# Decision Theory

## ① The loss function

Given a parameter $\theta$ living in some space $\Theta$, we would like to construct estimators of $\theta$ based on some data.

There are various ways of constructing estimators and thus our main task is to compare them.

To do this we introduce the notion of the loss function which measures to some extent how far the estimator is from the parameter $\theta$.

Thus $L(\theta, \hat{\theta})$ is the loss function.

Ex 1) $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$ (probably the most important one)

2) $L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$

3) $L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|^p$, $p > 0$ the $L^p$-loss

4) $L(\theta, \hat{\theta}) = \begin{cases} 0 & \text{if } \theta = \hat{\theta} \\ 1 & \text{if } \theta \neq \hat{\theta} \end{cases}$ zero-one loss

5) $L(\theta, \hat{\theta}) = \int \log \frac{f(x; \theta)}{f(x; \hat{\theta})} f(x; \theta) \, dx$

which is called the Kullback-Leisler loss

Now assume we harvested some data $X_1, \ldots, X_n$ The estimator $\hat{\theta}$ is a function of $X_1, \ldots, X_n$.

**Definition** The risk of an estimator is

$$R(\theta, \hat{\theta}) = E_\theta [L(\theta, \hat{\theta})] = \int L(\theta, \hat{\theta}(x)) f(x; \theta) \, dx$$

where $f(x; \theta)$ is the density of the sample $X_1, \ldots, X_n$.

② Comparison of estimators and risks

In principle, we would like to have an estimator with the smallest possible risk. However the risk actually depends on the parameter $\theta$ range.
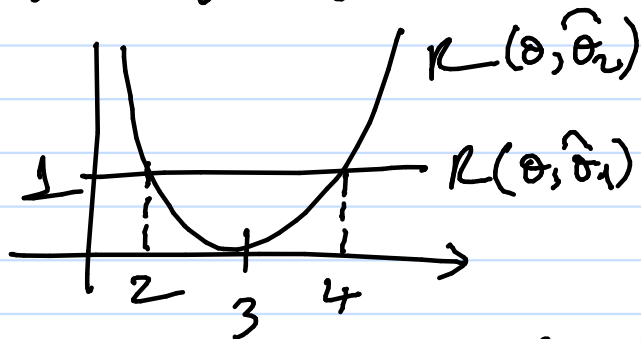
Ex $X \sim N(\theta, 1)$ and $\hat{\theta}_1 = X$, $\hat{\theta}_2 = 3$. Which of these is better in terms of quadratic loss?

$$L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$$

and $R(\theta, \hat{\theta}_1) = E[(X-\theta)^2] = Var(x) = 1$

$R(\theta, \hat{\theta}_2) = E[(3-\theta)^2] = (3-\theta)^2.$

Thus depending on the range of $\theta$ we have one or the other being smaller



Therefore for $\theta \in [2,4)$, $\hat{\theta}_2$ is a better estimator, but for $\theta \notin [2,4)$, $\hat{\theta}_1$ is better.

Thus the risk functions are not comparable everywhere.

Ex2 This is one important example as it
introduces the key concept of Bayes estimator

Take $X_1, \ldots, X_n \sim$ Bernoulli($p$). Consider the
square error loss and let $\hat{p} = \bar{X}$. This
has

$$R(p, \hat{p}) = E\left[(\hat{p} - p)^2\right] = E\left[(\bar{X} - p)^2\right] =$$

$$= \text{Var}(\bar{X}) = \frac{p(1-p)}{n}.$$

This $\hat{p}_1$ is the MLE and also the MVUE for $p$.
Another estimator is

$$\hat{p}_2 = \frac{\sum X_i + \alpha}{\alpha + \beta + n} \qquad \text{for } \alpha, \beta > 0$$

This is obtained via the Bayes' theory of
estimators taking the prior parameter $p$ to be
distributed uniformly in $[0,1]$. See the appendix
for more on this.

