

Maximum Likelihood Review

Maximal likelihood estimation is a method of estimating the parameters of a statistical model given data.

By statistical model on a sample space X we mean a set of distributions (actually measures) on X . If we write $PM(X)$ for the space of all possible distributions (measures) over X then a model is a subset $M \subset PM(X)$. The elements of M are indexed by a parameter θ with values in a parameter space T , that is,

$$M = \{P_\theta \mid \theta \in T\}$$

where each P_θ is a member of the set $PM(X)$.

A model is parametric if T is finite dimensional. Usually $T \subset \mathbb{R}^d$ for some $d \in \mathbb{N}$. If $\dim(T) = \infty$ then M is a non-parametric model.

The canonical problem of statistics is to take observations x_1, \dots, x_n , $x_i \in X$, which we model as random variables X_1, \dots, X_n which we assume are drawn from P_θ , i.e.

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P_\theta \quad \theta \in T$$

and use them to tell us something about the value of θ .

Estimators

Assume there is a sample x_1, x_2, \dots, x_n of n iid observations coming from a true but unknown distribution $P_0(\cdot)$.

Let us assume that $P_0 \in \mathcal{M}$, i.e. $\exists \theta_0 \in \mathcal{T}$ st. $P_{\theta_0} = P_0$. Let us also assume that \mathcal{T} is finite dimensional.

We would like an estimator $\hat{\theta} = f(x_1, \dots, x_n)$ which is close to θ_0 .

An estimator is a statistic, i.e. a function of the data, that is used to infer the value of an unknown parameter in a statistical model. Some estimators can be quite complex.

Let's have some important props of estim's: Letting $X = \{x_1, \dots, x_n\}$ and $\hat{\theta} = f(X)$ assumed 1-dimensional

$$1) \text{MSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta_0)^2]$$

- note that this expectation is with respect to samples (populations of size n) drawn from P_0 .
this is the mean or expected squared error; small if the estimator is good

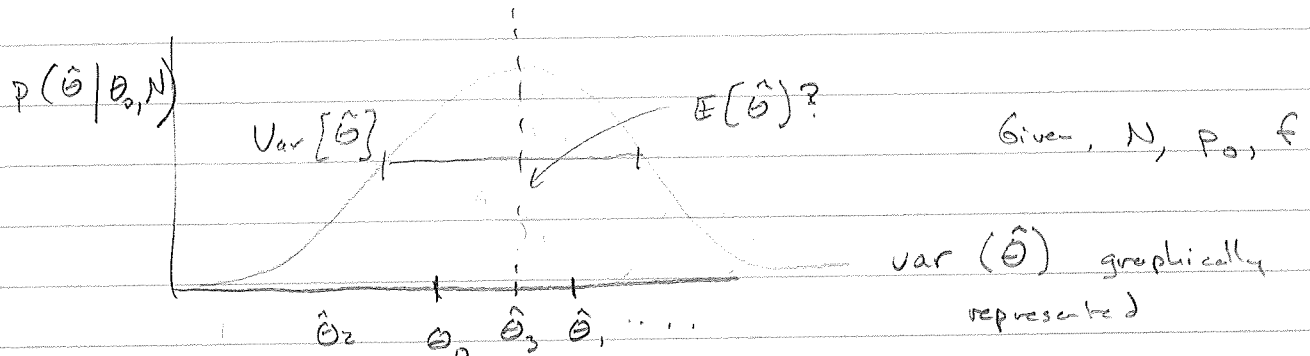
$$2) \text{var}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \mathbb{E}(\hat{\theta}))^2]$$

this is, for a given sample size (and P_0) the characteristic "spread" of the estimator

Datasets / samples

Estimator $\hat{\theta}$

$$\begin{array}{l}
 X_1 = \{x_1 = \dots x_2 = \dots x_3 = \dots x_n = \dots\} \\
 X_2 = \\
 X_3 = \\
 \vdots \\
 \vdots
 \end{array}
 \quad , \quad
 \begin{array}{l}
 \hat{\theta}_1 = \\
 \hat{\theta}_2 = \\
 \hat{\theta}_3 = \\
 \vdots \\
 \vdots
 \end{array}$$

with $x_n \sim \text{iid } P_0$, $P_\theta(x) = \prod P_0(x)$ 

3) Bias $B(\hat{\theta}) = E(\hat{\theta}) - \theta_0$

is the distance between the average of estimates over all samples of a given size and θ_0 .

