# Inferensteori

raazesh.sainudiin@math.uu.se : subject 1MS035

Kollas w/ pp hrs:
Theoretical foundation for datasceintists
(HT19)

## Sets, maps & Numbers

### Sets

A set is a collection of distinct elements.

e.g. $\{0, 0\}$. We give names to sets, e.g. $A = \{0, 0\}$.

$A' = \{0, 0, 0\}$ is not a set (multiset)

$\emptyset = \{\}$ is the empty set.

An element belongs to (or doesn't belong to) a set and we write $\in$ or $\notin$ respectively. e.g. $0 \in \{0, 0\}$ and $0 \notin \{a, 0\}$. etc. etc. det här kan du ..

$a \in \emptyset$

### Set operations

We can add elements to an existing set by union operation.

e.g. $\{\} \cup \{x\} = \{x\}$

$\{0, 0\} \cup \{x\} = \{0, 0, x\}$

$\{0, 0\} \cup \{0\} = \{0, 0\}$

**Def** $A \cup B := \{x \mid x \in A \text{ or } x \in B\}$

Intersection: $A \cap B := \{x \mid x \in A \text{ and } x \in B\}$

Set difference: $A \setminus B := \{x \mid x \in A \text{ and } x \notin B\}$

complement

Given a universal set $U$, $A^c := \{x \mid x \notin A\} = U \setminus A$

# Maps

A map or a function associates each element in the set called domain with exactly one element in the set range (codomain)

## Formally

A function is a specific kind of relation between elements in the domain and range

## Inverse Image/map/function

The inverse image of a function $f: X \to Y$ is $f^{(-1)}: Y \to X$ where

$f^{(-1)}(y) = \{x \mid x \in X \text{ and } f(x) = y\}$ or more generally, for any $B \subseteq Y$

$f^{(-1)}(B) = \{x \in X \mid f(x) \in B\}$ and $f^{(-1)} = \sigma_x(Y) \to \sigma_x(X)$

     collection of subsets of Y

## Note

In this course $N = \{1, 2, 3, ...\}$ and $Z_+ = Z_{\geq 0} = \{0, 1, 2, 3, ...\}$

# Probability

## Language

An experiment is an activity that produces distinct observable outcomes.
The set of such outcomes is called the sample space of the experiment,
denoted by $\Omega$.

An event is a subset of the sample space.

Probability is a function

$$P: \{events\} \to [0,1] \quad \text{where } P \text{ satisfies}$$

1) $\forall \text{ event } a, \ 0 \leq P(a) \leq 1$
2) $P(\Omega) = 1$
3) $A \cap B = \emptyset \Rightarrow P(A \cup B) = P(A) + P(B)$
4) $P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$ where $A_i$ are pairwise disjoint

## Motivation for axioms:

idea of long-term relative frequency of independent experiments (experiments/trials)
If we repeat an experiment a large number of times, the fraction of times
the event A occurs will be close to $P(A)$.

Formally, let $N(A,n)$ be the number of times A occurs in the first $n$ trials.

$$P(A) = \lim_{n \to \infty} \frac{N(A,n)}{n}$$

axiom: i) $0 \leq \frac{N(A,n)}{n} \leq 1$

ii) $\frac{N(\Omega,n)}{n} = \frac{n}{n} = 1$ ← something happens every time

iii) $A \cap B = \emptyset \Rightarrow N(A \cap B, n) = N(A,n) + N(B,n)$

iv) If this is ignored the mathematics you need is much harder

ex1 Tossing a fair coin. Can construct reals by dyadic partitioning.

