

This Week :

Dr. Osonde Osoba

- Confidence Intervals (C-I)
 - Pollster's C-I
 - Hypothesis Testing
- Estimation
 - Properties: Bias, MSE, Consistency
 - Method of moments
 - Maximum Likelihood Estimation (MLE)
 - Maximum A Posteriori (MAP)
 - Expectation - Maximization (EM)

- Maximum Entropy (Max Ent)

- Minimum Mean-Squared Error (MMSE)

H/W:

L-G: 8.24 (ab), 8.28 (ab), 8.29 (a), 8.30 (a)

4.147 - 4.149

8.39 - 8.41

G: 6.27, 6.28

Last week, we covered CLT:

$$\{X_i\}_{i=1}^n : \text{iid} \quad \sigma_x^2 < \infty$$

$$Z_n = \text{std}(\bar{X}_n)$$

$$\Rightarrow Z_n \xrightarrow{d} Z \sim N(0,1)$$

Qu 1: How good is the CLT CDF approximation for a given sample size, n ?

i.e characterize $\sup_z |F_n(z) - \Phi(z)|$ as $n \rightarrow \infty$

Ans: Berry - Essen Theorem

$$\exists c > 0 \quad \forall z \quad \forall n \quad \text{when } E[|X|^3] < \infty$$

such that

$$\boxed{|F_n(z) - \Phi(z)| \leq \frac{c}{\sqrt{n}} \cdot \frac{E[|X - \mu|^3]}{\sigma^3}}$$

Convergence rate $\propto 1/\sqrt{n}$

where $F_n(z)$ is the cdf of $Z_n = \text{std}(\bar{X}_n)$.

$$c \in [0.4097, 0.71]$$

[$\Phi(z)$ is the standard normal CDF]

11-3

Qu 2: Can we use the convergence of $Z_n = \text{std}(\bar{X}_n)$ to infer a probable range of values for μ ?

i.e find $I = [a, b]$ such that $P(\mu \in I) = (1 - \alpha)$.

Ans:

Consider the statement

$$P(Z_n \in I) = 1 - \alpha$$

$n \rightarrow \infty \Rightarrow$ By CLT we can use the CDF $\Phi(z)$ to estimate this probability. Φ is symmetric around 0.

So we expect $I = [-z_{\alpha/2}, z_{\alpha/2}]$ where the value of $z_{\alpha/2}$ depends on α and comes from analyzing $\Phi(z)$

$$1 - \alpha = P(Z_n \in [-z_{\alpha/2}, z_{\alpha/2}])$$

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

$$= P(-z_{\alpha/2} \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2})$$

$$= P(-z_{\alpha/2} \cdot \sigma/\sqrt{n} \leq \bar{X}_n - \mu \leq z_{\alpha/2} \cdot \sigma/\sqrt{n})$$

$$= P(\bar{X}_n - z_{\alpha/2} \cdot \sigma/\sqrt{n} \leq \mu \leq \bar{X}_n + z_{\alpha/2} \cdot \sigma/\sqrt{n})$$

$$1 - \alpha = P\left(\mu \in \left[\bar{X}_n - \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}, \bar{X}_n + \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}\right]\right)$$

So our range

$$I = \left[\bar{X}_n - \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}, \bar{X}_n + \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}\right]$$

I is called the $(1-\alpha)100\%$ Confidence Interval for the mean μ .

μ is a constant and \bar{X}_n is a r.v. So the statement $P(\mu \in I) = 1-\alpha$ does not have the typical "measure of preimage" meaning. since μ is not a measurable function. The randomness is in I .

Confidence Intervals are an example of a random set.

$$\text{i.e. } I: \Omega \rightarrow S \subset 2^{\mathbb{R}}$$

In general, the C-I of an estimate $\hat{\theta}_n$ for θ is:

$$\hat{\theta}_n \pm E$$

where E is the margin of error (M.O.E).

e.g. \bar{X}_n estimating μ :

$$\bar{X}_n \pm E$$

$$E = \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}} \quad \left(\sqrt{\text{Var}(\bar{X}_n)} = \text{standard error} \right)$$

(**)

\Rightarrow To get a C-I for μ at the $(1-\alpha)100\%$ level

with E M.O.E. we need

$$n = \frac{(z_{\alpha/2})^2 \cdot \sigma_x^2}{E^2}$$

Pollster's C-I:

e.g. "63% of Americans believe '—'."

usually at 95% confidence level and $E = \pm 3$ percentage pt

This is really an attempt to estimate p in

$X_k \sim \text{Bernoulli}(p)$ where X_k encodes the binary ^{popn} behavior.