Statisticians like unbiased estimators,
 i.e. $B(\hat{\theta}) = 0$

Note for now the relationship:

$$MSE(\hat{\theta}) = \text{var}(\hat{\theta}) + (\text{Bias}(\hat{\theta}))^2$$

proof

$$\begin{aligned} MSE(\hat{\theta}) &= E[(\hat{\theta} - \theta_0)^2] \\ &= E[(\hat{\theta} - E[\hat{\theta}] + E[\hat{\theta}] - \theta_0)^2] \\ &= E[(\hat{\theta} - E[\hat{\theta}])^2] \\ &\quad (A) \quad + E[(\hat{\theta} - E[\hat{\theta}])(E[\hat{\theta}] - \theta_0)] \\ &\quad (B) \quad + E[(E[\hat{\theta}] - \theta_0)(\hat{\theta} - E[\hat{\theta}])] \\ &\quad + E[(E[\hat{\theta}] - \theta_0)^2] \end{aligned}$$

noting that $E[\hat{\theta}]$, θ_0 are constant means that (A) = $(E[\hat{\theta}] - E[\hat{\theta}])(E[\hat{\theta}] - \theta_0) = 0$ and same for (B). So

$$\begin{aligned} MSE(\hat{\theta}) &= E[(\hat{\theta} - E[\hat{\theta}])^2] + (E[\hat{\theta}] - \theta_0)^2 \\ &= \text{VAR}(\hat{\theta}) + (\text{BIAS}(\hat{\theta}))^2 \end{aligned}$$

We will return to this later.

Let's look at a simple estimator of a population mean and variance (think heights if you like). Assume the true distribution is $N(\mu_0, \sigma_0^2)$. Let N be the sample size and $X = \{x_1, \dots, x_N\}$, $x_i \sim \text{iid } N(\mu_0, \sigma_0^2)$

Let's consider some estimators: for μ

$$A) \hat{\mu}_A = f_A(x_1, \dots, x_N) = 0$$

Maximum Likelihood \rightarrow

$$B) \hat{\mu}_B = f_B(x_1, \dots, x_N), \hat{\mu}_B = \underset{\mu}{\text{argmax}} L(X; \mu)$$

$$C) \hat{\mu}_C = f_C(x_1, \dots, x_N) = \lambda \mu_1 + (1-\lambda) \hat{\mu}_B$$

some "guess"
 $\mu_1 \approx \mu_0$

$$A) \text{Var}(\hat{\mu}_A) = 0$$

$$\text{BIAS}(\hat{\mu}_A) = \hat{\mu}_0$$

$$\text{MBE}(\hat{\mu}_A) = \hat{\mu}_0^2$$

B) This is the maximum likelihood estimator where we have a parametric model

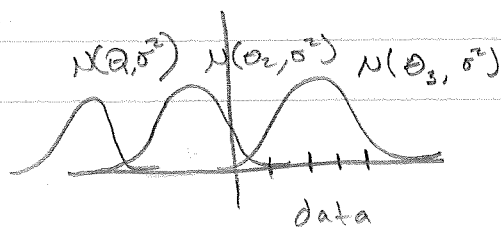
$$X \sim N(\mu, \sigma^2), \theta = \{\mu, \sigma^2\}$$

and

Likelihood function \rightarrow

$$L(X; \theta) = \prod_{i=1}^N P(x_i | \theta)$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$



which θ is preferred?

ML Estimation - by easy example

If we want to extremize $L(x; \theta)$ we look for

$$\frac{\partial L(x; \theta)}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L(x; \theta)}{\partial \theta^2} \leq 0$$

to maximize $L(x; \theta)$ w.r.t. θ

Important $\frac{\partial}{\partial \theta} L(x; \theta)$ is almost always nasty so it is usual to work with $\frac{\partial}{\partial \theta} \log L(x; \theta)$ which has the same max because \log is monotonically increasing. Often the resulting maximization is easier

Claim

$$\operatorname{argmax}_{\theta} \log(L(\theta)) = \operatorname{argmax}_{\theta} (L(\theta))$$

Check 0 points

$$0 = \frac{\partial \log L(\theta)}{\partial \theta}$$

$$= \frac{1}{L(\theta)} \frac{\partial L(\theta)}{\partial \theta}$$

$$= \frac{\partial L(\theta)}{\partial \theta} \quad \checkmark$$

$$\frac{\partial \log L(\theta)}{\partial \theta} = 0 \Leftrightarrow \frac{\partial L(\theta)}{\partial \theta} = 0$$