$\Omega = \{HT\}$

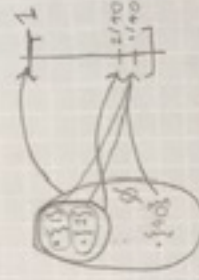$\begin{cases} \{H\} \\ \{T\} \\ \{HT\} \\ \emptyset \end{cases}$

Bernoulli random variable. Bernoulli($\theta$)

$2^\Omega = F_\Omega$

ex2/ NZ Lotto (40 balls)

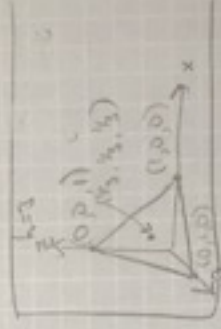Label of the 1st ball that comes out:

$\Omega = \{1, 2, \ldots, 40\}$

De Moivre random variable,

$2^\Omega = F_\Omega$

① DeMoivre$(\frac{1}{t}, \frac{1}{t}, \ldots, \frac{1}{t})$ = DeMoivre$(\theta_1, \theta_2, \ldots, \theta_L)$

Here we have $\theta_1 = \theta_2 = \ldots = \theta_L = \frac{1}{40}$, an Equi-probable DeMoivre variable

So, what is P("even number") = P($\{2, 4, \ldots, 40\}$) =

= P($\{2\} \cup \{4\} \cup \ldots \cup \{40\}$) = P($\{2\}$) + P($\{4\}$) + ... + P($\{40\}$) =

$= 20 \cdot \frac{1}{40} = \underline{\frac{1}{2}}$

(means repeatedly)

Properties:

1. $P(A) = 1 - P(A^c)$

2. A, B events. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

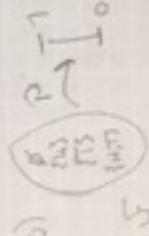The domain of $P$ is called a sigma field or sigma algebra,

It's denoted by $F(\Omega)$ or $F_\Omega$ or just $F$ if $\Omega$ is clear from context.

We see that

i) $\Omega \in F$,  ii) $A \in F \Rightarrow A^c \in F$,  iii) $A_1, A_2, \ldots \in F \Rightarrow \bigcup_{i=1}^{\infty} A_i \in F$

by Kolmogorov.

☞ Bernoulli experimental (coin flippr)

$$\Omega = \{1, 0\} \quad P \xrightarrow{\nearrow} 1$$
$$\mathcal{F} = \qquad \qquad 0$$
$$\{H,T\}$$

The triple $(\Omega, \mathcal{F}_\Omega, P)$ is called the probability space.
If $P$ is a set of probabilities, then $(\Omega, \mathcal{F}_\Omega, P)$ is called a statistical experiment.

Recall

events $A, B$ are independent $\Leftrightarrow P(A \cap B) = P(A) \cdot P(B)$

The product experiment

$X_1, \dots, X_n \overset{IID}{\sim} P_\theta$, then $P_\theta(\{X \in A, X_2 \in A, \dots, X_n \in A\}) = \prod_{i=1}^{n} P_\theta(X_i \in A)$
(indep. ident. dist.)

Conditional probability

$$Q(A|B) := \frac{P(A \cap B)}{P(B)} \quad \text{provided } P(B) > 0$$

is conditional probability a prob.? Yes basically $B$ is new $\Omega$.

☞ Construct my random graph from tossing unfair coins.

Toss a coin with $P(H) = \theta$ iid $n^2$ times.

Graph:

Let $V = \{v_1, \dots, v_n\}$ be a set of vertices and let $E \subseteq V^2$ and $|E| = k$

(n vertices, $k$ edges). Let

$$A = \begin{bmatrix} & v_1 & v_2 & \cdots & v_n \\ v_1 & & & & \\ \vdots & & & & \\ v_n & & & & \end{bmatrix}$$

adjacency matrix
where $A_{ij} = \begin{cases} 1, & \text{if } H \\ 0, & \text{if } T \end{cases}$ in coin toss

so, $E(|E|) = |V|^2 \theta$

# Prelude to Decision theory

## Estimation in parametric models



$\Omega$

$X$ ← Data ~ Random variable (R.V.)
probability prob./prob. (cont.)
function (disc.) density function (cont.)
distribution

$\downarrow$ $\mathcal{L}(X)$

## Problem

Let $X$ be a RV with values in $\mathcal{X}$ and $\mathcal{L}(X)$ is law of $x$.

