

Lesson 17

Statistics

So far, we have seen that
we can calculate probabilities
using densities of r.v.'s.

However, in real life, we usually
have a set of data, whose
density is unknown. In order
to calculate probabilities,
we have to know the density

with associated with the data.

Let's Start from the simplest

case, where we know that

the data X_1, X_2, \dots, X_n is

i.i.d, and we know the

type of distribution of X_i (e.g.
we know that it is $\text{Pois}(\lambda)$),

but we don't know its mean.

We learned that

$$\bar{X} =$$

Can be viewed as an
"estimate" of the mean
 $E[X]$. In particular, the
LLNs guarantees that
for large ~~n~~ m's.



When one wishes to calculate the mean of a distribution, one can generate a sample X_1, \dots, X_n and take $\frac{\sum X_i}{n}$ as an

estimate of the mean.

We also learned that
 \bar{X} is distributed normally,
according to the CLT.

We are interested in
constructing the so-called
Confidence Intervals (C.I.'s)
for \bar{X} , the estimate of
 $\mu = E[X]$.

Roughly speaking, this is an
interval that contains
 $\mu = E[X]$ with a
high Probability

More precisely, we fix a level of confidence, $1-\alpha$, where α , the level of doubt, is typically a small number.

We replace the estimate

\bar{X} with two estimators

L, U , s.t.

$$P(L \leq \mu \leq U) \geq 1-\alpha$$

Remark: Since we model data X_1, \dots, X_n as ^{i.i.d} r.v.'s, and estimates of parameters of F_X are regarded as functions of X_1, \dots, X_n , estimates are also r.v.'s themselves.

Therefore, \bar{X} , L , and U are r.v.'s whose distributions are dependent on μ .

$[L, U]$ is called

the $(1-\alpha)$ confidence intervals

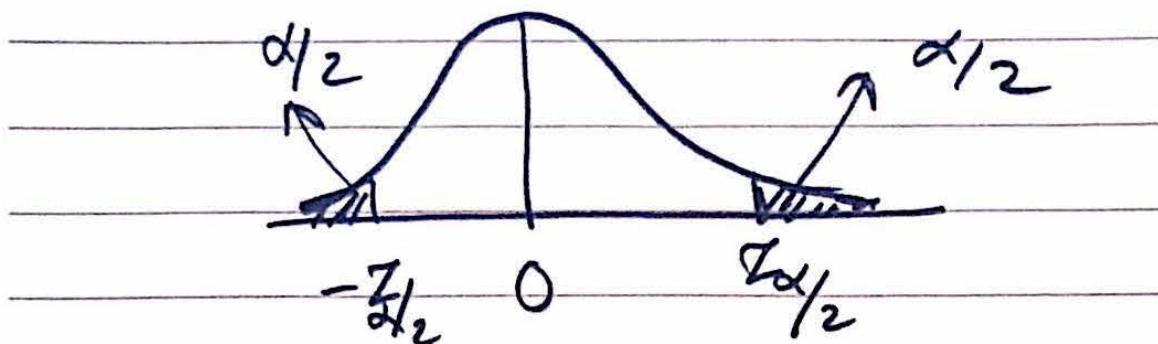
Because \bar{X} is approximately distributed according to a Gaussian, and because we want a symmetric interval

around \bar{X} , as we follow

the following method:

~~Remember~~ Observe that

for a standard normal Z :



$$1-\alpha = P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2})$$

$$= P$$

Exercise: Build a "one-sided"

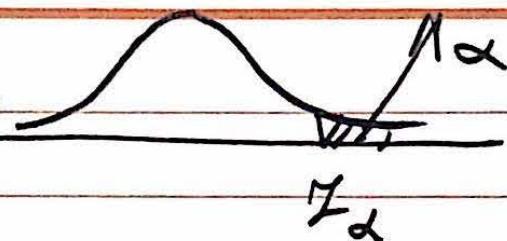
$1-\alpha$ Confidence Interval

$[L, +\infty)$ or $(-\infty, U]$,

using the same methodology.

Hint: Use the following

$$P(Z \leq z_\alpha) = 1-\alpha$$



$$P(Z \geq -z_\alpha) = 1-\alpha$$

Interpretation of C.I's

is subtle. We are tempted
to interpret the concept of
a C.I by "the true
parameter lies in the

C.I with probability $1-\alpha$
(e.g. 95%).

Here, μ is deterministic and
 L, U are r.r.'s, so based
on different experiments

L, U may have different realizations. In other words, randomness is met in μ , but is in L, U !

The precise meaning is
the following:

Confidence Interval ~~for~~

The mean Using Sample

Mean & Sample Standard Deviation

We assumed to have σ ,

which is NOT ALWAYS

a realistic assumption.

If X_i 's are i.i.d $N(\mu, \sigma^2)$,

even if we don't have

a large sample, the

quantity

$T =$

is distributed according to

a Student t distribution,

with $n-1$ degrees of freedom

The pdf of T is

$$f_T(x) \propto \left(1 + \frac{x^2}{n-1} \right)^{-\frac{n}{2}}$$

$$\frac{1}{\sqrt{n-1}} \frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{1}{2})}$$

$$-\infty < x < \infty$$

The $1-\alpha$ C.I. is

where

s_x : Sample Standard Deviation

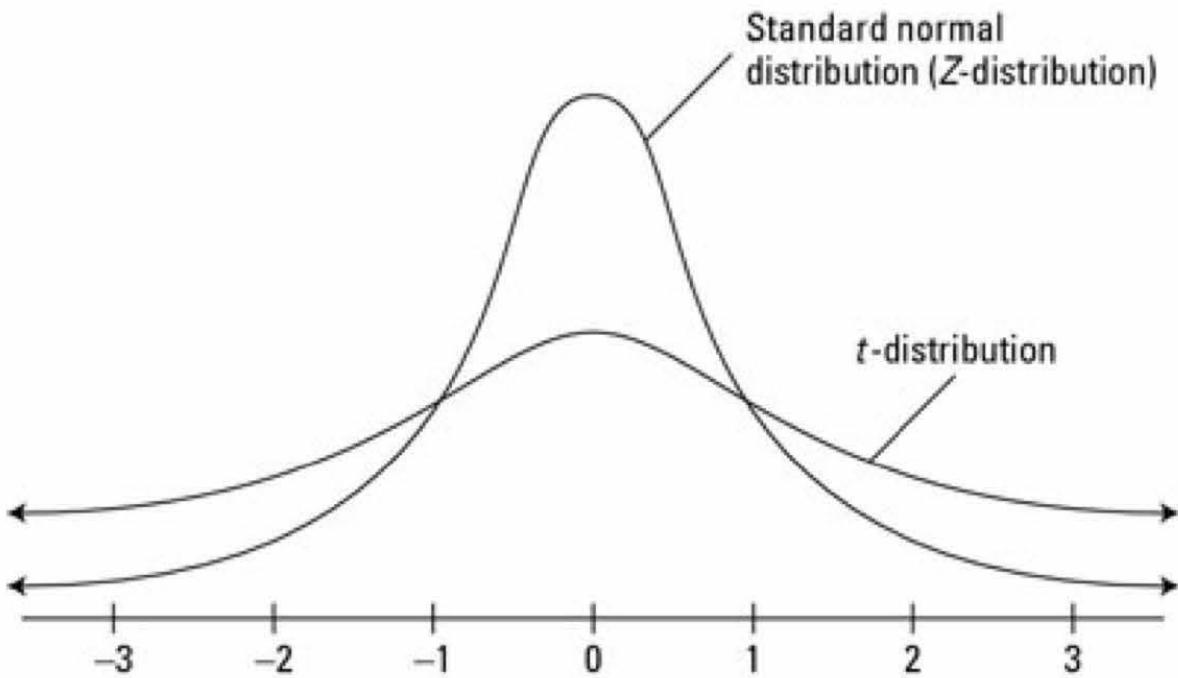
$$s_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

The t-distribution looks like

a normal distribution i.e. $N(0,1)$.

it is bell-shaped & symmetric

but has heavier tails.



