

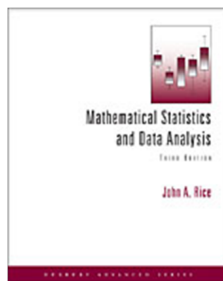
ST122: Probability and Statistics II

Waleed A. Yousef, Ph.D.,

Human Computer Interaction Lab.,
Computer Science Department,
Faculty of Computers and Information,
Helwan University,
Egypt.

March 24, 2019

Lectures follow Rice, “*Mathematical Statistics and Data Analysis*”, 3rd edition, Duxbury:



ISBN 0-534-39942-8

Course Objectives

- Developing rigorous treatment.
- Building intuition and insight.
- Linking to real life problems.
- Coding and scientific computing.

Contents

Contents	iii
Introduction: Statistical Inference in a Nutshell	iv
6 Distributions Derived from the Normal Distribution	1
6.1 Introduction	2
6.2 χ^2 , t , and F Distributions	3
6.3 Sample Mean, Sample Variance, and Sampling from Normal Distribution	9
6.3.1 Basic Concepts of Random Samples	9
6.3.2 Sampling from the Normal Distribution	15
8 Estimation of Parameters and Fitting of Probability Distributions	22
8.1 Introduction:	
Estimation in a Nutshell	23
8.2 The Method of Moments	26
8.3 The Method of Maximum Likelihood	37
8.3.1 Large Sample Theory for MLE	49
8.4 The Bayesian Approach to Parameter Estimation	60
8.4.1 Large Sample Theory of Bayesian Inference	68
8.5 Assessing Estimators, Efficiency, and the Cramér-Rao Lower Bound	69
8.5.1 Mean Squared Error (MSE) Criterion	69
8.5.2 Best Unbiased Estimator	74
8.5.3 Asymptotic Relative Efficiency (ARE)	94

Introduction: Statistical Inference in a Nutshell

Point estimate - different estimators - assessing estimators - large sample theory

Hypothesis testing.

Interval estimation.

Bayesian approach vs. Frequentist approach

Chapter 6

Distributions Derived from the Normal Distribution

6.1 Introduction

This Chapter discusses 3 probability distributions that frequently occur in Statistics: χ^2 , t , and F Distributions.

Remember that if $V \sim \text{Gamma}(\alpha, \lambda)$, then

$$f(v) = \frac{\lambda^\alpha}{\Gamma(\alpha)} v^{\alpha-1} e^{-\lambda v}, \quad v \geq 0,$$

$$M(t) = (1 - t/\lambda)^{-\alpha},$$

$$E[V] = \alpha/\lambda,$$

$$\text{Var}[V] = \alpha/\lambda^2.$$

And if V_1, \dots, V_n are i.i.d $\text{Gamma}(\alpha, \lambda)$, then

$$M_{\sum_i V_i}(t) = (1 - t/\lambda)^{-n\alpha},$$

$$\sum_i V_i \sim \text{Gamma}(n\alpha, \lambda).$$

6.2 χ^2 , t , and F Distributions

Definition 1 If $Z \sim N(0, 1)$, then $U = Z^2$ is called *chi-square distribution with 1 degree of freedom*; i.e., $U \sim \chi_1^2$. It is easy to show that (see Lec. notes Ch. 2):

$$f_U(u) = \frac{1}{\sqrt{2\pi}} u^{-1/2} e^{-u/2}.$$

Notice that:

$$\chi_1^2 \equiv \text{Gamma}\left(\frac{1}{2}, \frac{1}{2}\right),$$

Also:

$$\begin{aligned} X &\sim N(\mu, \sigma^2), \\ \frac{X - \mu}{\sigma} &\sim N(0, 1), \\ \left(\frac{X - \mu}{\sigma}\right)^2 &\sim \chi_1^2. \end{aligned}$$