The idea here is that the risk of $\hat{p}_2$ is

$$E\left[(\hat{p}_2 - p)^2\right] = E\left[\left(\frac{Y + \alpha}{\alpha + \beta + n} - p\right)^2\right] =$$

$Y = \sum\limits_{i=1}^{n} X_i$

$$= \text{Var}\left(\frac{Y + \alpha}{\alpha + \beta + n}\right) + \left(E\left[\frac{Y + \alpha}{\alpha + \beta + n}\right] - p\right)^2$$

$$= \frac{1}{(\alpha + \beta + n)^2} \text{Var}(Y) + \frac{(np + \alpha - p(\alpha + \beta + n))^2}{(\alpha + \beta + n)^2}$$

$$= \frac{n p(1-p)}{(\alpha + \beta + n)^2} + \frac{((1-p)\alpha - p\beta)^2}{(\alpha + \beta + n)^2}$$

$$= \frac{np(1-p) + ((1-p)\alpha - p\beta)^2}{(\alpha + \beta + n)^2}$$

$$= \frac{np - np^2 + (\alpha - p(\alpha + \beta))^2}{(\alpha + \beta + n)^2}$$

If we choose this to be independent of $p$, then we need to make sure that

$$np - np^2 + (\alpha - p(\alpha+\beta))^2 =$$

$$= np - np^2 + \alpha^2 - 2\alpha(\alpha+\beta)p + p^2(\alpha+\beta)^2$$

$$= \alpha^2 + p(n - 2\alpha(\alpha+\beta)) + p^2((\alpha+\beta)^2 - n)$$

Thus we need to choose $\alpha + \beta = \sqrt{n}$ &
$$\alpha(\alpha+\beta) = \frac{n}{2}$$

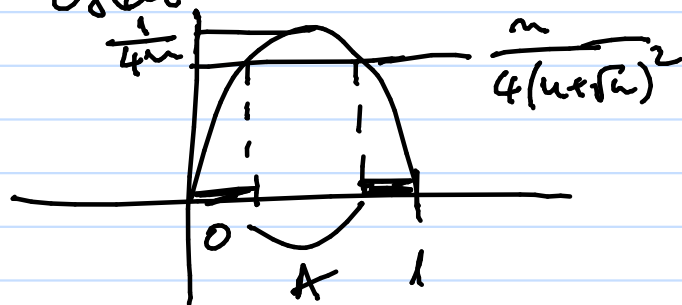which solves for $\alpha = \beta = \frac{\sqrt{n}}{2}$. This leads

to $\quad \widehat{p_2} = \dfrac{Y + \sqrt{n}/2}{n + \sqrt{n}}$ and the risc is

$$R(\widehat{p_2}, p) = \frac{\alpha^2}{(\alpha+\beta+n)^2} = \frac{n}{4(n+\sqrt{n})^2} < \frac{1}{4n}$$

which is independent of $p$.

Thus if we draw the dependence on $p$ we obtain



As we can see the first estimator is smaller on the set $A$

Notice that as $n \to \infty$, $\dfrac{n}{4(n+\sqrt{n})^2} \approx \dfrac{1}{4n}$.

One natural way of defining the risk is

$$\bar{R}(\hat{\theta}) = \sup_{\theta} R(\theta, \hat{\theta})$$

and the Bayes risk as

$$r(f, \hat{\theta}) = \int R(\theta, \hat{\theta}) f(\theta) d\theta$$

where $f(\theta)$ is a prior for $\theta$.

In the previous case we have

$$\bar{R}(\hat{p}_1) = \sup_{p} \frac{p(1-p)}{n} = \frac{1}{4n}$$

$$\bar{R}(\hat{p}_2) = \frac{n}{4(n+\sqrt{n})^2} < \frac{1}{4n}.$$

As we can see this means that $\hat{p}_2$ is better than $\hat{p}_1$ in general, however in absolute terms $\hat{p}_1$ is better only for $p$ outside a region which shrinks to 0 around ½.

Therefore $\hat{p}_1$ could be the preferred estimator though.

The Bayes' risk with the uniform prior $f$ for $\hat{p}_1$ and $\hat{p}_2$ are given by

$$r(f, \hat{p}_1) = \int_0^1 R(p, \hat{p}_1) dp = \int_0^1 \frac{p(1-p)}{n} dp = \frac{1}{6n}$$

$$r(f, \hat{p}_2) = \int_0^1 \frac{n}{4(n+\sqrt{n})^2} dp = \frac{n}{4(n+\sqrt{n})^2} > \frac{1}{6n}$$

for $n$ large enough, $(n \geq 20)$ Thus this suggests $\hat{p}_1$ is better than $\hat{p}_2$.