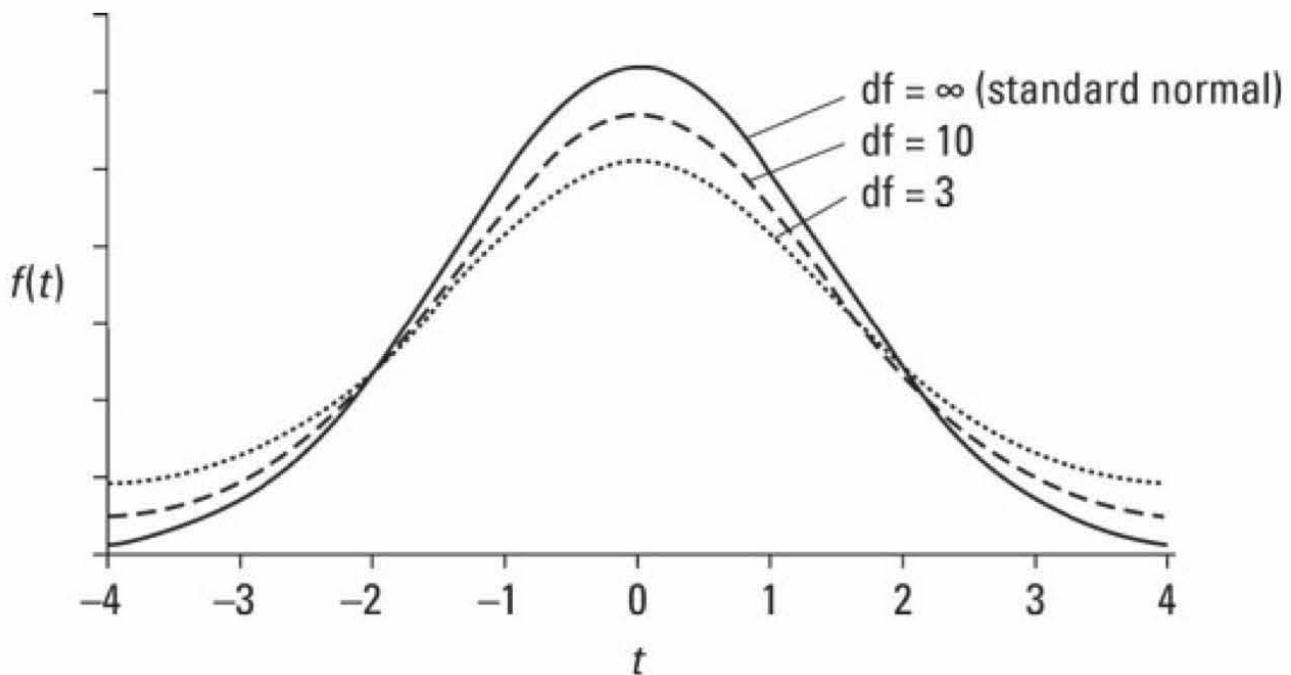
t-distributions are members of
a family of distributions
parameterized by their degrees
of freedom.

$t_{n-1, \alpha/2}$ is a cut off point

that signifies a tail probability

of $\alpha/2$. We usually use

tables to determine $t_{n-1, \alpha/2}$



Example: Using the CLT to
estimate proportions / Pollster
Confidence Intervals

Assume that we are interested
to estimate the proportion of

- a population that has some
characteristic, e.g.,
- proportion of USC students
who study engineering
 - proportion of Californians who support
measure X

This is an attempt to
 estimate p for $X_i \sim \text{Ber}(p)$
 $(i=1, 2, \dots, N \rightarrow \text{population size})$

by using a sample mean

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$$

where

n : sample size

$\sum_{i=1}^n X_i$: number of samples

that have the desired

characteristic

For large n $\hat{p} \sim N(\mu_X, \frac{\sigma_X^2}{n})$

where $\mu_X =$

$$\sigma_X^2 =$$

The $(1-\alpha)$ C.I for \hat{p} is

The "Margin of Error" is

defined as:

For $\alpha = 0.05$, determine the sample size n such that the Margin of Error is less than or equal to 3%

Handwriting practice lines. The page features a decorative scalloped border and ten sets of horizontal lines for handwriting practice.

Confidence Intervals for

Frequentist Hypothesis Testing

C.I.s provide a basis for

frequentist test of hypotheses.

The result of a hypothesis

testing procedure depends on

whether or not a test statistic

falls inside a $(1-\alpha)$ C.I.

A hypothesis is usually a

claim about a parameter of

a population (e.g., mean, variance, median, etc) to which we don't have access, using a statistic (a function of n i.i.d. samples from that population) with

some level of confidence.

The null hypothesis is a statement about a population parameter ~~is~~ that is sought to be rejected by using

the statistic as evidence.
If the statistic is too unlikely to have come from a population with hypothesized parameter, it is considered

as evidence that rejects

the null hypothesis.

Procedure for testing two-sided hypotheses about the mean of a population:

(i) Formulate competing

hypotheses about μ_x

$$\left\{ \begin{array}{l} H_0: \mu_x = \mu_0 \text{ (Null hypothesis)} \\ H_1: \mu_x \neq \mu_0 \text{ (Alternative hypothesis)} \end{array} \right.$$

(ii) Specify the "two-sided"

$$(1-\alpha) \text{ C.I. for } Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

which is

$$[-Z_{\alpha/2}, Z_{\alpha/2}]$$

(iii) obtain the test statistic

\bar{x} from the ~~samples~~

(iv) Reject H_0 (= support H_1) if

$$\frac{\bar{x} - \mu_0}{\sigma_x / \sqrt{n}} \notin [-z_{\alpha/2}, z_{\alpha/2}]$$

Fail to reject H_0

$$\text{if } \frac{\bar{x} - \mu_0}{\sigma_x / \sqrt{n}} \in (-z_{\alpha/2}, z_{\alpha/2})$$

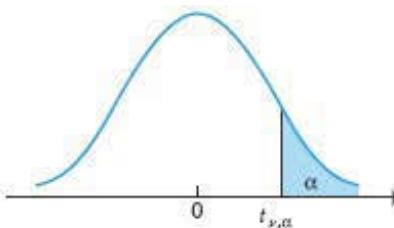
Exercise: Based on the one-sided C.I.'s you found in the previous exercise, construct a procedure for testing the following one-sided hypotheses.

hypotheses:

$$\left\{ \begin{array}{l} H_0: \mu \leq \mu_0 \\ H_1: \mu > \mu_0 \end{array} \right. \text{ or } \left\{ \begin{array}{l} H_0: \mu \geq \mu_0 \\ H_1: \mu < \mu_0 \end{array} \right.$$

Exercise: The average cost of a hotel room in Chicago is said to be \$168 per night. A random sample of 25 hotels resulted in $\bar{x} = \$172.5$ and $s = \$15.4$. Test at

$$\alpha = 0.05.$$



For selected probabilities, α , the table shows the values $t_{v,\alpha}$ such that $P(t_v > t_{v,\alpha}) = \alpha$, where t_v is a Student's t random variable with v degrees of freedom. For example, the probability is .10 that a Student's t random variable with 10 degrees of freedom exceeds 1.372.