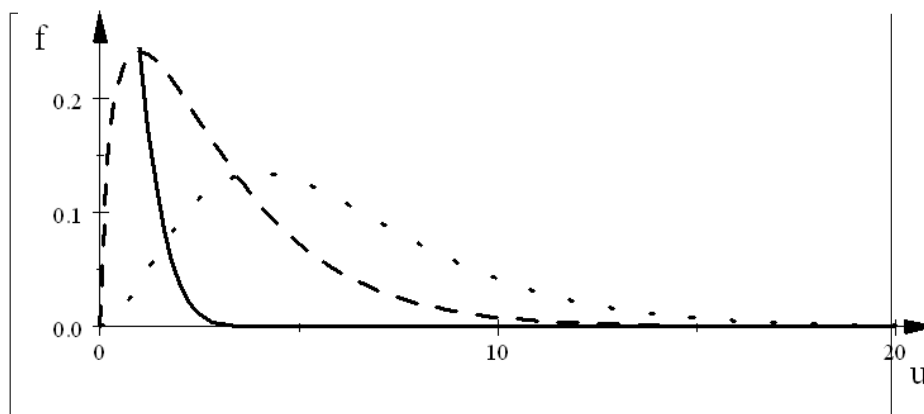
Definition 2 If U_1, \dots, U_n are i.i.d χ_1^2 r.v. then $V = \sum_i U_i$ is called chi-square distribution with n degrees of freedom; i.e., $V \sim \chi_n^2$.

Notice that $U_i \sim \text{Gamma}(\frac{1}{2}, \frac{1}{2})$, then

$$V \sim \text{Gamma}(n/2, 1/2),$$

$$f_V(v) = \frac{1}{2^{n/2} \Gamma(n/2)} v^{n/2-1} e^{-v/2},$$

$$E[V] = n, \text{ Var}[V] = 2n.$$



solid: $n = 1$, dashed: $n = 3$, dotted: $n = 6$

Suppose that U and V are indep, and

$$W = U + V.$$

If $U \sim \chi_m^2$, $V \sim \chi_n^2$ then (obviously)

$$W = \chi_m^2 + \chi_n^2 = \chi_{m+n}^2,$$

Also, if $W \sim \chi_k^2$ and $V \sim \chi_n^2$ then

$$\chi_k^2 = U + \chi_n^2$$

$$M_{\chi_k^2} = M_U M_{\chi_n^2},$$

$$\begin{aligned} M_U &= \frac{M_{\chi_k^2}}{M_{\chi_n^2}} \\ &= \frac{(1-2t)^{-k/2}}{(1-2t)^{-n/2}} = (1-2t)^{-(k-n)/2} \end{aligned}$$

$$U \sim \chi_{(k-n)}^2.$$

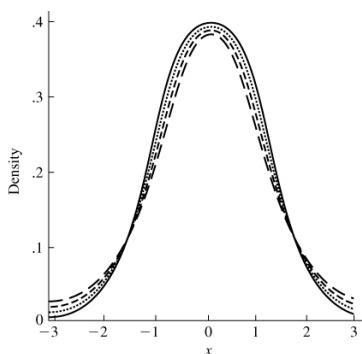
Definition 3 (Student's t Distribution) :

If $Z \sim N(0, 1)$, $U \sim \chi_n^2$, and Z, U are indep. then $T = Z/\sqrt{U/n}$ is called t distribution with n degrees of freedom; i.e., $T \sim t_n$. (prove that:)

$$f_T(t) = \frac{\Gamma((n+1)/2)}{\sqrt{n\pi}\Gamma(n/2)} \left(1 + \frac{t^2}{n}\right)^{-(n+1)/2},$$

$$E[T] = 0, \quad n \geq 2,$$

$$\text{Var}[T] = \frac{n}{n-2}, \quad n \geq 3.$$



- The smaller n the thicker tail.
- The figure shows $t_5, t_{10}, t_{30} (\approx N(0, 1))$
- $t_1 \equiv \text{Cauchy}(0, 1)$.

Definition 4 (Snedecor's F Distribution) :

Let $U \sim \chi_m^2$ and $V \sim \chi_n^2$, and U, V are indep. Then, $W = (U/m) / (V/n)$ is called F distribution with m, n degrees of freedom; i.e., $W \sim F_{m,n}$. (prove that:)

$$f_W(w) = \frac{\Gamma\left(\frac{m+n}{2}\right)}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} \left(\frac{m}{n}\right)^{\frac{m}{2}} w^{\frac{m}{2}-1} \left(1 + \frac{m}{n}w\right)^{-\frac{(m+n)}{2}},$$

$$E[W] = n / (n - 2), \quad n \geq 3.$$

$$\text{Var}[W] = 2 \left(\frac{n}{n-2}\right)^2 \frac{(m+n-2)}{m(n-2)}, \quad n \geq 5.$$

It is obvious that if $U \sim t_n$, then $U^2 \sim F_{1,n}$.

