

18.650
Statistics for Applications

Chapter 2: Parametric Inference

The rationale behind statistical modeling

- ▶ Let X_1, \dots, X_n be n independent copies of X .
- ▶ The goal of statistics is to learn the distribution of X .
- ▶ If $X \in \{0, 1\}$, easy! It's $\text{Ber}(p)$ and we only have to learn the parameter p of the Bernoulli distribution.
- ▶ Can be more complicated. For example, here is a (partial) dataset with number of siblings (including self) that were collected from college students a few years back: 2, 3, 2, 4, 1, 3, 1, 1, 1, 1, 1, 2, 2, 3, 2, 2, 2, 3, 2, 1, 3, 1, 2, 3, ...
- ▶ We could make no assumption and try to learn the pmf:

x	1	2	3	4	5	6	≥ 7
$\mathbb{P}(X = x)$	p_1	p_2	p_3	p_4	p_5	p_6	$\sum_{i \geq 7} p_i$

That's 7 parameters to learn.

- ▶ Or we could assume that $X \sim \text{Poiss}(\lambda)$. That's 1 parameter to learn!

Why we are doing Statistical modelling?

In practice, if we have data, if we observe a bunch of points, there are many ways we can think of this list of numbers. Like discrete, or continuous, or ≥ 1 , etc..

What statistical modelling is doing is to try to compress this information that could actually describe in a very naive way.

$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} P$ (Same distribution)

The goal of statistical modelling is to understand what is the distribution (P) we actually have here?

* Now we can start making Assumption's on this distribution P

Ex: P is Pdf, smooth (Assumption)

these Assumption's make our life simpler, Because by making successive Assumption's we are reducing the D.O.F of this space of distribution's

Ex: Assumption : if $x_i \in \{1, 0\} \Rightarrow x_i \sim \text{Ber}(p)$

↓
we just need
to find p

Ex: Assumption : if $x_i \in \mathbb{N} \Rightarrow x_i \sim \text{Poisson}(\lambda)$
we need to
find λ

Assumption: $x_i \in \{1, 2, 3, 4, 5, 6, >7\}$

\Rightarrow discrete distribution

$\Rightarrow x_i \sim \text{PMf}$

x	1	2	3	4	5	6	>7
$P(x=x)$	$P(x=1)$	$P(x=2)$	$P(x=3)$.	.	.	P_7
	p_1	p_2	p_3	p_4	p_5	p_6	p_7

$\underbrace{\qquad\qquad\qquad}_{\text{we are trying to estimate}} \text{these values}$

The ultimate goal of statistic is to say what distribution our data come from? Because that's basically the best we are going to be able to do.

Modelling 101: The purpose of modelling is to construct the space of possible distribution's to a subspace that's actually plausible, but much simple for us to estimate.

Statistical model (1)

Formal definition

Let the observed outcome of a statistical experiment be a sample X_1, \dots, X_n of n i.i.d. random variables in some measurable space E (usually $E \subseteq \mathbb{R}$) and denote by \mathbb{P} their common distribution. A *statistical model* associated to that statistical experiment is a pair

$$(E, (\mathbb{P}_\theta)_{\theta \in \Theta}),$$

where:

- ▶ E is *sample space*; *The set in which X lives*
- ▶ $(\mathbb{P}_\theta)_{\theta \in \Theta}$ is a family of probability measures on E ;
- ▶ Θ is any set, called *parameter set*. *The set in which θ lives.*

Statistical Experiment:

It's a pair $(E, \{P_\theta\}_{\theta \in \Theta})$
↑
Probability distribution

↓

$$\Theta = (\mu, \sigma^2) \Rightarrow N(\mu, \sigma^2) \text{ (2 dim)}$$

$$\Theta = \{P_1, P_2, \dots, P_6, P_{7,8}\}$$

(7 dim)

$\text{Pois}(\theta)_{\theta \in \mathbb{R}}$

What's important here is once they gave us
 Θ we know exactly all the probabilities
associated with random variable. we know its
distribution perfectly.

* The purpose of the statistical model is to
come we estimate the parameter, we actually
know exactly what distribution it has.

Statistical model (2)

- ▶ Usually, we will assume that the statistical model is *well specified*, i.e., defined such that $\mathbb{P} = \mathbb{P}_\theta$, for some $\theta \in \Theta$.
- ▶ This particular θ is called the true parameter, and is unknown: The aim of the statistical experiment is to *estimate* θ , or check its properties when they have a special meaning ($\theta > 2?$, $\theta \neq 1/2?$, ...)
- ▶ For now, we will always assume that $\Theta \subseteq \mathbb{R}^d$ for some $d \geq 1$: The model is called *parametric*.