PROBABILITY OF EXCEEDING THE CRITICAL VALUE						
v	0.10	0.05	0.025	0.01	0.005	0.001
1	3.078	6.314	12.706	31.821	63.657	318.313
2	1.886	2.920	4.303	6.965	9.925	22.327
3	1.638	2.353	3.182	4.541	5.841	10.215
4	1.533	2.132	2.776	3.747	4.604	7.173
5	1.476	2.015	2.571	3.365	4.032	5.893
6	1.440	1.943	2.447	3.143	3.707	5.208
7	1.415	1.895	2.365	2.998	3.499	4.782
8	1.397	1.860	2.306	2.896	3.355	4.499
9	1.383	1.833	2.262	2.821	3.250	4.296
10	1.372	1.812	2.228	2.764	3.169	4.143
11	1.363	1.796	2.201	2.718	3.106	4.024
12	1.356	1.782	2.179	2.681	3.055	3.929
13	1.350	1.771	2.160	2.650	3.012	3.852
14	1.345	1.761	2.145	2.624	2.977	3.787
15	1.341	1.753	2.131	2.602	2.947	3.733
16	1.337	1.746	2.120	2.583	2.921	3.686
17	1.333	1.740	2.110	2.567	2.898	3.646
18	1.330	1.734	2.101	2.552	2.878	3.610
19	1.328	1.729	2.093	2.539	2.861	3.579
20	1.325	1.725	2.086	2.528	2.845	3.552
21	1.323	1.721	2.080	2.518	2.831	3.527
22	1.321	1.717	2.074	2.508	2.819	3.505
23	1.319	1.714	2.069	2.500	2.807	3.485
24	1.318	1.711	2.064	2.492	2.797	3.467
25	1.316	1.708	2.060	2.485	2.787	3.450
26	1.315	1.706	2.056	2.479	2.779	3.435
27	1.314	1.703	2.052	2.473	2.771	3.421
28	1.313	1.701	2.048	2.467	2.763	3.408
29	1.311	1.699	2.045	2.462	2.756	3.396
30	1.310	1.697	2.042	2.457	2.750	3.385
40	1.303	1.684	2.021	2.423	2.704	3.307
60	1.296	1.671	2.000	2.390	2.660	3.232
100	1.290	1.660	1.984	2.364	2.626	3.174
∞	1.282	1.645	1.960	2.326	2.576	3.090
v	0.10	0.05	0.025	0.01	0.005	0.001

Parameter Estimation

Assume that X_1, \dots, X_n are i.i.d.

r.r.'s and $X_i \sim f(x_i | \theta)$ where θ

is an unknown parameter

(or parameter vector). The conditional

notation here does not necessarily

mean that θ is a r.r. It means

that the distribution (pdf/pmf)

depends on θ .

We wish to use X_1, X_2, \dots, X_n to estimate θ .

More precisely, we want to find a "statistic" T

$$T = T(X_1, X_2, \dots, X_n)$$

which is a deterministic function of X_1, \dots, X_n such that $\theta \approx T$.

Usually, we denote this T as $\hat{\theta}_n$.

We have two goals:

(i) Find a "good" $\hat{\theta}$

(ii) Understand what "good"

means

There are three measures

for determining what "good"

means:

(a) Bias

$$\text{Bias}(\hat{\theta}_n, \theta) =$$

If $\text{Bias}(\hat{\theta}_n, \theta) = 0$, we say

$\hat{\theta}_n$ is an unbiased estimate
for θ .

If $\lim_{n \rightarrow \infty} \text{Bias}(\hat{\theta}_n, \theta) = 0$,

then $\hat{\theta}_n$ is called

"asymptotically unbiased"

(b) Mean-Squared Error

$$\text{MSE}(\hat{\theta}_n, \theta) =$$

Theorem (The Bias-Variance

Decomposition) : For a single estimate $\hat{\theta}_n$,

$$MSE(\hat{\theta}_n, \theta) = \text{Var}(\hat{\theta}_n) + \text{Bias}^2(\hat{\theta}_n, \theta)$$

Proof)

Remark: For the same mean-squared error, there is a compromise between the variance of the estimator and its bias.

One cannot decrease the

Variance and the bias

Simultaneously In Machine

Learning, this is called the

bias-variance trade-off.

(c) Consistency

$\hat{\theta}_n$ is called consistent iff

$$\hat{\theta}_n \xrightarrow{i.p.} \theta$$



(Population) Median:

For a random variable X ,

m is a median iff

$$P(X \leq m) \leq \frac{1}{2} \text{ and } P(X > m) > \frac{1}{2}$$

Sample Median:

Assume that X_1, X_2, \dots, X_n are

i.i.d r.v.'s. The sample median

$$\therefore m_n = \text{median}(X_1, \dots, X_n).$$

Exercise : Show that \bar{m}_n is

an unbiased estimate of $E[X]$

if f_x is symmetric about $\mu_x = E[X]$

The Bias-Variance Decomposition for Convergence of Random Variables (Estimates)

Convergence in Mean-Squared

Sense is important in science

and engineering.

For an estimator, one can prove

Consistency (Convergence in

probability) using the bias-

Variance decomposition.

$$\text{MSE}(\hat{\theta}_n, \theta) = \text{Var}(\hat{\theta}_n) + \text{Bias}^2(\hat{\theta}_n, \theta)$$

If

$$\lim_{n \rightarrow \infty} \text{Bias}(\hat{\theta}_n, \theta) = 0$$

i.e. $\hat{\theta}_n$ is asymptotically unbiased

and

$$\lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}_n) = 0$$

then $\lim_{n \rightarrow \infty} \text{MSE}(\hat{\theta}_n, \theta)$

$$= \lim_{n \rightarrow \infty} E[(\hat{\theta}_n - \theta)^2] = 0$$

i.e. $\hat{\theta}_n \xrightarrow{\text{m.s.}} \theta$, which

implies $\hat{\theta}_n \xrightarrow{\text{i.p.}} \theta$

i.e., $\hat{\theta}_n$ is a consistent estimator of θ .

Mean-Squared LLN
(Theorem)
Assume that X_1, X_2, \dots, X_n are i.i.d and $E[X^2] < \infty$. Then

$$\bar{X}_n \xrightarrow{\text{m.s.}} E[X]$$

Proof) The bias variance decomposition.

(Important) Exercise: Assume that

i.i.d r.v.'s X_1, X_2, \dots, X_n are

similarly distributed and

$X_i \sim N(0, 1/\mu_i^2)$. Define $Y_n = \sqrt{n} X_n$

Does Y_n converge in probability

to 0?

We use multiple methods to

solve this problem:

1. *W*hat *is* *the* *name* *of* *the* *city* *you* *live* *in*?

2. *W*hat *is* *the* *name* *of* *the* *country* *you* *live* *in*?

3. *W*hat *is* *the* *name* *of* *the* *state* *you* *live* *in*?

4. *W*hat *is* *the* *name* *of* *the* *town* *you* *live* *in*?

5. *W*hat *is* *the* *name* *of* *the* *post* *office* *you* *use*?

6. *W*hat *is* *the* *name* *of* *the* *street* *you* *live* *on*?

7. *W*hat *is* *the* *name* *of* *the* *house* *you* *live* *in*?

8. *W*hat *is* *the* *name* *of* *the* *room* *you* *live* *in*?

9. *W*hat *is* *the* *name* *of* *the* *bed* *you* *sleep* *in*?

10. *W*hat *is* *the* *name* *of* *the* *table* *you* *eat* *at*?

11. *W*hat *is* *the* *name* *of* *the* *chair* *you* *sit* *in*?

12. *W*hat *is* *the* *name* *of* *the* *sofa* *you* *lie* *on*?

13. *W*hat *is* *the* *name* *of* *the* *car* *you* *drive*?

14. *W*hat *is* *the* *name* *of* *the* *boat* *you* *sail*?

15. *W*hat *is* *the* *name* *of* *the* *plane* *you* *fly*?

1. *Thickened*

2. *Thickened*

3. *Thickened*

4. *Thickened*

5. *Thickened*

6. *Thickened*

7. *Thickened*

8. *Thickened*

9. *Thickened*

10. *Thickened*

Method of Moments Estimation (MME)

The basic idea of MME

is to use the law of large numbers to estimate parameters

of a distribution

$$\bar{X}_n \xrightarrow{\text{a.s.}} \mu$$

- Assume $\underline{\theta} = (\theta_1, \dots, \theta_k) \in \mathbb{R}^k$
as parameters

- Denote $\mu_j = E[X^j | \underline{\theta}]$
 $j = 1, \dots, k$

(Assuming that all moments exist)

- Obviously $\mu_j = \mu_j(\theta_1, \dots, \theta_k)$
- If the inverse functions exist

$$\theta_j = h_j(\mu_1, \dots, \mu_k)$$

Denote $\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n x_i^j$

By SLNN $\hat{\mu}_j \xrightarrow{\text{a.s.}} \mu_j$

as $n \rightarrow \infty$

The MME of $\underline{\theta}$ is then:

$$\hat{\theta}_j = h_j(\hat{\mu}_1, \dots, \hat{\mu}_k)$$

Theorem: If h_j is continuous,

then $\hat{\theta}_j \xrightarrow{a.s.} \theta_j$ as $n \rightarrow \infty$

Example: Assume that $X \sim U(0, \theta)$.