Also, if $U \sim F_{n,m}$ then $U^{-1} \sim F_{m,n}$.

Summary (with terse notation):

$$N(0, 1)^2 \sim \chi_1^2,$$

$$\sum_{i=1}^n N(0, 1)^2 \sim \chi_n^2,$$

$$\chi_m^2 + \chi_n^2 \sim \chi_{m+n}^2,$$

$$N(0, 1) / \sqrt{\chi_n^2 / n} \sim t_n,$$

$$(\chi_m^2 / m) / (\chi_n^2 / n) \sim F_{m,n},$$

$$t_n^2 \sim F_{1,n}.$$

Example 5 *If X_1, X_2, X_3 are iid $N(0, 1)$, what is the dist. of*

$$\frac{X_1}{\sqrt{(X_1^2 + X_2^2 + X_3^2) / 3}}$$

6.3 Sample Mean, Sample Variance, and Sampling from Normal Distribution

6.3.1 Basic Concepts of Random Samples

Definition 6 *The r.v. X_1, \dots, X_n are called a **random sample of size n from the population F** if X_1, \dots, X_n are i.i.d from F ; and hence:*

$$f_{X_1 \dots X_n}(x_1, \dots, x_n) = \prod_i f(x_i).$$

$$\begin{array}{ccccccc} & & X_1 & X_2 & \dots & X_n & \\ F & \xrightarrow{\text{Sample}_1} & x_1, & x_2, & \dots & x_n & \\ F & \xrightarrow{\text{Sample}_2} & x_1, & x_2, & \dots & x_n & \\ & \vdots & & & & & \end{array}$$

We focus in our study on infinite populations; Ch. 7 is about finite populations.

Definition 7 Let X_1, \dots, X_n be a random sample of size n , and $T(x_1, \dots, x_n)$ be a real- (or vector-) valued function whose domain includes the sample space of (X_1, \dots, X_n) . Then the r.v. $Y = T(X_1, \dots, X_n)$ is called a statistic.

Definition 8 The sample mean, sample variance, and sample standard deviations are statistics defined as:

$$\begin{aligned}\bar{X} &= \frac{1}{n} \sum_i X_i \\ S^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2, \\ S &= \sqrt{S^2},\end{aligned}$$

Observed values will be denoted by \bar{x} , s^2 , and s .

		X_1	X_2	\dots	X_n	$\bar{X} = \frac{1}{n} \sum_i X_i$
F	<u>Sample₁</u>	$x_1,$	$x_2,$	\dots	x_n	$\bar{x} = \frac{1}{n} \sum_i x_i$
F	<u>Sample₂</u>	$x_1,$	$x_2,$	\dots	x_n	$\bar{x} = \frac{1}{n} \sum_i x_i$
	\vdots					

Lemma 9 For any numbers x_1, \dots, x_n :

$$\begin{aligned}\min_a \sum_i (x_i - a)^2 &= \sum_i (x_i - \bar{x})^2, \\ \sum_i (x_i - \bar{x})^2 &= \sum_i x_i^2 - n\bar{x}^2.\end{aligned}$$

Proof. : is identical to $\arg\min_c E(Y - c)^2 = E[Y]$.

$$\begin{aligned}\sum_i (x_i - a)^2 &= \sum_i ((x_i - \bar{x}) + (\bar{x} - a))^2 \\ &= \sum_i (x_i - \bar{x})^2 + \sum_i (\bar{x} - a)^2 \\ &\quad + 2 \sum_i (x_i - \bar{x})(\bar{x} - a) \quad (\sum_i x_i = n\bar{x}) \\ &= \sum_i (x_i - \bar{x})^2 + \sum_i (\bar{x} - a)^2,\end{aligned}$$

which is minimized by choosing $a = \bar{x}$.