Two things

- ① The purpose of the statistical model is once we estimate the parameter, we exactly know what distribution it has
- ② we could potentially have several parameters that give us the same distribution, (that would still be fine, because we could estimate one guy, or we could estimate the other guy, we would still recover underlying distribution of our data)

The problem is that this could create really annoying, theoretical problems, like things don't work, the algorithm's won't work. etc..

what we typically assume in the model is identified.

Voca blury:

Well Specified: for our observation's X

there exists $\exists \theta \in \Theta$ such that X follows $P_\theta \Rightarrow X \sim P_\theta$

(This is a strong Assumption) because

"All model's are wrong, But some of them are useful"

\Rightarrow mean's that may it's not true that this Poisson distribution we assumed for # siblings for college student is not perfectly correct.

\Rightarrow when we make this Assumption, we are actually assuming that the data really comes from Poisson distribution.

$$X \sim P_\theta \quad \theta^* \text{ or } \theta = \text{true parameter}$$

The Aim of this statistical experiment is to estimate θ , so that once we actually plug in θ in P_θ , we would know that the prob that my r.v takes any value.

Statistical model (3)

Examples

1. For n Bernoulli trials:

$$\left(\{0, 1\}, (\text{Ber}(p))_{p \in (0, 1)} \right).$$

2. If $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$, for some unknown $\lambda > 0$:

$$(\mathbb{R}_+^*, (\text{Exp}(\lambda))_{\lambda > 0}).$$

3. If $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poiss}(\lambda)$, for some unknown $\lambda > 0$:

$$(\mathbb{N}, (\text{Poiss}(\lambda))_{\lambda > 0}).$$

4. If $X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$, for some unknown $\mu \in \mathbb{R}$ and $\sigma^2 > 0$:

$$\left(\mathbb{R}, (\mathcal{N}(\mu, \sigma^2))_{(\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+^*} \right).$$

what we are trying to do as a statistician is to inject as much knowledge about the question and about the problem that we can so that the data has to do the minimal job and henceforth, we actually need less data.

Parametric Model (Parametric statistic)

$\theta \in \mathbb{R}^d$: where d Number of Parameter's

Non Parametric Model: infinite number of Parameter's to estimate.

Cannot be represented by vectors

For Ex: function

Model $P_f = \left\{ \begin{array}{l} \text{distribution of } f \quad f > 0 \\ f = 1 \end{array} \right\}$

Some Statistical models:

- ① Bernoulli trials $\Rightarrow (\{0, 1\}, \{\text{Ber}(p)\}_{p \in (0,1)})$
- ② Exponential $\Rightarrow ([0, +\infty), \{\text{Exp}(\lambda)\}_{\lambda > 0})$
- ③ Poisson $\Rightarrow (\mathbb{N}, \{\text{Pois}(\lambda)\}_{\lambda > 0})$
- ④ Siblin's model $\Rightarrow (\{1, 2, 3, \dots, 7\}, \{P(x=k) = p_k\}_{k=1 \text{ to } 7})$

⑤ Gaussian $(\mathbb{R}, \{N(\mu, \sigma^2)\}_{\substack{\mu \in \mathbb{R} \\ \sigma^2 \in \mathbb{R}^+}})$

$$\textcircled{5} = \{(u, \sigma^2) : u \in \mathbb{R}, \sigma^2 \in \mathbb{R}^+\}$$

Identification

$$\Theta \xrightarrow{\quad} \mathbb{P}_\theta \quad (\text{injective})$$

The parameter θ is called *identified* iff the map $\theta \in \Theta \mapsto \mathbb{P}_\theta$ is injective, i.e.,

$$\theta = \theta' \Rightarrow \mathbb{P}_\theta = \mathbb{P}_{\theta'}.$$

$$\theta \neq \theta' \implies \mathbb{P}_\theta \neq \mathbb{P}_{\theta'}$$

Examples

$$\nexists \quad \mathbb{P}_\theta = \mathbb{P}_{\theta'} \implies \theta = \theta'$$

1. In all four previous examples, the parameter was identified.
2. If $X_i = \mathbb{1}_{Y_i \geq 0}$, where $Y_1, \dots, Y_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$, for some unknown $\mu \in \mathbb{R}$ and $\sigma^2 > 0$, are unobserved: μ and σ^2 are not identified (but $\theta = \mu/\sigma$ is).

Parameter estimation (1)

Idea: Given an observed sample X_1, \dots, X_n and a statistical model $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$, one wants to *estimate* the parameter θ .