\Rightarrow The sample mean is:

$$\hat{p}_n = \frac{1}{n} \sum_{k=1}^n X_k$$

and for large n $\hat{p}_n \stackrel{d}{\approx} N(\mu_x, \sigma_x^2/n)$ by CLT.

$$\mu_x = p; \quad \sigma_x^2 = p(1-p)$$

[We have assume $\{X_k\}$ is iid and "representative"]

$\Rightarrow (1-\alpha)100\%$ C-I for p

$$\hat{p}_n \pm z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

Now $p \in [0,1] \Rightarrow 0 \leq p(1-p) \leq 1/4$

MOE = ± 3 pct pts

$$\alpha = 5\% \Rightarrow \boxed{z_{\alpha/2} = 1.96}$$

$$\alpha = 1\% \Rightarrow \boxed{z_{\alpha/2} = 2.58}$$

$$n = \left(\frac{1.96}{(0.03)4} \right)^2$$

$$n = \left(\frac{2.58}{(0.03)4} \right)^2$$

$\Rightarrow n \geq 1068$ samples

$\Rightarrow n \geq 1849$ samples

Confidence Intervals are the basis for frequentist hypothesis testing (HT). The result of a HT depend on whether a test statistic falls inside or outside a $(1-\alpha)100\%$ C-I.

e.g. [a basic Hypothesis Test]: (for μ_x)

① Formulate competing hypothesis about μ_x

$$H_0: \mu_x = \mu_0 \quad [\text{null hypothesis}]$$

$$H_{alt}: \mu_x \neq \mu_0 \quad [\text{alternate hypothesis}]$$

② Specify $(1-\alpha)100\%$ C-I for $\bar{X} - \mu_0$:

$$(1-\alpha)100\% \text{ C-I} : \left[-z_{\alpha/2} \sigma / \sqrt{n}, z_{\alpha/2} \sigma / \sqrt{n} \right]$$

③ Calculate test statistic \bar{x} for observed data.

④ Reject H_0 if $\bar{x} \notin (1-\alpha)100\%$ C-I

"Do not reject" H_0 / accept H_{alt} if

$$\bar{x} \in (1-\alpha)100\% \text{ C-I}$$

C-Is for p , σ^2 , σ_1^2 / σ_2^2 determine other types of Hypothesis tests for other parameters.

Tests of Hypotheses

Hypotheses	Assumptions	Critical Region
$H_0: \mu = \mu_0$ $H_1: \mu > \mu_0$	$N(\mu, \sigma^2)$ or n large, σ^2 known	$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \geq z_\alpha$
$H_0: \mu = \mu_0$ $H_1: \mu > \mu_0$	$N(\mu, \sigma^2)$ σ^2 unknown	$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \geq t_\alpha(n-1)$
$H_0: \mu_x - \mu_y = 0$ $H_1: \mu_x - \mu_y > 0$	$N(\mu_x, \sigma_x^2)$ $N(\mu_y, \sigma_y^2)$ σ_x^2, σ_y^2 known	$z = \frac{\bar{x} - \bar{y} - 0}{\sqrt{\sigma_x^2/n + \sigma_y^2/m}} \geq z_\alpha$
$H_0: \mu_x - \mu_y = 0$ $H_1: \mu_x - \mu_y > 0$	Variances unknown, large samples	$z = \frac{\bar{x} - \bar{y} - 0}{\sqrt{s_x^2/n + s_y^2/m}} \geq z_\alpha$
$H_0: \mu_x - \mu_y = 0$ $H_1: \mu_x - \mu_y > 0$	$N(\mu_x, \sigma_x^2)$ $N(\mu_y, \sigma_y^2)$ $\sigma_x^2 = \sigma_y^2$, unknown	$t = \frac{\bar{x} - \bar{y} - 0}{s_p \sqrt{1/n + 1/m}} \geq t_\alpha(n+m-2)$ $s_p = \sqrt{\frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2}}$
$H_0: \mu_D = \mu_x - \mu_y = 0$ $H_1: \mu_D = \mu_x - \mu_y > 0$	X and Y normal, but dependent	$t = \frac{\bar{d} - 0}{s_d/\sqrt{n}} \geq t_\alpha(n-1)$
$H_0: \sigma^2 = \sigma_0^2$ $H_1: \sigma^2 > \sigma_0^2$	$N(\mu, \sigma^2)$	$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2} \geq \chi_\alpha^2(n-1)$
$H_0: \sigma_x^2/\sigma_y^2 = 1$ $H_1: \sigma_x^2/\sigma_y^2 > 1$	$N(\mu_x, \sigma_x^2)$ $N(\mu_y, \sigma_y^2)$	$F = \frac{s_x^2}{s_y^2} \geq F_\alpha(n-1, m-1)$
$H_0: p = p_0$ $H_1: p > p_0$	$b(n, p)$ n is large	$z = \frac{y/n - p_0}{\sqrt{p_0(1-p_0)/n}} \geq z_\alpha$
$H_0: p_1 - p_2 = 0$ $H_1: p_1 - p_2 > 0$	$b(n_1, p_1)$ $b(n_2, p_2)$	$z = \frac{y_1/n_1 - y_2/n_2 - 0}{\sqrt{\left(\frac{y_1 + y_2}{n_1 + n_2}\right)\left(1 - \frac{y_1 + y_2}{n_1 + n_2}\right)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \geq z_\alpha$