Find the MME of θ .

Lemma
~~Assume~~: For i.i.d r.v.'s that have finite first and second moments are estimated as

$$\hat{\mu}_1 = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}_n$$

$$\hat{\mu}_2 = \frac{1}{n} \sum_{i=1}^n x_i^2$$

■ the following holds.

$$\hat{\mu}_2 = \frac{1}{n} \sum_{i=1}^n x_i^2 = \frac{1}{n} \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) + \bar{x}_n^2$$

Proof:

Remark: We use the following notation

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

and

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Therefore, the results

of the previous lemma

can be written as:

$$\hat{\mu}_2 = \frac{n-1}{n} S^2 + \bar{X}_n^2 = \hat{\sigma}^2 + \bar{X}_n^2$$

In fact S^2 and $\hat{\sigma}^2$ are estimates of $\text{Var}(X)$.

Exercise: Assume that X_i 's

are i.i.d r.v.'s with

$$E[X] = \mu \text{ and } \text{Var}(X) = \sigma^2.$$

Calculate $E[S^2]$ and $E[\hat{\sigma}^2]$.

Which one is an unbiased

estimate of σ^2 ?

Example: Assume that

$X \sim N(\mu, \sigma^2)$. Find the MME

of μ and σ^2 .

Exercise : Are S_x^2 and $\hat{\sigma}^2$

consistent estimators of σ^2 ?

Assume that X_1, X_2, \dots are i.i.d

and $\sigma_x^2 < \infty$ and $E[X^4] < \infty$

Maximum Likelihood

Estimation (MLE)

Assume that X_1, \dots, X_n are

i.i.d r.v.'s with $X_i \sim f(x|\theta)$.

Then the joint pdf/pmf of

X_i 's is

$$f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n | \theta) =$$

We also call the joint

distribution the 'likelihood function'

$$L(\theta | \underline{x}) = L(\theta | x_1, \dots, x_n)$$

↑ ↑
observations

=



Likelihood function evaluated

at x_1, \dots, x_n , the observations

Example:

Intuition: Toss a biased coin

with $p = P(\{H\})$. We know

that $p \in \{.2, .5\}$. If the coin

is tossed 10 times and

9 heads are observed, what

is your estimate of p ?

$\hat{p} = .2$ or $\hat{p} = .5$? Why?

$L(p = .5 | \text{observation})$

$$L(p=2 | \text{observation})$$

The above example presents
a motivation for MLE.

Def) $\hat{\theta}_{MLE} = \operatorname{argmax} L(\theta | \underline{x})$

That is $L(\hat{\theta}_{MLE} | \underline{x})$

$$= \max_{\theta \in \Theta} L(\theta | \underline{x})$$

In practice, it is usually more efficient to consider the "log-likelihood" function

$$\ell(\theta | \underline{x}) = \log_e L(\theta | \underline{x})$$

=

because \log_e is an increasing

function, we can also say

$$\hat{\theta}_{MLE} = \operatorname{argmax} \ell(\theta | \underline{x})$$

Remark: The log likelihood

contains sum of individual likelihoods,

and it is often easier to

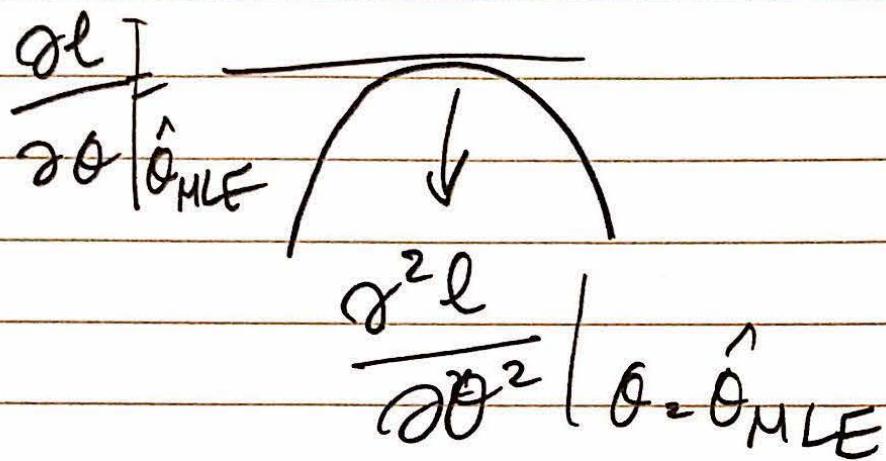
work with.

When l is differentiable with

respect to θ we should have

$$\frac{\partial l}{\partial \theta} \Big|_{\theta = \hat{\theta}_{MLE}} =$$

And to make sure $\hat{\theta}_{MLE}$ maximizes the \log likelihood function we should have



Example: Assume that $X \sim \text{Ber}(p)$.

Find \hat{p}_{MLE}

Example: Assume that

$X \sim N(0, \sigma^2)$. Find the
MLE of σ^2 .

Example: Assume that

$X \sim U(0, \theta)$. Find the Maximum Likelihood Estimate of θ

Properties of MLEs:

(1) Consistency

$$\hat{\theta}_{MLE} \xrightarrow{i.p.} \theta$$

2) Asymptotic Normality

$$\hat{\theta}_{MLE} \xrightarrow{d} N(\theta, \frac{1}{n} I(\theta))$$

where $I(\theta)$ is called the Fisher Information for θ .

$$(3) \quad \hat{g}(\theta)_{MLE} = g(\hat{\theta}_{MLE})$$

In general, when $\underline{\theta}$ is a

parameter vector e.g. $\underline{\theta} = [\theta_1, \dots, \theta_K]^T$

the gradient vector should

be zero at $\hat{\theta}_{MLE}$:

$$\nabla_{\underline{\theta}} l(\underline{\theta}) \Big|_{\underline{\theta} = \hat{\underline{\theta}}_{MLE}} = 0$$

where $\nabla_{\underline{\theta}} l(\underline{\theta}) = \begin{bmatrix} \frac{\partial l}{\partial \theta_1} \\ \vdots \\ \frac{\partial l}{\partial \theta_k} \end{bmatrix} = 0$

Also, the second order derivative conditions ~~conditions~~ are stated in terms of the Hessian Matrix:

$$\nabla_{\underline{\theta}}^2 l(\underline{\theta}) = [h_{ij}]$$

where

$$h_{ij} = \frac{\partial^2 l(\underline{\theta})}{\partial \theta_i \cdot \partial \theta_j}$$

For ~~the~~ $\hat{\underline{\theta}}_{MLE}$ to maximize

$l(\underline{\theta})$, the Hessian Matrix has

to be ~~positive~~ negative definite, i.e.

$$\nabla_{\underline{\theta}}^2 l(\underline{\theta}) \Big|_{\underline{\theta} = \hat{\underline{\theta}}_{MLE}} < 0$$

Recall that a matrix
~~x~~ symmetric

is negative definite if all
of its eigenvalues ~~are~~
strictly negative ~~all~~

Definition (Fisher Information)

If $X \sim f(x|\theta)$, and X_1, \dots, X_n are
i.i.d r.v.'s whose distribution
is $f(x|\theta)$, the Fisher Information
for θ is defined as

$$I(\theta) = \mathbb{E}_{X|\theta} \left[\left(\frac{\partial \log f(X|\theta)}{\partial \theta} \right)^2 \right]$$

$$= -\mathbb{E}_{X|\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta) \right]$$

Example: Assume that

$X \sim \text{Ber}(p)$. Calculate the Fisher

Information $I(p)$.

Example: We saw that the MLE is asymptotically normal,
 i.e. $\sqrt{n}(\hat{\theta}_n - \theta)$ ~~\xrightarrow{d}~~ $\sim N\left(0, \frac{1}{nI(\theta)}\right)$.

Using this result, find the asymptotic distribution of the MLE of p for $\text{Ber}(p)$ and compare the results to those of the CLT.

A blank page featuring ten horizontal ruling lines. The lines are evenly spaced and extend across the width of the page. The entire page is enclosed within a thick red rectangular border.