$$\begin{aligned}\sum_i (x_i - a)^2 &= \sum_i (x_i - \bar{x})^2 + \sum_i (\bar{x} - a)^2 \\ \sum_i (x_i - \bar{x})^2 &= \sum_i x_i^2 - n\bar{x}^2. \quad (a \stackrel{set}{=} 0)\end{aligned}$$

Notice that: both forms are $O(n)$; however this form requires only one for loop for execution! ■

HW: Write a computer program, and find its complexity (where a step is a multiplication), for calculating

$$S_1 = \sum_{i=1}^n \sum_{j=1}^n x_i x_j,$$

$$S_2 = \sum_{i=1}^n \sum_{j \neq i}^n x_i x_j.$$

Can you do a mathematical trick to reduce their complexities to $O(n)$. !!!

Theorem 10 (Distribution-Free Properties) :

1. $E[\bar{X}] = \mu,$
2. $\text{Var}[\bar{X}] = \sigma^2/n,$
3. $E[S^2] = \sigma^2.$

Proof. 1 and 2 are proven before. For 3,

$$\begin{aligned} E[S^2] &= E\left[\frac{1}{n-1} \sum_i (X_i - \bar{X})^2\right] \\ &= \frac{1}{n-1} E\left[\sum_i X_i^2 - n\bar{X}^2\right] \\ &= \frac{1}{n-1} \left(\sum_i E[X_i^2] - nE[\bar{X}^2]\right) \\ &= \frac{1}{n-1} \left(n(\sigma^2 + \mu^2) - n\left(\frac{\sigma^2}{n} + \mu^2\right)\right) = \sigma^2, \end{aligned}$$

which is **unbiased estimator** for σ^2 . ■

Lemma 11 Let X_1, \dots, X_n be a r.s. from a population with mgf $M(t)$, then

$$M_{\bar{X}}(t) = [M(t/n)]^n.$$

Proof. done before in CLT (just 2 lines). ■

Example 12 Let X_1, \dots, X_n be a r.s. from $N(\mu, \sigma^2)$, then

$$\begin{aligned} M(t) &= \exp(\mu t + \sigma^2 t^2 / 2), \\ M_{\bar{X}}(t) &= \left[\exp\left(\mu \frac{t}{n} + \sigma^2 \left(\frac{t}{n}\right)^2 / 2\right) \right]^n, \\ &= \exp\left(\mu t + \frac{\sigma^2}{n} t^2 / 2\right), \\ \bar{X} &\sim N\left(\mu, \frac{\sigma^2}{n}\right). \end{aligned}$$

We know that $E[\bar{X}] = \mu$ and $\text{Var}[\bar{X}] = \sigma^2/n$. But what is new is that \bar{X} is itself Normal. **We could have found it by transformation:** $Z = X_1 + X_2$. If $X_i \sim \text{Cauchy}(0, 1)$, prove that $\bar{X} \sim \text{Cauchy}(0, 1)$ as well!!

6.3.2 Sampling from the Normal Distribution

Theorem 13 Let X_1, \dots, X_n be r.s. form $N(\mu, \sigma^2)$

1. $\bar{X} \sim N(\mu, \sigma^2/n)$,
2. \bar{X} and $(X_2 - \bar{X}, \dots, X_n - \bar{X})$ are indep,
3. \bar{X} and S^2 are indep,
4. $(n-1) S^2 / \sigma^2 \sim \chi_{n-1}^2$.

Intuition before proof:

Meaning of \bar{X} and $(X_2 - \bar{X}, \dots, X_n - \bar{X})$ are indep?

Suppose $X_i \sim \text{Bernouli}(1/2)$, and we get a sample where $\bar{X}_{10} = 1$. Obviously, $X_i = 1$.