Estimation :

The basic idea here is to use random samples (i.e iid samples) to infer parameters/information about the underlying population. So we have

$$\{X_i\}_{i=1}^n : \text{iid samples}$$

and we define a statistic (a function of the data)

$$\hat{\theta}_n = g(X_1, X_2, \dots; X_n).$$

Suppose θ is an parameter of the population distribution.

The goal is to specify statistics $g(\cdot)$ that yield "optimal" estimators $\hat{\theta}_n$ for the estimand θ .

There are 3 basic tools for judging estimator "quality":

(a) Bias :

$$B(\hat{\theta}_n; \theta) = |E[\hat{\theta}_n] - \theta|$$

if $B(\hat{\theta}_n; \theta) = 0$ we say $\hat{\theta}_n$ is an unbiased estimate for θ .

If $B(\hat{\theta}_n; \theta) \neq 0$ but $B(\hat{\theta}_n; \theta) \rightarrow 0$ as $n \rightarrow \infty$ then $\hat{\theta}_n$ is asymptotically unbiased.

Else $\hat{\theta}_n$ is biased.

(b) Mean Squared-Error ($MSE(\hat{\theta}_n, \theta)$)

$$MSE(\hat{\theta}_n; \theta) = E[(\hat{\theta}_n - \theta)^2]$$

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

[Variance - Bias decomposition for MSE]

(c) Consistency:

$\hat{\theta}_n$ is a consistent estimate for θ

$$\Leftrightarrow \hat{\theta}_n \xrightarrow{P} \theta$$

$$\Leftrightarrow \forall \varepsilon > 0 \quad \lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| > \varepsilon) = 0$$

Maximum Likelihood Estimation (MLE)

Assuming iid $\{X_i\}_{i=1}^n$ with common pdf $f(x; \theta)$ parametrized by an unknown variable θ .

The goal of MLE is to use $\{X_i\}_{i=1}^n$ to estimate

θ . The MLE approach is to find the value of θ that maximizes the probability of observing

the sample $\vec{x} = \{x_i\}_{i=1}^n$. This requires the

idea of a likelihood $l(\theta) = f_{\vec{x}}(\vec{x}; \theta)$. This

is the pdf of \vec{X} considered as a function of

θ instead of \vec{x} .

$$\hat{\theta}^{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \{ f_{\vec{x}}(\vec{x}; \theta) \}$$

$$\vec{X} = \{X_i\}_{i=1}^n \quad \text{iid}$$

$$\Leftrightarrow \hat{\theta}^{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \left\{ \prod_{i=1}^n f_x(x_i; \theta) \right\}$$

$$\hat{\theta}^{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \left\{ \prod_{i=1}^n L(\theta) \right\}$$