A blank page featuring ten horizontal ruling lines. The lines are evenly spaced and extend across the width of the page. The entire page is enclosed within a thick red rectangular border.

Question : If $\hat{\theta}_n$ is an unbiased estimator of θ , is it possible to improve the precision (= variance) of $\hat{\theta}$ with adding more samples?

Theorem: (The Cramér-Rao Bound)

Assume that $\hat{\theta}_n$ is an unbiased estimator of θ . Then

$$\text{Var}(\hat{\theta}_n) \geq \frac{1}{n I(\theta)}$$

Remark: The Cramér-Rao bound is a lower bound for the variance of any unbiased estimator $\hat{\theta}_n$.

Maximum A-Posteriori

Estimate (MAP)

This estimator makes

the 'Bayesian Assumption'

that the unknown parameter

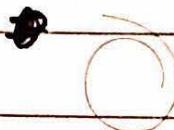
θ is also a r.v. with
pdf $h(\theta)$, which is called
the "prior".

Consequently, the best
estimate for θ is the one

that maximizes the
"posterior pdf" $\hat{f}(\theta | \underline{x})$,

i.e.

$$\hat{\theta}_{MAP} = \operatorname{argmax} (\hat{f}(\theta | \underline{x}))$$



However

$$g(\theta|x) = \frac{f(x|\theta) h(\theta)}{\int f(x|\theta') h(\theta') d\theta'}$$

the denominator is a constant,

$$\text{therefore } g(\theta|x) \propto f(x|\theta) h(\theta)$$

and

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{arg\,max}}$$

Note that

$$(i) \hat{\theta}_{MAP} = \hat{\theta}_{MLE} \text{ when }$$

$$h(\theta) = \text{constant} \propto 1$$

which means MLE is the same as MAP when the prior ~~estimate~~ pdf of θ is uniform.

(ii) The prior pdf encodes subjective knowledge about θ .

(iii) MLE and MAP heavily rely on the parametric form of ~~of~~ $f(x|\theta)$.

In practice, $f(x|\theta)$ may be difficult to find or too complicated to optimize.

This often happens when X is a complicated / Corrupted

function/version of a simple r.v.

i.e., i.e. $X = T(Z)$.

To calculate MLEs in

such situations, the

celebrated Expectation

Maximization (EM) Algorithm
is used.

Example: Assume that X_1, \dots, X_n
are i.i.d $N(\mu, 1)$ r.v/s and
 $\mu \sim N(\mu_0, 1)$, where μ_0 is a constant
What is the MAP estimate
of μ ?

Exercise: Show that the MLE

of μ when $X \sim N(\mu, 1)$ is

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

Methods that assume a particular form for the distribution of a r.v. and try to estimate the parameters of that particular distribution

from a "random sample,"
i.e. i.i.d. r.v.'s X_1, \dots, X_n
can be called parametric methods.

Here, we focus on the so-called non-parametric methods, that are generic and do not assume parametric forms for probability

distributions.

Histograms

Given i.i.d X_1, \dots, X_n , assume

that we select m intervals called

bins, $[e_j, e_{j+1})$ where

$$e_1 < \dots < e_{m+1} \quad (\text{edge sequence})$$

where $e_1 \leq \min_i X_i$

$$e_{m+1} \geq \max_i X_i$$

When $e_{m+1} = \max_i X_i$, we use

the bin $[e_m, e_{m+1}]$ instead of

$[e_m, e_{m+1})$ so that no data
is lost.

The histogram count for

the j^{th} bin $[e_j, e_{j+1})$ is

$$H_j = \sum_{i=1}^n I_{[e_j, e_{j+1})}(x_i)$$

which is the number of samples
that reside in $[e_j, e_{j+1})$.

Observe that $I_{[e_j, e_{j+1})}(x_i)$

is a Bernoulli r.v. The

probability of success \uparrow is

calculated as

$$E[I_{[e_j, e_{j+1})}(x_i)] = P(e_j \leq x_i < e_{j+1})$$

Since X_i are i.i.d., $I_{[e_j, e_{j+1}]}(X_i)$ are also i.i.d for each j .

The sample mean of

$I_{[e_j, e_{j+1}]}(X_i)$ is

$$\frac{1}{n} \sum_{i=1}^n I_{[e_j, e_{j+1}]}(X_i) = \frac{H_j}{n}$$

For large n , $\frac{H_j}{n} \approx E[I_{[e_j, e_{j+1}]}(X_i)]$

according to the LLN, so

$$P(e_j \leq X_i \leq e_{j+1}) \approx \frac{H_j}{n}$$

which is a straightforward way

to estimate the distribution of X in $[e_j, e_{j+1}]$ by a constant.

If X_i 's are integer valued,

it is convenient to use $e_j = j - \frac{1}{2}$

and $e_{j+1} = j + \frac{1}{2}$, i.e. to center

the bins on the integers,

Then $\frac{H_j}{n} \approx P(X_i=j)$

If we are planning to draw

a density function over the histogram;

$$\frac{H_j}{n} \approx P(e_j \leq X_i \leq e_{j+1}) = \int_{e_j}^{e_{j+1}} f_X(x) dx$$

$$\approx f_X(c_j) \Delta x_j \quad \text{where } c_j = \frac{e_j + e_{j+1}}{2}$$

$$\Rightarrow f_X(c_j) = \frac{H_j}{n \Delta x_j}$$

Kernel Density Estimation

The empirical distribution places

the probability mass $\frac{1}{n}$ to each $X_i = x_i$ that is observed.

In order to have a smooth

estimate of the pdf of X , one

can imagine distributing it in

a neighborhood of $X_i = x_i$, according

to some pdf. This approach

is called Kernel Density Estimation (KDE).

KDE:

1. Choose a pdf K , the kernel.

Typically, K is a symmetric pdf centered at origin. $N(0,1)$ and $U[-0.5, 0.5]$ ~~are~~, Triangle $[-1, 1]$ are

common choices.

2. At each x_i , center a copy of the kernel

$$\frac{1}{h} K\left(\frac{x-x_i}{h}\right)$$

This will control the distribution

of the $\frac{1}{n}$ probability assigned to x_i . h is called the smoothing parameter, the window width, or the bandwidth.

3. The KDE estimate of f_X is

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x-x_i}{h}\right)$$

$$= \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$$

Remark: The R function density performs KDE.
Remark: The difficult choice

in KDE is choosing h .

Empirical Distribution

Definition: Assume that X_1, X_2, \dots, X_n are i.i.d. samples with common CDF F_X . The empirical distribution function of X_1, X_2, \dots, X_n is

defined as

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[X_i, +\infty)}(x)$$

where I_A is the indicator function of set A.

Informally, the empirical Cdf

is increased by $\frac{1}{n}$ at each

X_i .

Goodness -of- fit

To measure how well a fitted distribution resembles

the sample data, we can

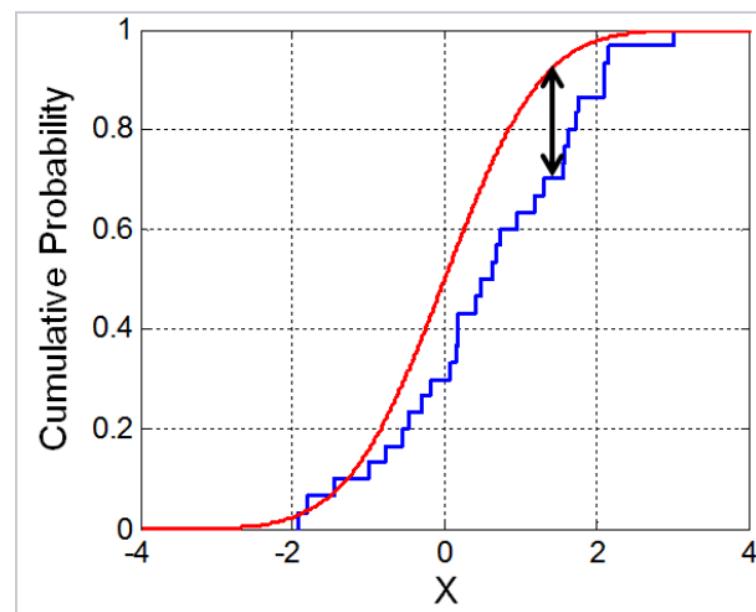
use the Kolmogorov-Smirnov (K-S) test statistic :

$$D_n = \sup_x | \hat{F}_n(x) - F(x) |$$

empirical
distribution

fitted
distribution

$$D_n = \sup_x |F_n(x) - F(x)|$$



The Golivenko - Cantelli Theorem

Theorem: Assume that X_1, X_2, \dots, X_n are i.i.d. samples with common cdf $F_X(x)$ and the empirical cdf $\hat{F}_n(x)$. Then the

Kolmogorov-Smirnov Statistic

$$D_n = \sup_x |\hat{F}_n(x) - F(x)|$$

Converges to Zero almost surely

for all x , i.e.

$$\text{P}(\lim_{n \rightarrow \infty} D_n \neq 0) = 0.$$

Remark: $\hat{F}_n(x)$ is a random variable at each x . D_n is also a r.v.