Aside from normality, observe that

$$\sum_i \left(X_i - \bar{X} \right) = 0,$$

which means we have only $(n - 1)$ differences:

$$\begin{aligned} \left(X_1 - \bar{X} \right) &= - \sum_{i=2}^n \left(X_i - \bar{X} \right), \\ S^2 &= \frac{1}{(n-1)} \sum_i \left(X_i - \bar{X} \right)^2 \\ &= \frac{1}{(n-1)} \left[\left(X_1 - \bar{X} \right)^2 + \sum_{i=2}^n \left(X_i - \bar{X} \right)^2 \right] \\ &= \frac{1}{(n-1)} \left[\left(\sum_{i=2}^n \left(X_i - \bar{X} \right) \right)^2 + \sum_{i=2}^n \left(X_i - \bar{X} \right)^2 \right] \end{aligned}$$

Matlab Code 6.1:

figure; hold on;

% Change 'Normal' to 'Exp'

x=random('Normal', 0, 1, 1000, 10);

xbar=mean(x, 2);

s=std(x, 0, 2);

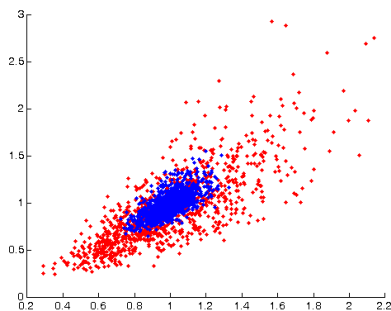
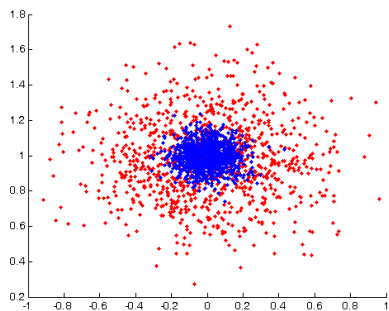
plot(xbar, s, '.r')

x=random('Normal', 0, 1, 1000, 100);

xbar=mean(x, 2);

s=std(x, 0, 2);

plot(xbar, s, '.b')



Proof. the mgf is given by

$$\begin{aligned}
&= M(s, t_2, \dots, t_n) \\
&= E \left[\exp \left(s\bar{X} + t_2 (X_2 - \bar{X}) + \dots + t_n (X_n - \bar{x}) \right) \right] \\
&= E \left[\exp \left(\sum_{i=1}^n \frac{s}{n} X_i + \sum_{i=2}^n t_i (X_i - \bar{X}) \right) \right] \\
&= E \left[\exp \left(\sum_{i=1}^n \left(\frac{s}{n} + (t_i - \bar{t}) \right) X_i \right) \right] \quad (t_1 = 0) \\
&= E \left[\exp \left(\sum_{i=1}^n a_i X_i \right) \right] \quad (a_i = \frac{s}{n} + (t_i - \bar{t})) \\
&= \prod_i M_{X_i}(a_i) \\
&= \prod_i \exp \left(\mu a_i + \frac{\sigma^2}{2} a_i^2 \right) \\
&= \exp \left[\mu \sum_i a_i + \frac{\sigma^2}{2} \sum_i a_i^2 \right] \\
&= \exp \left[\mu s + \frac{\sigma^2}{2} \left(\frac{s^2}{n} + \sum_i (t_i - \bar{t})^2 \right) \right] \\
&= \exp \left(\mu s + \frac{\sigma^2}{2n} s^2 \right) \exp \left(\frac{\sigma^2}{2} \sum_i (t_i - \bar{t})^2 \right),
\end{aligned}$$

the two factors are the mgf of \bar{X} and $(X_2 - \bar{X}, \dots, X_n - \bar{X})$. Hence they are independent and since $S = S(X_2 - \bar{X}, \dots, X_n - \bar{X})$: \bar{X} and S are independent.

Now

$$\begin{aligned}\sum_i \left(\frac{X_i - \mu}{\sigma} \right)^2 &= \frac{1}{\sigma^2} \sum_i \left[(X_i - \bar{X}) + (\bar{X} - \mu) \right]^2 \\ &= \frac{1}{\sigma^2} \sum_i (X_i - \bar{X})^2 + \frac{1}{\sigma^2} \sum_i (\bar{X} - \mu)^2 \\ &= \frac{1}{\sigma^2} \sum_i (X_i - \bar{X})^2 + \left(\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \right)^2\end{aligned}$$

$$W = U + V \quad (U, V \text{ indep.})$$

$$\chi_n^2 = U + \chi_1^2$$

$$U \sim \chi_{n-1}^2. \quad (n-1 \text{ df})$$



Lemma 14

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}.$$

Proof.