The $(L(\theta))$ log-likelihood is often easier to optimize

$$L(\theta; \vec{x}) = \ln f_{\vec{x}}(\vec{x}; \theta) = \ln L(\theta)$$

i.e.
$$\hat{\theta}^{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \{ L(\theta; \vec{x}) \}$$

ex: $\{X_i\}_{i=1}^n$: iid ; $f_x(x; a) = 2ax \cdot \exp(-ax^2)$; $x \in \mathbb{R}^+$

$$\Rightarrow L(a; \vec{x}) = \log \left(\prod_{i=1}^n f_x(x_i; a) \right) = \sum_{i=1}^n \ln(2ax_i \cdot \exp(-ax_i^2))$$

$$L(a; \vec{x}) = \sum_{i=1}^n [\ln 2 + \ln a + \ln x_i - ax_i^2]$$

$$\frac{\partial L(a; \vec{x})}{\partial a} = \sum_{i=1}^n \frac{1}{a} - x_i^2 = \frac{n}{a} - \sum_{i=1}^n x_i^2$$

$$L'(a; \vec{x}) = 0$$

$$\Rightarrow n / \hat{a}^{\text{MLE}} = \sum_{i=1}^n x_i^2$$

$$\therefore \hat{a}^{\text{MLE}} = \frac{n}{\sum_{i=1}^n x_i^2}$$

The general MLE procedure for data $\{X_i, Y_i\}_{i=1}^n$:

i/ Formulate $L(\theta; \vec{X})$:

$$L(\theta; \vec{X}) = \ln f_{\vec{X}}(\vec{X}; \theta)$$

ii/ Solve for critical point of $L(\theta; \vec{X})$:

find $\hat{\theta}^{MLE}$ such that $\left. \frac{\partial L}{\partial \theta} \right|_{\theta = \hat{\theta}^{MLE}} = 0$

iii/ Make sure $\hat{\theta}^{MLE}$ is a maximum crit. pt. for $L(\theta; \vec{X})$

i.e. $\left. \frac{\partial^2 L}{\partial \theta^2} \right|_{\theta = \hat{\theta}^{MLE}} < 0$

ML Estimates have 3 properties:

i/ $\hat{\theta}^{MLE}$ is consistent:

$$\hat{\theta}^{MLE} \xrightarrow{P} \theta$$

ii/ $\hat{\theta}^{MLE}$ obey the Invariance Principle (IP):

$$g(\theta)^{MLE} = g(\hat{\theta}^{MLE})$$

iii/ $\hat{\theta}^{MLE}$ is asymptotically normal:

$$\hat{\theta}^{MLE} \xrightarrow{d} N\left(\theta, \frac{1}{J_n(\theta)}\right)$$

$J_n(\theta)$: Fisher Information for θ .

11-11

Ex: $f(x; \theta) = \theta x^{\theta-1} \quad x \in (0, 1) \quad \theta \in \mathbb{R}^+$
 $\{X_i\}_{i=1}^n : \text{iid} \quad \text{(a) } \hat{\theta}^{\text{MLE}} \quad \text{(b) } E[\hat{\theta}^{\text{MLE}}]$

$$\Rightarrow L(\theta; \vec{x}) = \ln \left(\prod_{i=1}^n \theta x_i^{\theta-1} \right)$$

$$= \sum_{i=1}^n [\ln \theta + (\theta-1) \ln x_i]$$

$$= n \ln \theta + (\theta-1) \sum_{i=1}^n \ln x_i$$

$$\partial L / \partial \theta = n / \theta + \sum_{i=1}^n \ln x_i$$

$$\partial L / \partial \theta \Big|_{\theta=\hat{\theta}} = 0 \Rightarrow n / \hat{\theta} = - \sum_{i=1}^n \ln x_i$$

$$\Rightarrow \boxed{\hat{\theta} = \frac{-n}{\sum_{i=1}^n \ln x_i}}$$

In more general case:

$$\vec{\theta} = \{\theta_k\}_{k=1}^p$$

ie $\nabla_{\vec{\theta}} L(\vec{\theta}) \Big|_{\vec{\theta}=\hat{\theta}} = \vec{0}$

$$\hat{\theta}^{\text{MLE}} = \underset{\vec{\theta}}{\text{arg max}} \{L(\vec{\theta}; \vec{x})\}$$

e.g. $X_i \sim N(\mu, \sigma^2)$

$$\vec{\theta} = (\mu, \sigma^2)^T$$

$$\Rightarrow \hat{\theta}_n^{\text{MLE}} = \begin{pmatrix} \bar{X}_n \\ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \end{pmatrix}$$

We can view the log-likelihood as a statistic:

$$L(\theta; \vec{X}) = \sum_{i=1}^n \ln f_X(X_i; \theta)$$

How much does each extra sample of X improve the ^{average} precision of estimators for θ ?