Remark: The G-K Theorem is sometimes called the

Fundamental theorem of statistics, and is a basis for the Kolmogorov-Smirnov test.

Assume that multiple distributions are estimated from i.i.d samples using different (parametric or non-parametric) methods. Which distribution is "the best?"

The Kolmogorov-Smirnov Test Criterion says that the KS statistic has to be computed for each distribution and the distribution with the minimum D_n has to be selected.

Goodness of fit as hypothesis

Testing

H_0 : All samples x_1, \dots, x_n

Come from the cdf F

H_1 : At least one of the

samples x_1, \dots, x_n does not

Come from the cdf F

Under the null hypothesis

(assuming H_0 is true)

the K-S test statistic

D_n follows a K-S

distribution, which has

a complicated form.

Therefore, to reject the null,

we have to observe a

sufficiently large D_n so that

we are in the rejection

region.

K-S distribution has a complicated form.

$$P\left(D_n < \frac{1}{2n} + v\right) = n! \int_{\frac{1}{2n}-v}^{\frac{1}{2n}+v} \int_{\frac{3}{2n}-v}^{\frac{3}{2n}+v} \cdots \int_{\frac{2n-1}{2n}-v}^{\frac{2n-1}{2n}+v} g(u_1, u_2, \dots, u_n) du_n \dots du_2 du_1 \quad 0 \leq v \leq \frac{2n-1}{2n}$$

$$g(u_1, u_2, \dots, u_n) = 1 \text{ over } 0 \leq u_1 \leq u_2 \leq \dots \leq u_n \leq 1,$$

And 0 otherwise

Birnbaum, Z.W. (1952), "Numerical Tabulation of the Distribution of Kolmogorov's Statistic for Finite Sample Size," Journal of the American Statistical Association, 47, 425–441

$$F_{D_6}(t) = \begin{cases} 0 & t < \frac{1}{12} \\ 46080t^6 - 23040t^5 + 4800t^4 - \frac{1600}{3}t^3 + \frac{100}{3}t^2 - \frac{10}{9}t + \frac{5}{324} & \frac{1}{12} \leq t < \frac{1}{6} \\ 2880t^6 - 4800t^5 + 2360t^4 - \frac{1280}{3}t^3 + \frac{235}{9}t^2 + \frac{10}{27}t - \frac{5}{81} & \frac{1}{6} \leq t < \frac{1}{4} \\ 320t^6 + 320t^5 - \frac{2600}{3}t^4 + \frac{4240}{9}t^3 - \frac{785}{9}t^2 + \frac{145}{27}t - \frac{35}{1296} & \frac{1}{4} \leq t < \frac{1}{3} \\ -280t^6 + 560t^5 - \frac{1115}{3}t^4 + \frac{515}{9}t^3 + \frac{1525}{54}t^2 - \frac{565}{81}t + \frac{5}{16} & \frac{1}{3} \leq t < \frac{5}{12} \\ 104t^6 - 240t^5 + 295t^4 - \frac{1985}{9}t^3 + \frac{775}{9}t^2 - \frac{7645}{648}t + \frac{5}{16} & \frac{5}{12} \leq t < \frac{1}{2} \\ -20t^6 + 32t^5 - \frac{185}{9}t^3 + \frac{175}{36}t^2 + \frac{3371}{648}t - 1 & \frac{1}{2} \leq t < \frac{2}{3} \\ 10t^6 - 38t^5 + \frac{160}{3}t^4 - \frac{265}{9}t^3 - \frac{115}{108}t^2 + \frac{4651}{648}t - 1 & \frac{2}{3} \leq t < \frac{5}{6} \\ -2t^6 + 12t^5 - 30t^4 + 40t^3 - 30t^2 + 12t - 1 & \frac{5}{6} \leq t < 1 \\ 1 & t \geq 1. \end{cases}$$



Fortunately, there are tables

for the critical values of

the distribution of K-S

test.

K-S is distribution free, i.e. it does not assume any particular distribution for the data.

Drawback: Only for one dimensional distributions

There are other goodness of fit tests including the Anderson-Darling test.

The Chi-Squared Goodness of Fit Test

This test is based on the histogram of data. The idea is that if the pmf or

pdf ~~of~~ estimated for a sample is too different from its histogram, it probably is not a good fit for the data.

In particular, if

$P(e_j \leq X_i \leq e_{j+1})$ calculated

according to the estimated

pmf/pdf is not close to

the histogram estimate

of $P(e_j \leq X_i \leq e_{j+1})$,

which is $\frac{h_j}{n}$, the

estimated pmf/pdf is

not a good fit.

It can be shown that

the quantity

$$\chi^2 = \sum_{j=1}^m \left(\frac{|H_j - np_j|^2}{np_j} \right)$$

has a Chi-squared distribution

with $m-1$ degrees of freedom

when n is large.

This is because

$$H_j - np_j$$

$$\sqrt{np_j}$$

has a standard normal distribution

when n is large.

Therefore, if

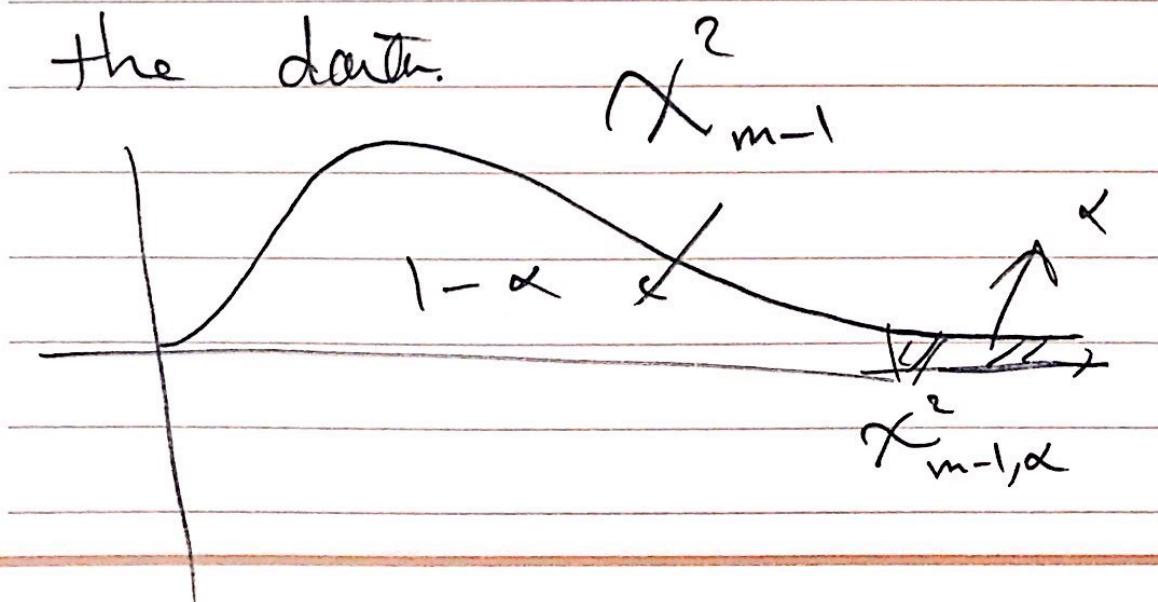
$$\chi^2 > \chi^2_{m-1, \alpha}$$

(i.e. if the discrepancy is too large), we reject the null

that the pdf approx we

found is a good fit for

the data.