Definitions

- ▶ *Statistic*: Any measurable¹ function of the sample, e.g., $\bar{X}_n, \max_i X_i, X_1 + \log(1 + |X_n|)$, sample variance, etc...
- ▶ *Estimator* of θ : Any statistic whose expression does not depend on θ .
- ▶ An estimator $\hat{\theta}_n$ of θ is *weakly* (resp. *strongly*) *consistent* iff

$$\hat{\theta}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P} \text{ (resp. a.s.)}} \theta \quad (\text{w.r.t. } \mathbb{P}_\theta).$$

¹Rule of thumb: if you can compute it exactly once given data, it is measurable. You may have some issues with things that are implicitly defined such as sup or inf but not in this class

We need some measure of performance of a given parameter, we need to be able to evaluate if eyeballing the problem is worse than actually collecting a large amount of data. To able to answer these question's we need to answer what accuracy mean's

Estimator : (Statistic)

An Estimator is a measurable function of data

$$\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$$

An Estimator is an Statistic which does not depend on θ .

An Estimator is said to be consistent, if when we collect more and more data, $\hat{\theta}$ is getting closer and closer to true parameter.

A Estimator $\hat{\theta}(x_1, x_2, \dots, x_n)$ is called weakly convergent if it converge in probability; strongly convergent if it is almost surely

Parameter estimation (2)

- ▶ *Bias* of an estimator $\hat{\theta}_n$ of θ :

$$\mathbb{E} \left[\hat{\theta}_n \right] - \theta.$$

- ▶ *Risk* (or *quadratic risk*) of an estimator $\hat{\theta}_n$:

$$\mathbb{E} \left[|\hat{\theta}_n - \theta|^2 \right].$$

Remark: If $\Theta \subseteq \mathbb{R}$,

"Quadratic risk = bias² + variance".

BIAS of a Estimator

$$E[\hat{\theta}] - \theta$$

if $Bia = E[\hat{\theta}] - \theta = 0 \Rightarrow \hat{\theta}$ is unbiased

what is it mean to be unbiased?

- * May be for this particular round of data we collected (x_1, x_2, \dots, x_n) , The Estimate $\hat{\theta}(x_1, x_2, \dots, x_n)$ we are pretty far from the true Estimator, But if I re do this experiment over and over again, on average all the values of my Estimator I got, this would be my true Parameter.

- * If I were to repeat this experiment, in average we would get this thing right.

Ex: Estimator \bar{x}_n

$$\begin{aligned} E[\bar{x}_n] &= E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{n}{n} E[x] \\ &= E[x] \end{aligned}$$

$$\Rightarrow E[\bar{x}_n] = E[x] = p$$

Another Estimator : X_1

$$\hat{\Theta}(x_1, x_2, \dots, x_n) = X_1$$

$$E[X_1] = P \quad (\text{unbiased})$$

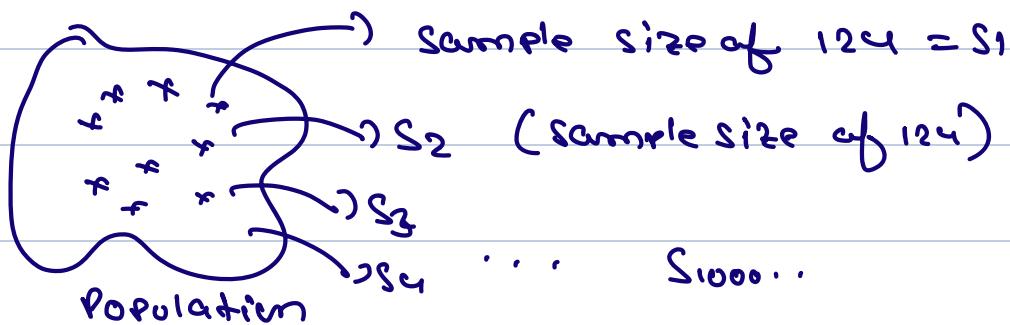
This should probably illustrate to us that Bias is not something that really is telling you the entire picture.

Ex: take only one observation (Same Bias)

Bias is telling us where are we (in average)
But surely not telling us what fluctuations we are getting.