[Assume $\hat{\theta}_n$ is unbiased]

The "Fisher Information" quantifies this idea:

$$J_n(\theta) = E_{X|\theta} \left[\left(\frac{\partial \ln f(\vec{X}|\theta)}{\partial \theta} \right)^2 \right] \stackrel{\text{under "regularity conditions"}}{=} E_{X|\theta} \left[\frac{-\partial^2 \ln f(\vec{X}|\theta)}{\partial \theta^2} \right]$$

[Precision of estimator $\hat{\theta}_n \propto 1/\text{var}(\hat{\theta}_n)$]

The Cramer-Rao Lower Bound:

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

establishes a lower bound for the variance of any unbiased estimator $\hat{\theta}_n$.

$J_n(\theta)$ quantifies how much "information" X has about θ on average. Any unbiased $\hat{\theta}_n$ cannot "squeeze out" $\Rightarrow J_n(\theta)$ worth of info out of X on average.

Maximum ^(MAP) A Posteriori Estimation:

Make the "Bayesian Assumption" that the unknown θ is also a rv with pdf $h(\theta)$ (prior).

Then the best estimate for θ is one that maximizes the posterior pdf $f(\theta|x)$.

This is the MAP Estimate $\hat{\theta}^{\text{MAP}}$.

$$\hat{\theta}^{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} (f(\theta|x))$$

$$\Leftrightarrow \hat{\theta}^{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} (g(x|\theta) \cdot h(\theta))$$

[By Bayes theorem]

i/ Recall that $\hat{\theta}^{\text{MAP}}$ is a Bayes-optimal estimate for the loss function $l(d; \theta) = 1 - \delta(\theta - d)$.

ii/ $\hat{\theta}^{\text{MAP}} = \hat{\theta}^{\text{MLE}}$ when $h(\theta) = \text{constant} \forall \theta$.

iii/ The prior pdf $h(\theta)$ encodes prior/subjective info about θ .

MLE/MAP rely heavily on the parametric form of $g(x|\theta)$. In practice, $g(x|\theta)$ might be difficult or too complicated to analyze. This often happens when X is a complicated/corrupted version of a simple rv Z .
 i.e. $X = T(Z)$

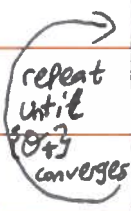
where $T(\cdot)$ is a corruption process.

The Expectation-Maximization ^(EM) algorithm is a method for handling such complicated MLE problems. The approach is to replace the complicated $\ln g(x|\theta)$ with a simpler surrogate:

$$Q(\theta|\theta_t) = E [\ln f(z|\theta) | X, \theta_t]$$

This is intended as a simpler $L(\theta)$ for X .

EM is traditionally done in 2 steps:

	E-step: $Q(\theta \theta_t) = E_{z X, \theta_t} [\ln f(z \theta) X, \theta_t]$
	M-step: $\theta_{t+1} = \operatorname{argmax} [Q(\theta \theta_t)]$

The procedure is guaranteed to converge to locally optimal estimates for θ .

All estimation methods so far have assumed a pdf family $f(x|\theta)$ for the observed data.

Suppose we want to match observed data $\{x_i\}_{i=1}^n$ to a pdf but we have no reason to pick any one pdf family over another. Then choose a pdf family with the maximum Entropy or highest disorder.

Entropy:

The entropy of X , $H(X)$, is the average information associated with X , $I(x) = \log\left(\frac{1}{P(x=x)}\right)$

$$\begin{aligned} H(X) &= E_x \left[\log \frac{1}{P(x=x)} \right] \quad [X \text{ discrete}] \\ &= \sum_k -\log(P(x=k)) \cdot P(x=k) \quad [\text{discrete}] \\ &= - \int \log f_x(x) \cdot f_x(x) dx \quad [\text{continuous}] \end{aligned}$$

e.g 1 $X \sim \text{Bernoulli}(p)$

$$\Rightarrow H(X) = -p \log p - (1-p) \log(1-p)$$

e.g 2 $X \sim \text{Uniform}(1, \dots, n)$

$$\Rightarrow H(X) = \sum_k -\frac{1}{n} \log\left(\frac{1}{n}\right) = -\frac{1}{n} \cdot \log\left(\frac{1}{n}\right) = -\log\left(\frac{1}{n}\right) = \log(n)$$

$$\log(y) \text{ concave} \Rightarrow E_x \left[\log\left(\frac{1}{P(x)}\right) \right] \leq \log \left[E_x \left[\frac{1}{P(x)} \right] \right] \quad (\text{Jensen's})$$

$$\therefore -H(X) \geq -\log(n) \quad (\Leftrightarrow) \quad \boxed{H(X) \leq \log(n)}$$

$$\text{i.e.} \quad \boxed{H(X) \leq H(U)} \quad \text{if } X, U \in \{1, \dots, n\} \\ U \sim \text{Uniform}$$