The Bootstrap

Assume that we have an i.i.d sample X_1, X_2, \dots, X_n and a statistic $T(X_1, X_2, \dots, X_n)$. If we have access to the common

distribution of X_i 's, we can derive the distribution of T . However, it is not usually the case. If we have the luxury of having multiple samples from the common

distribution, we can still compute $T(\cdot)$ from each sample to have a sample of $T(\cdot)$'s and somehow estimate the distribution of $T(\cdot)$, e.g. by KDE.

Sample 1: $X_1^{(1)} \quad X_2^{(1)} \dots \quad X_n^{(1)} \rightarrow T^{(1)}$

Sample 2: $X_1^{(2)} \quad X_2^{(2)} \dots \quad X_n^{(2)} \rightarrow T^{(2)}$

⋮
⋮
Sample m: $X_1^{(m)} \quad X_2^{(m)} \dots \quad X_n^{(m)} \rightarrow T^{(m)}$

Estimate F_T from $T^{(1)} \dots T^{(m)}$

Unfortunately, having multiple samples
is not always possible.

Therefore, we use a "resampling"
method called The Bootstrap.

The basic idea is to randomly

draw datasets from the sample
we have with replacement.

This is done B times, producing
 B bootstrap datasets. Then we
calculate $T(\cdot)$ for each of

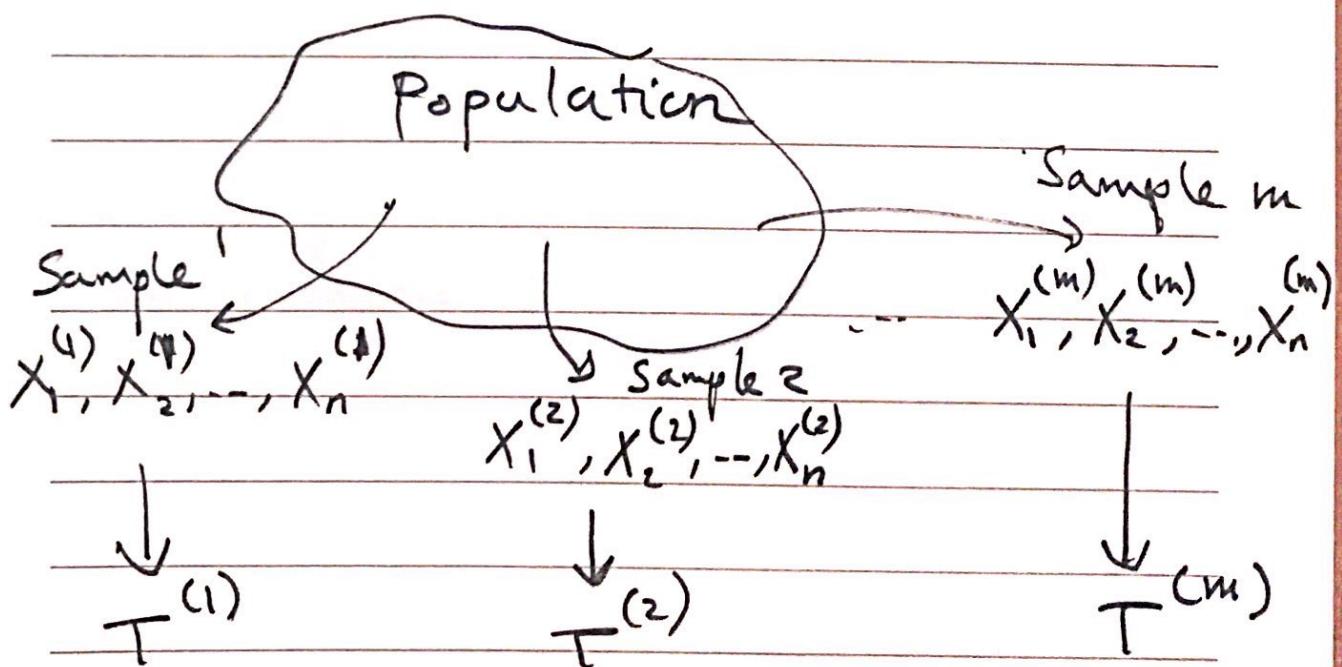
the B data sets and obtain
a bootstrap sample of $T(\cdot)$.
Then we can estimate the
distribution or statistical properties
of $T(\cdot)$ (such as mean, variance etc.)

from the bootstrap sample.

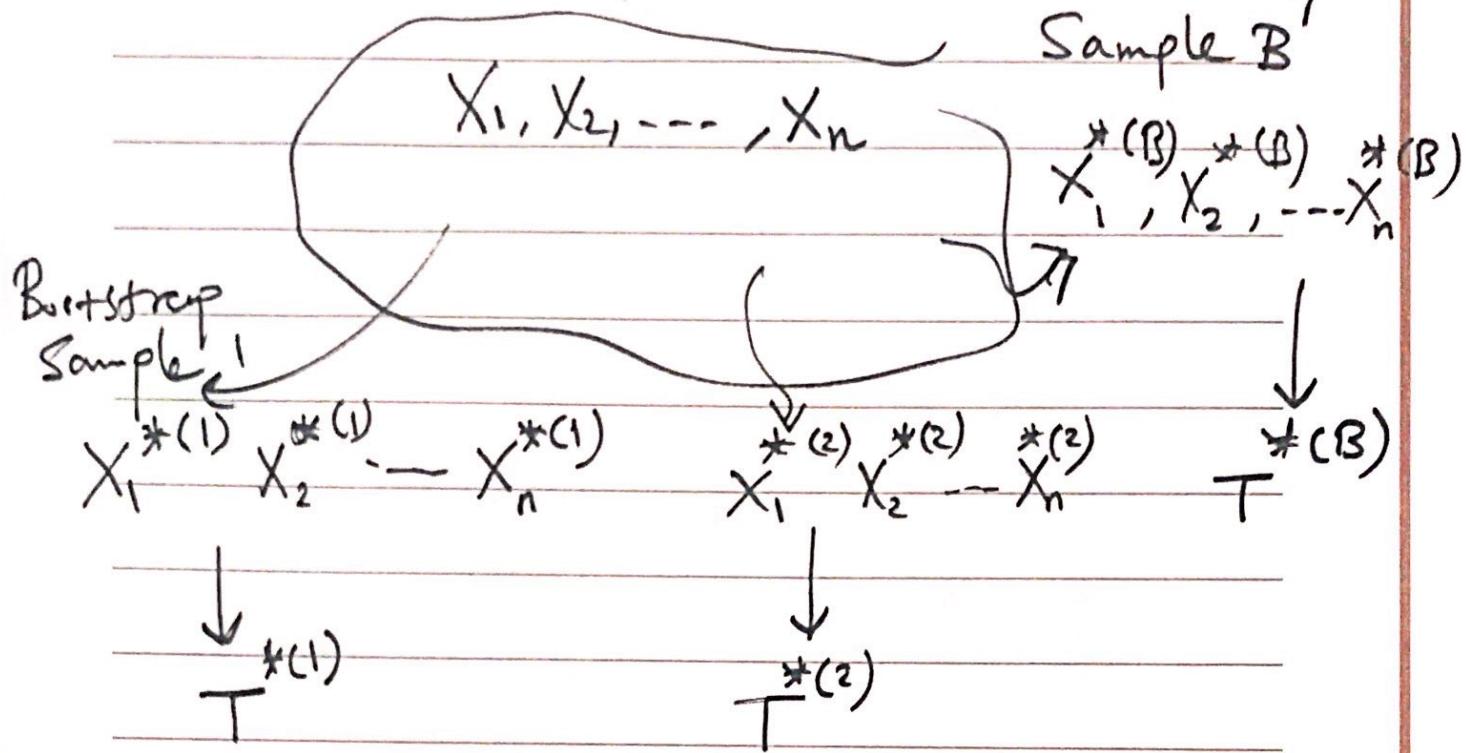
Remark: Sampling with replacement
from the i.i.d Sample X_1, \dots, X_n
is equivalent to sampling from
the empirical distribution function
of X_1, X_2, \dots, X_n .