When we start having fluctuations coming into the picture, we actually have to look at risk and quadratic risk of a estimator

* $x_1, x_2, \dots, x_n \sim \text{Ber}(p)$ assume $n=124$



If I were to repeat this 1000 times, every 1000 of those times I collect 124 data point's and then I do it again & again & again.

then in average the number I should get
should be closer to θ^* $\Rightarrow E[\hat{\theta}] = \underbrace{\theta^*}_{\text{same}} = \theta$
as $\doteq E[\hat{\theta}_n] - \theta$

Quadratic Risk \doteq L2 risk of $\hat{\theta}$

$$E[(\hat{\theta} - \theta)^2]$$

if $E[(\hat{\theta} - \theta)^2] \xrightarrow{n \rightarrow \infty} 0$, my Estimator $\hat{\theta}$ is

weakly consistent.

$$\mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2]$$

$$= \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2 + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta) \right]$$

$$= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2$$

$$= \text{Var}(\hat{\theta}) + \text{bias}^2$$

$$\Rightarrow \text{Quadratic Risk} = \text{Var}(\hat{\theta}) + \text{bias}^2$$

There is usually an inherent trade off b/w getting a small bias and small variance. if we reduce one too much then the other one gonna increase.

Example:

$$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} \text{Ber}(\theta)$$

$$\textcircled{1} \quad \hat{\theta} = \bar{x}_n \quad \mathbb{E}[\bar{x}_n] = \mathbb{E}[x]$$

$$\Rightarrow \text{Bias} = 0 \quad (\text{unbiased})$$

$$\begin{aligned} \text{Var}(\hat{\theta}) &= \text{Var}(\bar{x}_n) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \\ &= \frac{1}{n^2} \cdot n \cdot \sigma^2 \end{aligned}$$

$$= \frac{\sigma^2}{n} = \frac{\Theta(1-\theta)}{n}$$

$$\textcircled{1} \text{ Quadratic Risk} = \text{Var}(\hat{\theta}) + \text{Bias}^2$$

$$= \frac{\Theta(1-\theta)}{n} + 0$$

$$= \frac{\Theta(1-\theta)}{n}$$

This is just summarizing the performance of an Estimator in an I.I.D.

Three Estimators for θ

① \bar{X}_n

$$\text{Bias} = 0$$

$$\text{Var} = \frac{\Theta(1-\theta)}{n}$$

$$\text{Risk} = \frac{\Theta(1-\theta)}{n}$$

② 0.5

$$\text{Bias} = (0.5 - \theta)$$

$$\text{Var} = 0$$

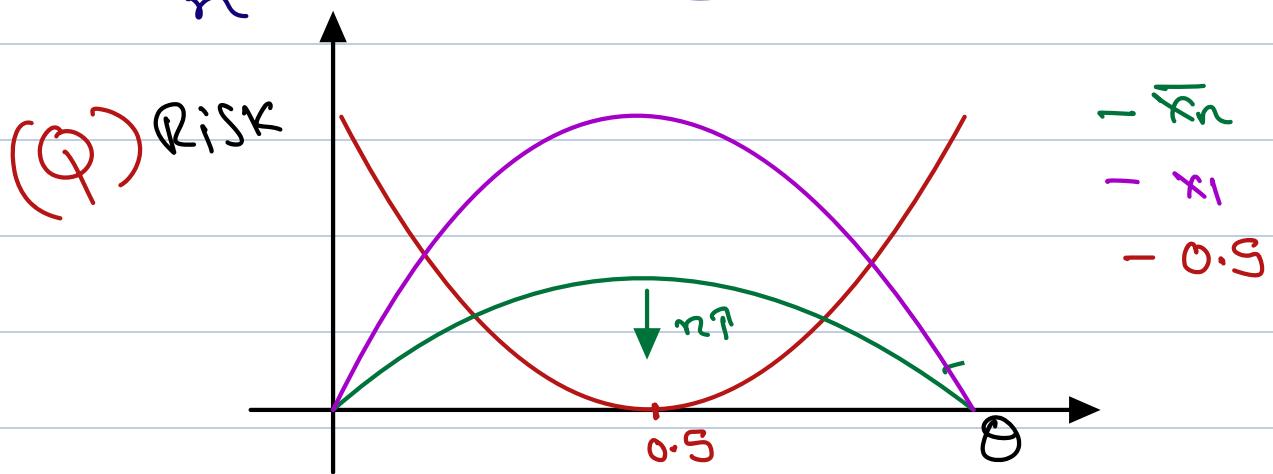
$$(0.5 - \theta)^2$$

③ x_1

$$\text{Bias} = 0$$

$$\text{Var} = \Theta(1-\theta)$$

$$\Theta(1-\theta)$$



in Non-Parametric estimation, all we do is Bias / variance trade off's all the time.

Confidence intervals (1)

Let $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$ be a statistical model based on observations X_1, \dots, X_n , and assume $\Theta \subseteq \mathbb{R}$.