$$\begin{aligned}\frac{\bar{X} - \mu}{S/\sqrt{n}} &= \frac{(\bar{X} - \mu) / (\sigma/\sqrt{n})}{(S/\sqrt{n}) / (\sigma/\sqrt{n})} \\ &= \frac{(\bar{X} - \mu) / (\sigma/\sqrt{n})}{S/\sigma} \\ &= \frac{(\bar{X} - \mu) / (\sigma/\sqrt{n})}{\sqrt{((n-1) S^2 / \sigma^2) / (n-1)}} \\ &= \frac{N(0, 1)}{\sqrt{\chi_{n-1}^2 / (n-1)}} = t_{n-1},\end{aligned}$$

used for inference about μ when σ is unknown.

$$\frac{\bar{X} - \mu}{\sigma} \sim N(0, 1)$$

used for inference about μ when σ is known. ■

Lemma 15 *If $X \sim N(\mu_X, \sigma_X)$, $Y \sim N(\mu_Y, \sigma_Y)$, and we have two samples X_1, \dots, X_m and Y_1, \dots, Y_n*

$$\frac{S_X^2 / \sigma_X^2}{S_Y^2 / \sigma_Y^2} \sim F_{m-1, n-1}.$$

Proof.

$$\begin{aligned} \frac{S_X^2 / \sigma_X^2}{S_Y^2 / \sigma_Y^2} &= \frac{((m-1) S_X^2 / \sigma_X^2) / (m-1)}{((n-1) S_Y^2 / \sigma_Y^2) / (n-1)} \\ &= \frac{\chi_{m-1}^2 / (m-1)}{\chi_{n-1}^2 / (n-1)} && \text{(Indep.)} \\ &= F_{m-1, n-1}, \end{aligned}$$

used for inference about σ_X^2 / σ_Y^2 . ■

Chapter 8

Estimation of Parameters and Fitting of Probability Distributions

8.1 Introduction:

Estimation in a Nutshell

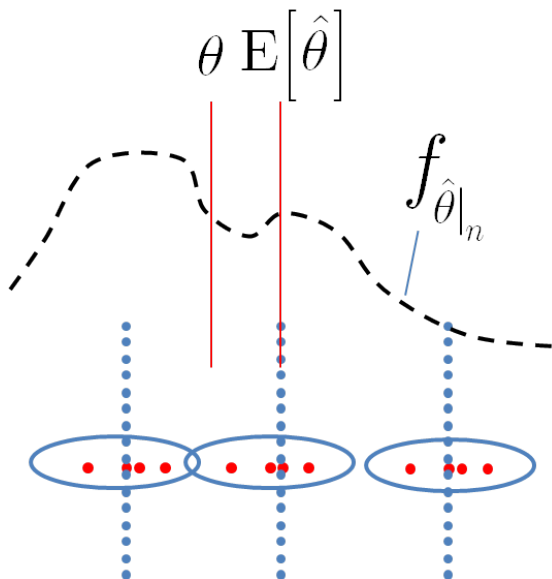
- Distributions depend on some population parameters; e.g., $N(\mu, \sigma^2)$, $Exp(\lambda)$, etc. Generally, we should write (e.g.):

$$f_X(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[\frac{-1}{2} (x - \mu)^2 / \sigma^2 \right]$$

- Obtaining data (values of a random sample) allows “estimating” these parameters.

Definition 16 *A point estimator is any function $W(X_1, \dots, X_n)$ of a sample; i.e., any statistic is a point estimator.*

- We can choose, e.g., $\hat{\sigma}^2 = \frac{1}{n} \sum_i (X_i - \bar{X})^2$ to be an estimator for σ^2 .
- $\frac{1}{n} \sum_i (x_i - \bar{x})^2$ is an estimate (realization).



- How to estimate θ “well” ($\hat{\theta}$)?
- What is $f_{\hat{\theta}}$ (**sampling distribution**)?
- What is $E[\hat{\theta}]$, $SD[\hat{\theta}]$ (**standard error**),...?
- How to estimate $\tau(\theta)$, e.g.:
 - σ^2 , the variance, for $N(\mu, \sigma^2)$.
 - $\alpha\lambda$, the mean, for $Gamma(\alpha, \lambda)$.

How to decide F_X before estimation?