Check curvature

$$\operatorname{sign} \left(\frac{\partial^2}{\partial \theta^2} \log L(\theta) \right) = \operatorname{sign} \left(\frac{\partial}{\partial \theta} \frac{1}{L(\theta)} \frac{\partial L(\theta)}{\partial \theta} \right)$$

$$= \operatorname{sign} \left(-\frac{1}{L(\theta)^2} \frac{\partial L(\theta)}{\partial \theta} \frac{\partial L(\theta)}{\partial \theta} + \frac{1}{L(\theta)} \frac{\partial^2 L(\theta)}{\partial \theta^2} \right)$$

$$= \operatorname{sign} \left(\frac{\partial^2 L(\theta)}{\partial \theta^2} \left(\frac{1}{L(\theta)} - \frac{1}{L(\theta)^2} \right) \right)$$

true if pos.

$$\frac{1}{L(\theta)} - \frac{1}{L(\theta)^2} \stackrel{?}{\geq} 0$$

$$L(\theta)^2 - L(\theta) \geq 0 \quad L(\theta)^2 \geq L(\theta) \quad \checkmark$$

Back to mean estimator, B

argmax $\log L(x; \theta)$

occurs at $\frac{\partial}{\partial \theta} \log L(x; \theta) = 0$ here $\theta = \mu$

$$\begin{aligned} \frac{\partial}{\partial \mu} \log L(x; \mu) &= \sum_n \frac{\partial}{\partial \mu} \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_n - \mu)^2}{2\sigma^2}} \\ &= \sum_n \frac{\partial}{\partial \mu} -\frac{(x_n - \mu)^2}{2\sigma^2} \end{aligned}$$

$$= \sum_n 2(x_n - \mu)$$

$$= 2 \left(\sum_n x_n - n \mu \right) = 0$$

$$\Rightarrow \hat{\mu}_B = \frac{\sum_{n=1}^N x_n}{N}$$

Bias of estimator B:

$$\begin{aligned} E[\hat{\mu}_B] &= E\left[\frac{\sum_{n=1}^N x_n}{N}\right] = \frac{1}{N} \sum_{n=1}^N E[x_n] \\ &= \frac{1}{N} \cdot N \cdot \mu = \mu \end{aligned}$$

$\hat{\mu}_B$ is unbiased, i.e. $\text{BIAS}(\hat{\mu}_B) = 0$

Variance of estimator:

$$\begin{aligned} \text{Var}[\hat{\mu}_B] &= \text{Var}\left[\frac{\sum_{n=1}^N x_n}{N}\right] = \frac{1}{N^2} \text{Var}\left[\sum_{n=1}^N x_n\right] \\ &= \frac{1}{N^2} \sum_{n=1}^N \text{Var}(x_n) \\ &= \frac{\sigma^2}{N} \end{aligned}$$

So $\text{MSE}(\hat{\mu}_B) = \text{VAR}(\hat{\mu}_B) = \frac{\sigma^2}{N}$ which, in the limit of $N \rightarrow \infty$, $\frac{\sigma^2}{N} \rightarrow 0$.

Estimation : the core of stats, along with frequentist inference.

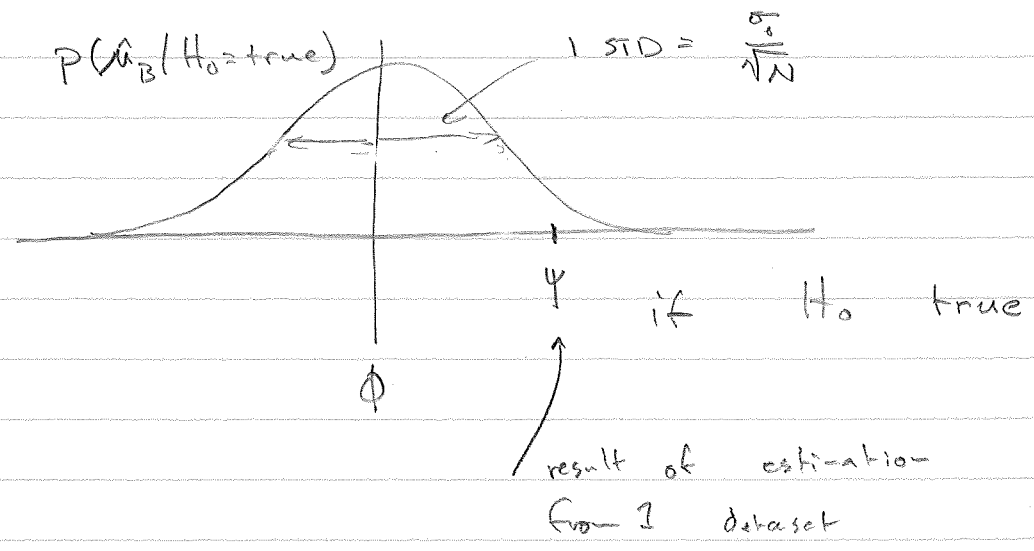
Assuming that we know σ_0^2 let's hypothesize that $X_n \sim i.i.d N(\phi, \sigma_0^2)$ then to test this hypothesis we run the following thought experiment. Assume that we use estimator B and that we draw an infinite number of sample populations of size N from the model above.