Assume $\mathcal{L}(X)$ is known up to a finite dimensional parameter $\Theta$ from the parameter space $\Theta$ (hypothesis):

$$\mathcal{L}(X) \in \{P_{\theta} \mid \theta \in \Theta\}, \text{ here we assume } \Theta \subseteq \mathbb{R}^d \text{ for } d \geq 1.$$

The decision problem is to estimate a function $g(\theta)$ based on a realisation of $X$.

Typically (WLOG)
$$X = (x_1, \ldots, x_n), \quad x_i \overset{iid}{\sim} x_1$$

$n$ is called sample size, (initially) $\mathcal{X}$ is countable or $\subseteq \mathbb{R}^n$.

## Def ①
A statistic $T$ is an arbitrary function of the observed R.V. $X$ (data)

## Def ②
As an estimator $T$ of $g(\theta)$ we admit any $T: \mathcal{X} \to g(\Theta)$.
statistic

$$\{g(\theta) \mid \theta \in \Theta\}$$

Concretely:



$\omega \cdot T(x) = t$   gives us an estimate of $g(\theta)$

$$T^{-1} = g(\theta)$$

## Indicator function

$$X_{(A)}(\omega)\mathbb{1}^{(\omega)}_{(A)} := \begin{cases} 1 \text{ if } \omega \in A \\ 0 \text{ if } \omega \notin A \end{cases} \quad \text{is } \text{Bernoulli},$$

$$\text{Law}(X_i) \sim \text{Bernoulli}(\theta) \cdot \text{Bin}(1,\theta)$$

Suppose the data vector $x^\circ(x_1,..,x_n)$ is a realization of $X \sim \bigotimes_{i=1}^{n} \text{Bern}(\theta)$,

i.e. $x \in \mathbb{Z}^\circ \{0,1\}^n$, for unknown but fixed $\theta \in \boxed{\theta} \circ [0,1]$

Note that $\mathbb{P}_\theta(X=x) = \prod_{i=1}^{n} \mathbb{P}_\theta(X_i = x_i) = \prod_{i=1}^{n} \theta^{x_i}(1-\theta)^{1-x_i} = \theta^{\sum_{i=1}^{n}x_i}(1-\theta)^{n-\sum_{i=1}^{n}x_i} \overset{\textcircled{*}}{=} \theta^{t}(1-\theta)^{n-t}$.

$$\boxed{\text{Let } T(x) = \sum_{i=1}^{n} x_i}$$

Thus, the prob. of data (i.e. $\mathbb{P}_\theta(X=x)$) only depends on the statistic $t$.

Now consider another statistic : (sample mean)

$$T(X) = \frac{1}{n}\sum_{i=1}^{n} X_i = \overline{X}_n.$$

Then $\textcircled{*}$ becomes $\theta^t (1-\theta)^{n-t} \overset{t=n\overline{x}}{=} \theta^{n\overline{x}}(1-\theta)^{n(1-\overline{x})} = \theta^{n\overline{x}}(1-\theta)^{n(1-\overline{x})}$

An estimator of $\theta$, $g(\theta)$, (say $g(\theta)=\theta$), should converge towards the "true" but unknown $\theta$ to be estimated, as the sample size $n \to \infty$

## Def

A sequence $T_i := T_n \underbrace{(X_1, X_2, ..., X_n)}_{X^{(n)}}$ of estimators (each based on a sample of size $n$) for a parameter $\theta$ is called (asymptotically) $\underline{\text{consistent}}$ if

$$\forall \varepsilon > 0, \theta \in \boxed{\theta} : \mathbb{P}_{\theta,n}\left(|T_n(X^{(n)}) - \theta| \geq \varepsilon \right) \to 0 \quad \text{as } n \to \infty \quad \text{or, in shorter notation:}$$

$$T_n = T_n(X^{(n)}) \xrightarrow[n\to\infty]{\mathbb{P}} \theta \quad \text{if } \text{Law}(X^{(n)}) = \mathbb{P}_\theta$$

ex (Bern. product ex.)