**Definition** A decision rule that minimizes
the Bayes risk is called a <u>Bayes rule/estimator</u>.
Formally $\hat{\theta}$ is a <u>Bayes rule</u> with respect to
$f$ if

$$r(f, \hat{\theta}) = \inf_{\tilde{\theta}} r(f, \tilde{\theta})$$

where the infimum is taken over all possible
estimators $\tilde{\theta}$.

    An estimator that minimizes the maximum
risk is called the <u>minimax</u> rule. Formally,
$\hat{\theta}$ is <u>minimax</u> if

$$\sup_{\theta} R(\theta, \hat{\theta}) = \inf_{\tilde{\theta}} \sup_{\theta} R(\theta, \tilde{\theta})$$

where the infimum is taken over all estimators
$\tilde{\theta}$.

<u>Bayes Estimators</u>

Assume $f$ to be a given fixed prior.
Then

$$f(\theta | x) = \frac{f(x|\theta) f(\theta)}{m(x)} = \frac{f(x|\theta) f(\theta)}{\int f(x|\theta) f(\theta) d\theta}$$

where $m(x)$ is the puf/pdf marginal of of $X$.

The posterior risk of an estimator $\hat{\theta}(x)$ is

$$r(\hat{\theta} | x) = \int L(\theta, \hat{\theta}(x)) f(\theta | x) d\theta$$

<u>Theorem</u> 1) $r(f, \hat{\theta}) = \int r(\hat{\theta}|x) m(x) dx$.
    2) If $\hat{\theta}(x)$ is the value of $\theta$ which minimizes
$r(\hat{\theta}|x)$, then $\hat{\theta}(x)$ is a Bayes estimator.

Pf 1) $r(f,\hat{\theta}) = \int R(\theta,\hat{\theta}) f(\theta) d\theta$

$\left|\begin{array}{l} f(x|\theta) f(\theta) = \\ = f(\theta|x) u(x) \end{array}\right.$

$= \int \left( \int L(\theta, \hat{\theta}(x)) f(x|\theta) dx \right) f(\theta) d\theta$

$= \int \int L(\theta, \hat{\theta}(x)) f(\theta|x) u(x) dx d\theta$

$= \int \left( \int L(\theta, \hat{\theta}(x)) f(\theta|x) d\theta \right) u(x) dx$

$= \int r(\hat{\theta}|x) u(x) dx$

2) If we take $\hat{\theta}(x)$ the minimizer of $r(\hat{\theta}|x)$

then $r(f,\tilde{\theta}) = \int r(\tilde{\theta}|x) u(x) dx \geq$

$\geq \int r(\hat{\theta}(x)|x) u(x) dx = r(f,\hat{\theta})$

and thus $\hat{\theta}$ is the Bayes estimator.

Theorem If $L(\theta,\hat{\theta}) = (\theta - \hat{\theta})^2$, then the
Bayes estimator is given by
$$\hat{\theta}(x) = \int \theta f(\theta|x) d\theta = \mathbb{E}[\theta|X=x].$$

The Bayes rule minimizes
$$r(\hat{\theta}|x) = \mathbb{E}_\theta \{ (\theta - \hat{\theta}(x))^2 | X=x \}$$
$$= \int (\theta - \hat{\theta}(x))^2 f(\theta|x) d\theta$$

and thus $\hat{\theta}(x) = \mathbb{E}[\theta|X=x]$.

Pf Indeed $r(\hat{\theta}|x) = \int (\theta - \hat{\theta})^2 f(\theta|x) d\theta$
and thus the value $\hat{\theta}(x)$ which minimizes this is
$\mathbb{E}[\theta|X=x]$,

**Appendix** How to compute the Bayes estimators.

Assume that we have $f(x|\theta)$ and we make the assumption that $\theta$ follows a certain distribution called <u>prior</u>.

Then the Bayes estimator of $\theta$ given the data $X$ is the <u>posterior distribution</u>

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

where $f(x) = \int f(x|\theta) f(\theta) d\theta$.