We already know that, in this case, $E[\hat{\mu}_B] = \phi$ and $Var[\hat{\mu}_B] = \frac{\sigma_0^2}{N}$

What can be demonstrated is that, for sufficiently large N and under mild regularity assumptions

$$\hat{\mu}_B \sim N(\phi, \frac{\sigma_0^2}{N})$$

Now, given a particular dataset X_D we compute $\hat{\mu}_B(X_D) = \psi$. We have a picture like this

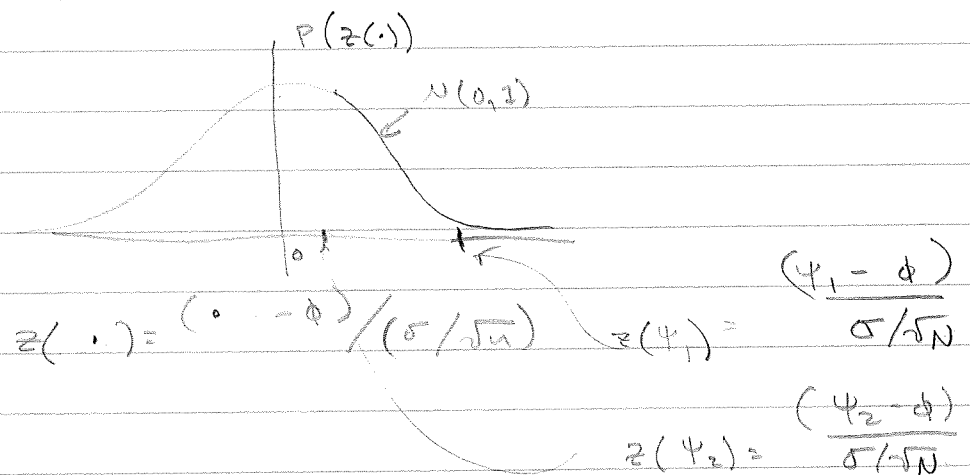


lecture 7

To think about confidence intervals and drawing conclusions about a null hypothesis it will help, in many cases to "normalize" our estimator to a z-score, i.e. a $N(0,1)$ variable.

If $\hat{\mu}_B \sim N(\phi, \sigma^2/n)$ then

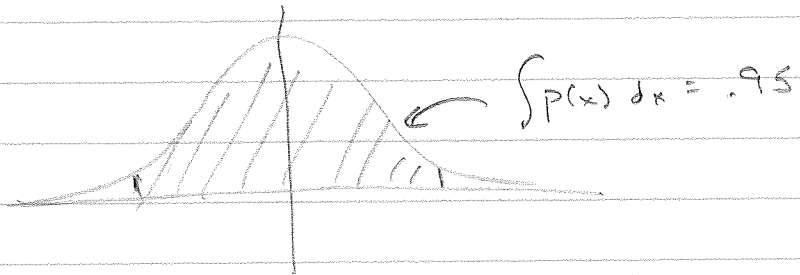
$$z(\cdot) = \frac{(\hat{\mu}_B - \phi)}{\sigma/\sqrt{n}} \sim N(0,1) \quad \text{Pf. trivial}$$



Frequentist inference rejects a null hypothesis when the value of an estimator computed on a real sample is surprising (i.e. low probability) under the null hypothesis.

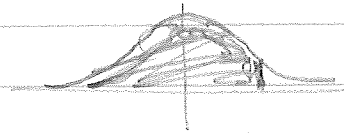
If, for instance, $z(\psi)$ falls within the 95% confidence interval then we fail to reject H_0 . What is the 95% confidence interval?