Also:

i/ Conditional Entropy:

$$H(Y|X) = E_{Y|X} [-\ln P(Y|X)]$$

$$[\text{Thm} : H(Y) \geq H(Y|X)]$$

ii/ Mutual Information:

$$I(X, Y) = \sum_x \sum_y P(x, y) \log \frac{P(x, y)}{P(x) \cdot P(y)}$$

$$I(X, Y) = E_{x, y} \left[\log \left(\frac{P(x, y)}{P(x) \cdot P(y)} \right) \right]$$

$$[I(X, Y) = 0] \Leftrightarrow [X, Y \text{ independent}]$$

iii/ Joint Entropy:

$$H(X, Y) = E_{x, y} [-\log P(x, y)]$$

$$I(X, Y) = H(X) - H(X|Y)$$

$$H(X, Y) = H(X) + H(Y|X)$$

The ^(Max Ent) maximum entropy principle says match the data to the pdf with the maximum entropy, subject to constraints on moments. i.e. find pdf of $\{x_i\}_{i=1}^n$ such that:

$$i/ \left\{ E[X^k] = \frac{1}{n} \sum_{i=1}^n x_i^k \right\}_{k=0}^m \leftarrow \text{match } m \text{ moments.}$$

ii/ $f_X(x)$ has the max entropy for X .

We already saw that $U \sim U(1, \dots, n)$ has the Max Ent on $\{1, \dots, n\}$ without any moment constraints ($m=0$)
 $m=1 \Rightarrow \left\{ \begin{array}{l} X \sim \text{exp} \\ \text{or } X \sim \text{geometric} \end{array} \right\}$; $m=2 \Rightarrow X \sim N(\cdot, \cdot)$

Dropping the Max Ent constraint (ii/) reduces this to the historical method of moments:

e.g. $\{X_i\}$: iid $\sim \gamma(\alpha, \theta)$

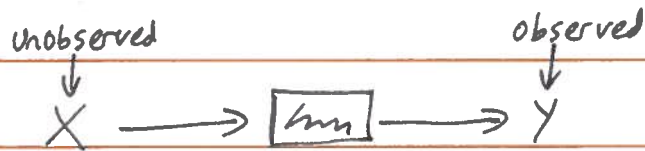
find (α, θ) such that:

$$E[X] = \alpha\theta = \frac{1}{n} \sum_i X_i$$

$$E[X^2] = \alpha\theta^2 + \alpha^2\theta^2 = \frac{1}{n} \sum_i X_i^2$$

Minimum Mean-Squared Error Estimation:

A different kind of Estimation problem:



Find an estimate for X using Y : $g(Y)$

such that $g(Y)$ minimizes the average mean-square estimation error $E[(X - g(Y))^2]$

$$\Rightarrow \boxed{\hat{X} = g(Y) = E[X|Y]} \quad \text{[MMSE]}$$

This means that the estimation error is orthogonal to Y (no further "information" of value in Y about X)

$$E[(g(Y) - X) \cdot Y] = 0$$

If we constrain $g(\cdot)$ to ^{also} be linear we get the Linear MMSE (LMMSE)

(MMSE = LMMSE) \Leftrightarrow (X, Y are jointly Gaussian)

\vec{X}, \vec{Y} JG

$$\Rightarrow \boxed{E[\vec{X}|\vec{Y}] = \vec{\mu}_X + \underline{K}_{xy} \underline{K}_{yy}^{-1} (\vec{Y} - \vec{\mu}_Y)} \quad \text{[multi-dim]}$$

$$\mu_x + \rho_{xy} \sigma_x / \sigma_y (y - \mu_y) \quad \text{[1-dim]}$$

$$\underline{K}_{x|y} = \underline{K}_{xx} - \underline{K}_{xy} \underline{K}_{yy}^{-1} \underline{K}_{yx} ; \quad \sigma_{x|y}^2 = \sigma_x^2 (1 - \rho_{xy}^2)$$