- From the physics of the problem. E.g., given number of calls in time units, the distribution is known to be *Poisson* (λ).
- Assumption; you need to validate it latter.

Why do we estimate parameters?

- Understanding (interpretation).
- Prediction.
- Simulation and data generation.

How do we choose estimators?

8.2 The Method of Moments

We estimate k^{th} moment by **sample moment**

$$\mu_k = E[X^k]$$
$$\hat{\mu}_k = \frac{1}{n} \sum_i X_i^k.$$

Then for population parameters θ_i , we have

$$\mu_1 = \mu_1(\theta_1, \dots, \theta_r),$$
$$\vdots$$
$$\mu_r = \mu_r(\theta_1, \dots, \theta_r).$$

We solve

$$\theta_1 = \theta_1(\mu_1, \dots, \mu_r),$$
$$\vdots$$
$$\theta_r = \theta_r(\mu_1, \dots, \mu_r).$$

And

$$\hat{\theta}_1 = \hat{\theta}_1(\hat{\mu}_1, \dots, \hat{\mu}_r),$$
$$\vdots$$
$$\hat{\theta}_r = \hat{\theta}_r(\hat{\mu}_1, \dots, \hat{\mu}_r).$$

Motivation behind method of moments

$$\hat{\mu}_k \xrightarrow{p} \mu_k.$$

Definition 17 *An estimator $\hat{\theta} = \hat{\theta}(n)$, which estimates θ , from a sample of size n is said to be consistent in probability if*

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

Example 18 $N(\mu, \sigma^2)$, and the mean and variance of any other distribution:

$$\hat{\mu}_1 = \frac{1}{n} \sum_i X_i = \bar{X},$$

$$\hat{\mu}_2 = \frac{1}{n} \sum_i X_i^2,$$

$$\mu_1 = E[X] = \mu,$$

$$\mu_2 = E[X^2] = \mu^2 + \sigma^2,$$

$$\mu = \mu_1,$$

$$\sigma^2 = \mu_2 - \mu_1^2,$$

$$\hat{\mu} = \hat{\mu}_1 = \bar{X},$$

$$\hat{\sigma}^2 = \hat{\mu}_2 - \hat{\mu}_1^2 = \frac{1}{n} \sum_i X_i^2 - \bar{X}^2 \quad (\widehat{\sigma^2})$$

$$= \frac{1}{n} \left(\sum_i X_i^2 - n\bar{X}^2 \right) = \frac{1}{n} \sum_i (X_i - \bar{X})^2$$

$$= \frac{n-1}{n} S^2,$$

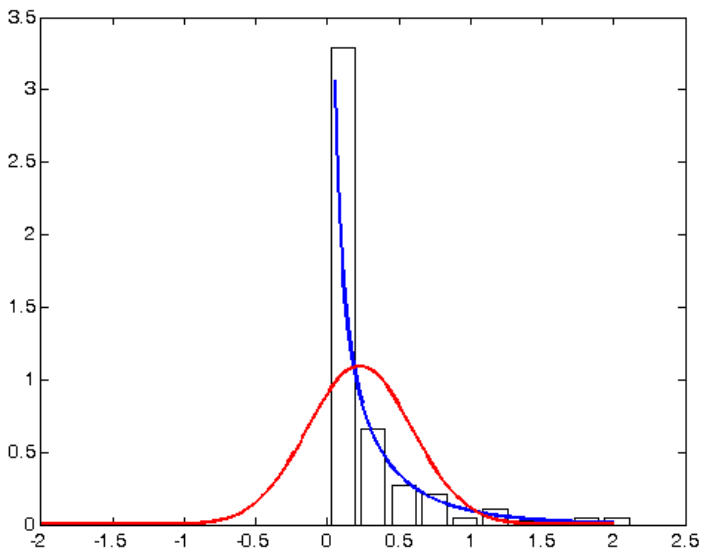
$$\hat{\mu} \sim N(\mu, \sigma^2/n),$$

$$\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-1}^2.$$

Example 19 : Analyzing real dataset for average amount of storms rainfall in Illinois.