Definition

Let $\alpha \in (0, 1)$.

- ▶ *Confidence interval (C.I.) of level $1 - \alpha$ for θ :* Any random (i.e., depending on X_1, \dots, X_n) interval \mathcal{I} whose boundaries do not depend on θ and such that:

$$\mathbb{P}_\theta [\mathcal{I} \ni \theta] \geq 1 - \alpha, \quad \forall \theta \in \Theta.$$

- ▶ *C.I. of asymptotic level $1 - \alpha$ for θ :* Any random interval \mathcal{I} whose boundaries do not depend on θ and such that:

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta [\mathcal{I} \ni \theta] \geq 1 - \alpha, \quad \forall \theta \in \Theta.$$

- * A Confidence interval is a random interval, The Bound's of the interval depends on random data.
- * \bar{X} is what the random thing that make fluctuate the confidence interval
- * Now we have interval and have its boundaries, Boundaries are not allowed to depend on unknown parameter (otherwise it's not a confidence interval)
Just like an Estimator, that depends on unknown parameter, is not a Estimator.
- * Confidence interval has to be something that we can compute once I collect data.

$$P_{\theta} [I \ni \theta] \geq 1-\alpha \quad \forall \theta \in \Theta$$

Probability that I contains θ

$\theta \in I$ (wrong) , $I \ni \theta$ (correct)

\Rightarrow The randomness is in I , θ is the true unknown value (deterministic)

The confidence interval is still a random variable

- * Now, if we start plugging in numbers instead of r.v. (x_1, x_2, \dots, x_n) , Ex: P(1,0,0,1,1,0, ..., in this case the random interval is actually for Ex: $[0.42, 0.65]$,
- * The Probability that $\theta \in [0.42, 0.65]$ is not $1-\alpha$, i.e. $P_\theta([0.42, 0.65] \ni \theta) \neq 1-\alpha$, It is either 0 (if not in there), 1 if its there.

Confidence intervals (2)

Example: Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(p)$, for some unknown $p \in (0, 1)$.

- ▶ LLN: The sample average \bar{X}_n is a strongly consistent estimator of p .

- ▶ Let $q_{\alpha/2}$ be the $(1 - \frac{\alpha}{2})$ -quantile of $\mathcal{N}(0, 1)$ and

$$\mathcal{I} = \left[\bar{X}_n - \frac{q_{\alpha/2} \sqrt{p(1-p)}}{\sqrt{n}}, \bar{X}_n + \frac{q_{\alpha/2} \sqrt{p(1-p)}}{\sqrt{n}} \right].$$

- ▶ CLT: $\lim_{n \rightarrow \infty} \mathbb{P}_p [\mathcal{I} \ni p] = 1 - \alpha, \quad \forall p \in (0, 1).$
- ▶ Problem: \mathcal{I} depends on p !

Confidence intervals (3)

Two solutions:

- ▶ Replace $p(1 - p)$ with $1/4$ in \mathcal{I} (since $p(1 - p) \leq 1/4$).
- ▶ Replace p with \bar{X}_n in \mathcal{I} and use Slutsky's theorem.

CLT:

$$\frac{\sqrt{n}(\bar{X} - P)}{\sqrt{P(1-P)}} \xrightarrow[n \rightarrow \infty]{d} N(0,1)$$

$$\Rightarrow \bar{X} \sim N(P, \frac{P(1-P)}{n})$$

Using CLT, the distribution of
 $\frac{\sqrt{n}(\bar{X} - P)}{\sqrt{P(1-P)}}$ approximately equal to
 $N(0,1)$

$$\Rightarrow \lim_{n \rightarrow \infty} P_p \left(\left| \frac{\sqrt{n}(\bar{X} - P)}{\sqrt{P(1-P)}} \right| \leq q_{\alpha/2} \right) = 1 - \alpha$$

$$\Rightarrow \lim_{n \rightarrow \infty} P_p \left(-q_{\alpha/2} \leq \frac{\sqrt{n}(\bar{X} - P)}{\sqrt{P(1-P)}} \leq q_{\alpha/2} \right) = 1 - \alpha$$

$$\Rightarrow \lim_{n \rightarrow \infty} P_p \left(P - \sqrt{\frac{P(1-P)}{n}} \cdot q_{\alpha/2} \leq \bar{X} \leq P + \sqrt{\frac{P(1-P)}{n}} \cdot q_{\alpha/2} \right) = 1 - \alpha$$

MIT OpenCourseWare

<https://ocw.mit.edu>

18.650 / 18.6501 Statistics for Applications

Fall 2016

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.