L(\theta, \delta(x)) f_X(x \mid \theta) \, dx \right\} \pi(\theta) \, d\theta \\
&= \int_\theta R(\theta, \delta)\pi(\theta) \, d\theta
\end{aligned}
$$

## Proof continued

- Suppose that the Bayes rule $\delta^*$ is inadmissible and dominated by $\delta_1$.
- Thus, in an open set $C$ of $\theta$, $R(\theta, \delta_1) < R(\theta, \delta^*)$ with $R(\theta, \delta_1) \leq R(\theta, \delta^*)$ elsewhere.
- Consequently, $\mathbb{E}\{L(\theta, \delta_1(X))\} < \mathbb{E}\{L(\theta, \delta^*(X))\}$ which is a contradiction to $\delta^*$ being the Bayes rule. $\qquad\square$

- The relationship between a Bayes rule with prior $\pi(\theta)$ and an admissible decision rule is even stronger.
- The following result was derived by Abraham Wald (1902-1950)

## Wald's Complete Class Theorem, CCT

In the case where the parameter space $\Theta$ and sample space $\mathcal{X}$ are finite, a decision rule $\delta$ is admissible if and only if it is a Bayes rule for some prior distribution $\pi(\theta)$ with strictly positive values.

- An illuminating blackboard proof of this result can be found in Cox and Hinkley (1974, Section 11.6).

- There are generalisations of this theorem to non-finite decision sets, parameter spaces, and sample spaces but the results are highly technical.

- We'll proceed assuming the more general result, which is that a decision rule is admissible if and only if it is a Bayes rule for some prior distribution $\pi(\theta)$, which holds for practical purposes.

So what does the CCT say?

1. Admissible decision rules respect the SLP. This follows from the fact that admissible rules are Bayes rules which respect the SLP. This provides support for using admissible decision rules.

2. If you select a Bayes rule according to some positive prior distribution $\pi(\theta)$ then you cannot ever choose an inadmissible decision rule.