It's, under the standard normal



ie. the value z for which
 $P(-z \leq Z \leq z) = 1 - \alpha = .95$ for $Z \sim N(0,1)$

This occurs at $\Phi(z) = P(Z \leq z) = 1 - \frac{\alpha}{2} = 0.975$
 where $\Phi(z)$ is the CDF of $N(0,1)$



$$\Phi(z) = \int_0^z N(x; 0, 1) dx$$

ie. $\Phi^{-1}(\Phi(z) = 0.975) = 1.96$ so

$$0.95 = 1 - \alpha = P(-z \leq Z \leq z) = P\left(-1.96 \leq \frac{\psi - \phi}{\sigma/\sqrt{n}} \leq 1.96\right)$$

and, thusly

$$0.95 = P\left(\psi - 1.96 \cdot \left(\frac{\sigma}{\sqrt{n}}\right) \leq \phi \leq \psi + 1.96 \cdot \left(\frac{\sigma}{\sqrt{n}}\right)\right)$$

What does this mean? It means that both ϕ , under our assumptions, is likely to be in this band but, perhaps more importantly, that ϕ and ψ can be flipped saying, in effect, that if the true parameter is ϕ , ψ outside of

$$\left(\phi - 1.96 \left(\frac{\sigma}{\sqrt{n}}\right), \phi + 1.96 \left(\frac{\sigma}{\sqrt{n}}\right) \right)$$

would be surprising and could be taken

as evidence (at 5% c.f.) to reject the null hypothesis.

Note that a p-value is the value of α one would have to choose to reject the null hypothesis given the observed sample, i.e. p-value is sol'n to (in this case)

$$\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) = \frac{\psi - \phi}{\sigma/\sqrt{n}} \quad \text{or} \quad 1 - \frac{\alpha}{2} = \Phi\left(\frac{\psi - \phi}{\sigma/\sqrt{n}}\right)$$

$$\text{p-value} = \alpha = 2\left(1 - \Phi\left(\frac{\psi - \phi}{\sigma/\sqrt{n}}\right)\right)$$

Revisiting all the way back we have estimators

	A	B	C
var	0	$\frac{\sigma^2}{N}$	
bias	$\hat{\mu}_A$	0	
mse	$\hat{\mu}_A^2$	$\frac{\sigma^2}{N}$	

Which suggests that a) no one estimator is always best and b) that we might be able to explore the spectrum of estimators and engineer estimators that work well for our problem.

For instance, what if we think we know the true "answer", is $\mu_1 \approx \mu_0$, like estimator C.

$$\lambda \hat{\mu}_1 + (1 - \lambda) \hat{\mu}_B$$

Estimator C

$$\text{Bias} (\lambda \mu_0 + (1-\lambda) \hat{\mu}_B)$$

$$= \mathbb{E} [\lambda \mu_0 + (1-\lambda) \hat{\mu}_B] - \mu_0$$

$$= \lambda \mu_0 + (1-\lambda) \mu_0 - \mu_0 = \lambda (\mu_1 - \mu_0)$$

$$\text{Var} (\lambda \mu_0 + (1-\lambda) \hat{\mu}_B)$$

$$= (1-\lambda)^2 \text{Var} (\hat{\mu}_B) = (1-\lambda)^2 \frac{\sigma^2}{N}$$

	A	B	C
var	0	$\frac{\sigma^2}{N}$	$(1-\lambda)^2 \frac{\sigma^2}{N}$
bias	$\hat{\mu}_A$	0	$\lambda (\mu_1 - \mu_0)$
mse	$\hat{\mu}_A^2$	$\frac{\sigma^2}{N}$	$(1-\lambda)^2 \frac{\sigma^2}{N} + \lambda^2 (\mu_1 - \mu_0)^2$

But, wait, this suggests that bias could lead to lower MSE, and here will who -

$$\frac{\sigma^2}{N} > (1-\lambda)^2 \frac{\sigma^2}{N} + \lambda^2 (\mu_1 - \mu_0)^2$$

$$\frac{\sigma^2}{N} (1 - (1-\lambda)^2) > \lambda^2 (\mu_1 - \mu_0)^2$$

$$\frac{\sigma^2}{N} > \frac{\lambda^2}{2\lambda - \lambda^2} (\mu_1 - \mu_0)^2$$

$$\frac{\sigma^2}{N} > \frac{1}{\frac{2}{\lambda} - 1} (\mu_1 - \mu_0)^2$$

$$\frac{\sigma^2}{N} > \frac{\lambda}{2-\lambda} (\mu_1 - \mu_0)^2$$

which could easily happen.

The take-home here is that maximum likelihood estimators are but one family of estimator and, further, that bias is not a bad word, in fact introducing it can lower the variance of your estimator and potentially the MSE of your estimator as well.

Frequentist inference involves computing and using confidence intervals for the sample variance of estimators.

Bias introduced in the form of regularization and priors may make for better estimation provided that the introduced bias is helpful.