The estimator $T_n(X_1,...,X_n) = \overline{X}_n$ is consistent for $\Theta$

## Proof

Since $X_1,...,X_n$ are iid $Be(\Theta)$ R.v.s, we have $E(X_1) = \Theta$ and the result follows from the law of large numbers.

So, consistency can be seen as a minimal requirement on estimators. But this still leaves a lot of consistent estimators to choose from.

A quantitative comparison of estimators is made possible by the approach of statistical decision theory.

We choose a loss function $Loss(t,\Theta)$ which measures the loss (inaccuracy) while unknown parameter $\Theta$ is estimated by $t$. $\underbrace{\qquad}_{\substack{\text{this est-or is a} \\ \text{realization of } T, \text{ i.e.} \\ \text{estimator } T(X^{(n)})}}$

Natural choices for loss when $\Theta \in \Theta \subseteq \mathbb{R}$

Absolute error: $Loss(t,\Theta) = |t - \Theta|$

Quadratic error: $Loss(t,\Theta) = (t-\Theta)^2$

or $\quad : Loss(t,\Theta) = \mathbb{1}_{[\delta,\infty)}(|t-\Theta|)$ for some $\delta > 0$ to emphasize the distance being less than $\delta$

## Note

Loss is a R.V. $= Loss(T(X),\Theta)$ needs to account for randomness.

## Def

The risk of an estimator $T$ at parameter $\Theta$ is

$R(t,\Theta) := E_\Theta(Loss(T(X),\Theta))$

$\uparrow$
Risk function of $T$

Note Risk might exist
since the expectation might not exist

## Idea

$$R(T^*, \theta) = \min_{T \in \mathcal{T}} R(T, \theta) \quad \text{for any fixed } \theta \in \Theta$$

where $\mathcal{T}$ is a class of estimators.

We find an estimator $T^*$ that minimizes the whole risk function simultaneously. If such a $T^*$ can be found, then it is called a uniformly best estimator (in $\mathcal{T}$).
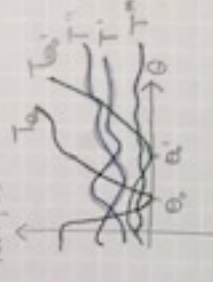
$$\underbrace{\text{(UBE)}}$$

In general, such a UBE won't exist

**Argument**

For each $\theta_0 \in \Theta$, consider $T_{\theta_0}(x) = \theta_0$

so $R(T_{\theta_0}, \theta_0) = 0 \Rightarrow$ best if $\theta_0$ is true



If $T^*$ was UBE, then it would have to compete with $T_{\theta_0}$

## Unbiased estimators

**Def**

Consider an estimator $T$ s.t. $E_\theta T$ exists.

Then $E_\theta(T) - g(\theta)$ is called bias of the estimator.

If $E_\theta(T) = g(\theta)$ for any $\theta \in \Theta$, then the the estimator $T$ is called unbiased for $g(\theta)$

## Def

A statistic $S$ is called sufficient for $\theta$ if $P_\theta(X \in B \mid S(X) = s)$ is independent of $\theta$, for all values of $s$ and all events $B$.

In words, conditional distribution of the data $X$ given that $S(X)$ takes any value does not depend on the parameter $\theta$



For discrete experiment

$$P_\theta(X \in B \mid S(x) = s) = \begin{cases} \dfrac{P_\theta(\{X \in B\} \cap \{S(x) = s\})}{P_\theta(S(x) = s)} & \text{, if } P_\theta(S(x) = s) > 0 \\[2mm] 0 & \text{, if } P_\theta(S(x) = s) = 0. \end{cases}$$

Let's clarify again:

Originally, the law of $X$ depended on $\Theta$ ( $\text{Law}(X) = P_\Theta(\cdot)$ ).