Remark: One can view the original sample, from which bootstrap samples are drawn, the new population, and bootstrap samples as the samples drawn from that population

Real World:



Bootstrap Work



In bootstrap world, the statistics

of the original sample are

considered as population

parameters.

Example: Calculating the estimated variance of a statistic using bootstrap.

$$\text{Var}(T(X_1, X_2, \dots, X_n)) = \frac{1}{B-1} \sum_{b=1}^B (T(\tilde{X}_1^{*(b)}, \tilde{X}_2^{*(b)}, \dots, \tilde{X}_n^{*(b)}) - \bar{T}^*)^2$$

where

$$\bar{T}^* = \frac{1}{B} \sum_{b=1}^B T(\tilde{X}_1^{*(b)}, \tilde{X}_2^{*(b)}, \dots, \tilde{X}_n^{*(b)})$$

Bootstrap Percentile Confidence

Intervals

Assume the sample X_1, \dots, X_n was drawn from a population.

We wish to construct a $(1-\alpha)$

Confidence Interval for $T(X_1, \dots, X_n)$

(say \bar{X}_n or S_n) using bootstrap:

1. Resample R_1, R_2, \dots, R_B from

X_1, X_2, \dots, X_n with replacement

$$R_i = \{X_1^{*(i)}, X_2^{*(i)}, \dots, X_n^{*(i)}\} \quad i=1, 2, \dots, B$$

Note: $X_j^{*(i)}$ is may be repetitive

2. Calculate

$$T^{*(i)} = T(X_1^{*(i)}, X_2^{*(i)}, \dots, X_n^{*(i)})$$

$$i = 1, 2, \dots, B$$

for each bootstrap sample R_i .

3. Order $T^{*(i)}$'s and call the

ordered version $U_{(1)}, U_{(2)}, \dots, U_{(B)}$

$$U_{(1)} \leq U_{(2)} \leq \dots \leq U_{(B)}$$

The $1-\alpha$ C.I. is

$$[U_{(a)}, U_{(b)}]$$

$$\text{where } a = \lfloor \frac{\alpha}{2} B \rfloor$$

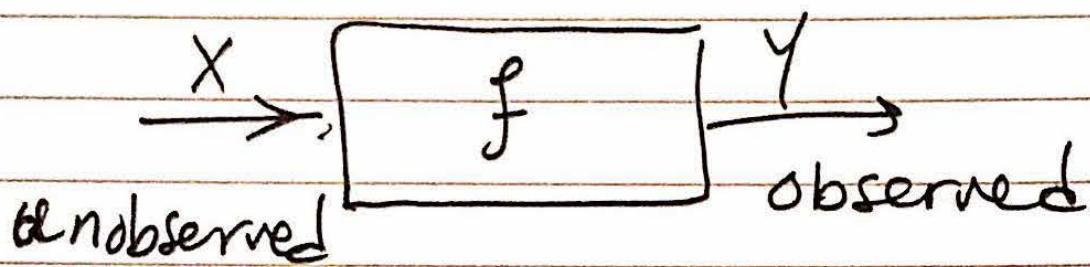
$$\text{and } b = \lfloor (1 - \frac{\alpha}{2}) B \rfloor$$

$$B=1000, \alpha=5\% \rightarrow a=25 \quad b=975$$

Minimum Mean-Squared

Error Estimation

(MMSE)



Find an estimate for X

using Y , $\hat{X} = g(Y)$, such

that $g(Y)$ minimizes the

average mean-square

estimation error $E[(\hat{X} - g(Y))^2]$

You may recall that
such an estimate is
the conditional expectation

$$\hat{X} =$$

[The estimation error is
orthogonal to Y , i.e. no
further information about Y
is in the error.

Linear Estimation of Random Vectors

Question: Given two random vectors $\underline{X}, \underline{Y}$, where \underline{X} is not observed but \underline{Y} is observed

what is the best estimate of \underline{X} in terms of a linear (affine) function of \underline{Y} ?

This needs clarification.

What does best mean?

The estimate with minimum mean-squared error,

$$\text{MSE} = \mathbb{E} [\|\underline{x} - A\underline{y} - b\|_2^2]$$

Let $\mathbb{E}[x] = m_x$ and $\mathbb{E}[y] = m_y$

Then

$$\begin{aligned} \text{MSE} &= \mathbb{E} \left[\underbrace{\|\underline{x} - m_x - A(\underline{y} - m_y)\|_2^2}_{+ m_x - A m_y - b \|_2^2} \right] \\ &\stackrel{?}{=} \end{aligned}$$

$$\begin{aligned} \text{MSE} &= \mathbb{E} [(\underline{z} + \underline{g})^T (\underline{z} + \underline{g})] \\ &= \mathbb{E} [\underline{z}^T \underline{z} + \underline{z}^T \underline{g} + \underline{g}^T \underline{z} + \underline{g}^T \underline{g}] \end{aligned}$$

$$= E \left[\| (\underline{X} - \underline{m}_x) - A (\underline{Y} - \underline{m}_y) \|^2 \right] \\ + E \left[\| \underline{m}_x - A \underline{m}_y - \underline{b} \|^2 \right]$$

Obviously, to minimize MSE,
 \underline{b} must be

Now, let's estimate \hat{A} that
minimizes

$$\begin{aligned} & \mathbb{E} [\|(\underline{X} - \underline{m}_X) - A(\underline{Y} - \underline{m}_Y)\|^2] \\ &= \mathbb{E} [\|\underline{\tilde{X}} - A\underline{\tilde{Y}}\|^2] \end{aligned}$$

This looks like the MMSE

problem, except $g(\underline{\tilde{Y}})$ is constrained
to be $A\underline{\tilde{Y}}$.

We use the Principle of
Orthogonality to solve this
problem.

The error $\tilde{X} - \tilde{A}\tilde{Y}$ has to be orthogonal to all linear functions of \tilde{Y} , $g(\tilde{Y}) = \tilde{B}\tilde{Y}$:

$$E[] =$$

To find the A that ~~maximizes~~
minimizes the error, let's apply
the function tr to both sides

$$\text{tr}(E[(\tilde{B}\tilde{Y})^T (\tilde{X} - \tilde{A}\tilde{Y})]) = 0$$

$\text{tr}(\cdot)$ commutes with $E[\cdot]$ (why?)

So,

$$E[\text{tr}((B\tilde{Y})^T (\tilde{X} - A\tilde{Y}))] = 0$$

on the other hand $\text{tr}(CD) = \text{tr}(DC)$

so

$$E[\text{tr}((\tilde{X} - A\tilde{Y})(B\tilde{Y})^T)] = 0$$

changing the order of $E[\cdot]$ and
 $\text{tr}(\cdot)$ again

which means A solves

$$AC_y = C_{xy}$$

and if C_y is non-singular

$$A = C_{xy}C_y^{-1}$$

Therefore the LMME of X is

terms of \underline{Y} is

$$C_{xy} C_y^{-1} (\underline{Y} - \underline{m}_y) + \underline{m}_x$$

Note that when \underline{X} and \underline{Y}

are jointly Gaussian, LMMSE estimate and MMSE coincide!

KNN Estimation of Conditional Expectation

We learned that $E[Y|X]$ is the MMSE estimate of Y given X . How to estimate it from

data?

Assume that we have n samples of Y , given $X=x$, and they are (conditionally) ~~independent~~ i.i.d. Then $E[Y|X=x]$ can be

approximated using LLNs:

$$\mathbb{E}[Y|X=x] \approx \frac{1}{n} \sum_{i=1}^n Y_i | X=x$$

However, in real situations, we may not have enough samples for estimating Y at $X=x$.

Therefore, we relax the above approximation to k other x 's that are "close" to x :

$$\mathbb{E}[Y|X=x] \approx \cancel{\frac{1}{n} \sum_{i=1}^n Y_i}$$

$$\frac{1}{n} \sum_{i=1}^n Y_i | X=x_1 \cup X=x_2 \cup \dots \cup X=x_k$$

Graphically