Typically, if we observe $x_1,...,x_n$, then

$$L_n(\theta) = f(x_1|\theta)---f(x_n|\theta).$$

Thus $f(\theta|x^n) = C_n L_n(\theta) f(\theta)$

for some constant $C_n$.

The mean of the posterior distribution is a good source of estimators for the parameter $\theta$.

Ex $X_1,...,X_n \sim$ Bernoulli $(p)$. If we take the uniform dist for $p$, then

$$f(p|x^n) \propto f(p) L_n(p) = p^S (1-p)^{n-S}$$

where $S = \sum_{i=1}^{n} x_i$. Thus the posterior dist of $p$ is Beta$(s+1, n-s+1)$ In general Beta$(\alpha,\beta)$ has density $\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}$

The mean of Beta$(\alpha, \beta)$ is

$$\int_0^1 \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x \cdot x^{\alpha-1}(1-x)^{\beta-1} dx =$$

$$= \int_0^1 \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha+1-1}(1-x)^{\beta-1} dx =$$

$$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot \frac{\Gamma(\alpha+1)\Gamma(\beta)}{\Gamma(\alpha+\beta+1)} =$$

$$= \frac{\alpha}{\alpha+\beta},$$

Thus we have $p|\hat{\alpha} \sim$ Beta$(s+1, m-s+1)$

The mean of this is $\bar{p} = \frac{s+1}{m-s+1+s+1} = \frac{s+1}{m+2}$.

Thus the estimator can be written as

$$\bar{p} = \frac{s+1}{m+2},$$

In general, if we take as prior Beta$(\alpha, \beta)$
then $\bar{p} = \frac{s+\alpha}{\alpha+\beta+m}$

We can write this as

$$\bar{p} = \frac{m}{\alpha+\beta+m} \cdot \frac{s}{m} + \frac{\alpha+\beta}{\alpha+\beta+m} p_0$$

where $p_0 = \frac{\alpha}{\alpha+\beta}$ is the prior mean.

EX. If $X_1, \ldots, X_m \sim N(\theta, \sigma^2)$ and we take the
prior to be $N(a, b^2)$, then $\theta | \bar{X} \sim N(\bar{\theta}, \tau^2)$
where $\bar{\theta} = w\bar{X} + (1-w) a$

$$w = \frac{1/s^2}{1/s^2 + 1/b^2}, \quad 1/\tau^2 = 1/s^2 + 1/b^2$$

$$s^2 = \frac{1}{m} \sum_{i=1}^{m}(X_i - \bar{X})^2.$$

Indeed $f(x|\theta) \propto e^{-\sum \frac{(x_i-\theta)^2}{2\sigma^2}}$

and then $f(x|\theta)f(\theta) \propto e^{-\sum_{i=1}^{m} \frac{(x_i-\theta)^2}{2\sigma^2}} e^{-\frac{(\theta-a)^2}{2b^2}}$

and this can be written as

$$e^{-\frac{\theta^2}{2}\left(\frac{m}{\sigma^2}+\frac{1}{b^2}\right) + \theta\left(\frac{m\bar{x}}{\sigma^2}+\frac{a}{b^2}\right)}$$

$$= e^{-\frac{\theta^2}{2\alpha^2} + \theta\beta} = e^{-\frac{\theta^2}{2\alpha^2} + 2\frac{\theta}{2\alpha^2}\cdot\alpha^2\beta - \frac{(\theta-\alpha^2\beta)^2}{\alpha^2}} \propto e^{-\frac{(\theta-\alpha^2\beta)^2}{2\alpha^2}}$$

with $\frac{1}{\alpha^2} = \frac{m}{\sigma^2}+\frac{1}{b^2}$ & $\beta = \frac{m\bar{x}}{\sigma^2}+\frac{a}{b^2}$

which means that $\theta|x \sim N(\alpha^2\beta, \alpha^2)$.

We can write the mean as

$$\bar{\theta} = w\bar{x} + (1-w)a$$

where $w = \frac{\alpha^2}{\sigma^2/m} = \frac{1/(\sigma^2/m)}{1/(\sigma^2/m)+1/b^2}$

which is a linear combination between $\bar{x}$, the MLE & $a$.