After the value of the sufficient statistic $S(X)$ is known, then the conditional law $P_\Theta(\cdot \mid S(X) = t)$ is no longer dependent on $\Theta$.

$$\underbrace{\phantom{xxxxxxxxxx}}_{g(\Theta)}$$

Since we are interested in making inferences about $\Theta$, the conditional law is uninteresting for our purpose, so we can "disregard it".

After taking $S(X)$ into account, the remaining randomness does not depend on $\Theta$ anymore.

Remark (exercise to prove formally):

The data itself is sufficient, i.e.

if $S(X) = X$, then for $x \approx \chi$, $\quad \Bbb{P}(X \in B \mid X = x) = \mathbb{I}_B(x) = \begin{cases} 1, & \text{if } x \in B \\ 0, & \text{if } x \notin B \end{cases}$

Prop

$1. \text{In } \overset{n}{\underset{i=1}{\bigotimes}} \text{Be}(\Theta)$ (⊗) e.g., the sample mean $\overline{X}_n$ is a sufficient statistic.

Proof

$n \overline{X}_n = \overset{n}{\underset{i=1}{\sum}} X_i \sim \text{Bin}(n, \Theta)$

Suppose $\overline{X}_n = t$ where $t$ is one of the possible values. This means $n\overline{X}_n = k$ for some $k \in \{0, 1, \ldots, n\}$. Then, for any $x = (x_1, \ldots, x_n)$,

$P_\Theta(X_1 = x_1, \ldots, X_n = x_n \mid n\overline{X}_n = k) = \dfrac{P_\Theta(X_1 = x_1, \ldots, X_n = x_n, n\overline{X}_n = k)}{P_\Theta(n\overline{X}_n = k)}$ (✗)

If $\Theta$ R.P.I.s $\overset{n}{\underset{i=1}{\sum}} x_i = k$, then (✗) is $\dfrac{\left(\Theta^k (1-\Theta)^{n-k}\right)}{\binom{n}{k}\Theta^k (1-\Theta)^{n-k}} = \dfrac{1}{\binom{n}{k}}$ which clearly is indep. of $\Theta$.

If for this $X$, $\sum_{i=1}^{n} x_i + h$, then the numerator is $\theta$ which is indep. of $\theta$, so $\theta$ is indep. of $\theta$. $\square$

Say we want to estimate $\theta$ in this example and we limit ourselves to estimators that are functions of the sufficient statistic $\overline{X}_n$:

$$T(x) = h(\overline{X}_n)$$

Additionally, suppose we limit ourselves further to unbiased estimators:

$$E(T(X)) - \theta = 0 \Rightarrow E(T(X)) = \theta$$

## Prop:

Under sufficiency and unbiasedness restrictions on allowed estimators of $\theta$ in our $\widehat{\theta} \; Be(\theta)$ exp., the only possible estimator is $\overline{X}_n$.

### Proof

$$0 = E_\theta h(\overline{X}_n) - \theta \overset{\downarrow}{=} E_\theta \left( h(\overline{X}_n) - \overline{X}_n \right) = \sum_{k=0}^{n} \underbrace{\left( h\left(\tfrac{k}{n}\right) - \tfrac{k}{n} \right)\binom{n}{k}}_{c_k} \theta^k (1-\theta)^{n-k} =$$

$$= (1-\theta)^n \sum_{k=0}^{n} c_k r^k \quad \text{for } r = \frac{\theta}{1-\theta}$$

Now, if $\theta \in [0,1)$, then $r \in [0,\infty)$, hence the above polynomial can only be zero if every $c_k$ is also zero.

Thus, $0 = \binom{n}{k}\left( h\left(\tfrac{k}{n}\right) - \tfrac{k}{n} \right) \overset{(I) \neq 0}{\Longrightarrow} h\left(\tfrac{k}{n}\right) = \tfrac{k}{n}$ for $k = 1, \ldots, n$ and $h(\overline{X}_n) = \overline{X}_n$ for all possible values of $\overline{X}_n$. $\square$

## Def

A statistic T is called **complete** if, for all $h: \mathbb{T} \to \mathbb{R}$:

$$\Big( E_\theta(h(T(X))) = 0 \text{ for all } \theta \in \Theta \Big) \Rightarrow$$

$$\Rightarrow \underbrace{P\big(h(T(X)) = 0\big) = 1 \text{ for all } \theta \in \Theta}_{\text{almost sure event}}$$

Intuitively, completeness means there is "no superfluous information" or "no redundancy" in the complete statistic T

Consider an event $T(X) \in B$ and suppose $P_\theta(T(X) \in B) = \alpha$ for $\alpha$ indep. of $\theta$.

By taking $h(t) = \mathbb{1}_B(t) - \alpha$, completeness $\Rightarrow P_\theta(\mathbb{1}_B(T(x)) = \alpha) = 1$ for all $\theta \in \Theta$ which means that $\alpha$ is either 0 or 1. Thus, for any event $T(X) \in B$ which has a non-trivial probability ($\neq 0, 1$), this prob. must depend on $\theta$.

## eg

$$\bigotimes_{i=1}^{n} Be(\theta)$$

The statistic $T(X) = (X_1, \bar{X}_n)$ is sufficient but not complete:

$h(X_1, \bar{X}_n) = X_1 - \bar{X}_n$ has $E(h) = 0$ but $h(X_1, \bar{X}_n)$ is not almost surely 0.

## Prop

$\bigotimes_{i=1}^{n} Be(\theta)$, $T(X) = \bar{X}_n$ is sufficient and complete.

## Proof

Suppose for some function $h$,

$E(h(\bar{X}_n)) = 0$

This means that $\sum_{k=0}^{n} h\left(\frac{k}{n}\right)\binom{n}{k}\theta^k (1-\theta)^{n-k} = 0 \quad \forall \theta \in [0,1]$.

Then, $h\left(\frac{k}{n}\right) = 0$ for $k \in \{0, 1, \ldots, n\}$ as per the earlier argument used to show that $\bar{X}_n$ is the only unbiased and sufficient estimator.

Limits of R.V.s : Chapter 8 in CSE Book.pdf 151-156

§1 - conv. of sequence of R.v.

$X_1, ..., X_n$ conv. to a R.v. X :

- In distribution

$X_n \xrightarrow{D} X$

- In probability

$X_n \xrightarrow{P} X$

- Markov's ineq., Chebyshev's ineq.

③ Suppose $T$ is a sufficient statistic with values in a set $\mathbb{T}$ and $S: \mathbb{T} \to \mathbb{S}$ is a one-to-one mapping with values in $\mathbb{S}$ (ie, there exists an inverse mapping $S^{[-1]}$ such that $S^{[-1]}(S(t)) = t$ for each $t \in \mathbb{T}$). Show that the statistic $S(T(x))$ is sufficient.

④ In class it was claimed that the statistic $T(x) = (x, \bar{X}_n)$ is sufficient for the $\bigotimes_{i=1}^{n}$ Bernoulli$(\theta)$ experiment. Prove this claim.

# Problem Set Week 3

① Show that the estimator $\bar{X}_n$ is consistent for $\theta$ in $\bigotimes_{i=1}^{n}$ Bernoulli($\theta$) experiment

② Suppose the data $X$ in a statistical experiment can take values in a countable set $\mathcal{X}$ (i.e, Law($X$) is discrete). In class it was claimed that the data itself are a sufficient Statistic (i.e, $T(X) = X$ is sufficient). Write down the argument that proves this claim (It can be a short paragraph)

⑤ Let $X_1, \ldots, X_n$ be independent and identically distributed with Law($X$) ~ Poisson($\lambda$), $\lambda \in (0, \infty)$ is unknown. Show that the sample mean $\bar{X}_n$ is a sufficient statistic.