Let's draw data points and normalized histogram (divide by its area):

$$\begin{aligned} Area &= \sum_i \Delta N_i \\ &= \Delta \sum_i N_i = \Delta n. \end{aligned}$$



From the mgf of Gamma we obtained

$$E[X] = \mu_1 = \frac{\alpha}{\lambda},$$

$$E[X^2] = \mu_2 = \frac{\alpha(\alpha + 1)}{\lambda^2},$$

Solve both equations for α and λ ,

$$\alpha = \lambda\mu_1$$

$$\mu_2 = \frac{\lambda^2\mu_1^2 + \lambda\mu_1}{\lambda^2},$$

$$= \mu_1^2 + \mu_1/\lambda,$$

$$\lambda = \frac{\mu_1}{\mu_2 - \mu_1^2},$$

$$\alpha = \frac{\mu_1^2}{\mu_2 - \mu_1^2},$$

$$\hat{\mu}_1 = \frac{1}{n} \sum x_i = 0.2244,$$

$$\hat{\mu}_2 = \frac{1}{n} \sum x_i^2 = 0.1836,$$

$$\hat{\lambda} = 1.6842,$$

$$\hat{\alpha} = 0.3779$$

$$f(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}$$

$$= 0.5178 x^{-0.6221} e^{-1.6842x}, x \geq 0$$

What would happen have if we fit $N(\mu, \sigma^2)$?

Matlab Code 8.1:

```
x=[];
x=[x; csvread('illinois60.txt')];
x=[x; csvread('illinois61.txt')];
x=[x; csvread('illinois62.txt')];
x=[x; csvread('illinois63.txt')];
x=[x; csvread('illinois64.txt')];

n=length(X) % will be 227
plot(x, zeros(length(x)), '.r')
[N, xout]=hist(x);
bar(xout, N/(n*(xout(2)-xout(1))), 'w'
    ); % normalize
hold on;
```

```

mul    = sum(x) / n                                % . 2 2 4 4
mu2    = sum(x.^2) / n                              % . 1 8 3 6
alpha=  mul^2 / (mu2-mul^2)                        % . 3 7 7 9
lmda   =  mul / (mu2-mul^2)                        % 1 . 6 8 4 2

```

```

z=0.05:.01:2;
y1=(lmda^alpha) / gamma(alpha) * z.^(
    alpha-1) .* exp(-lmda*z);
plot(z, y1, 'b', 'LineWidth', 2);

z=-2:.01:2;
y2=1 / (sqrt(2 * pi * (mu2-mul^2))) * exp(-(z
    -mul).^2 / (2 * (mu2-mul^2)));
plot(z, y2, 'r', 'LineWidth', 2);

```

Example 20 (*Binomial* (n, p))

$$\mu_1 = np,$$

$$\mu_2 = np(1-p) + (np)^2,$$

$$p = \frac{\mu_1}{n},$$

$$\mu_2 = \mu_1 \left(1 - \frac{\mu_1}{n}\right) + \mu_1^2$$

$$n = \frac{\mu_1^2}{\mu_1 - (\mu_2 - \mu_1^2)}$$

$$p = \frac{\mu_1 - (\mu_2 - \mu_1^2)}{\mu_1},$$

$$\hat{n} = \frac{\overline{X}^2}{\overline{X} - \frac{1}{n} \sum_i (X_i - \overline{X})^2},$$

$$\hat{p} = \frac{\overline{X} - \frac{1}{n} \sum_i (X_i - \overline{X})^2}{\overline{X}}.$$

- Sometimes the estimate will be negative!!
- In general, method of moments is a good start.

Example 21 (Cov(X, Y)) :

$$\begin{aligned}\sigma_X^2 &= E(X - \mu_X)^2 \\ &= E(X^2) - \mu_X^2 \\ &= \mu_{2X} - \mu_{1X}^2.\end{aligned}$$

$$\begin{aligned}\text{Cov}(X, Y) &= E(X - \mu_X)(Y - \mu_Y) \\ &= E[XY] - \mu_X\mu_Y \\ &= \mu_{11} - \mu_{1X}\mu_{1Y}\end{aligned}$$

$$\begin{aligned}\hat{\sigma}_X^2 &= \frac{1}{n} \sum_i X_i^2 - \bar{X}^2 \\ &= \frac{1}{n} \sum_i (X_i - \bar{X})^2.\end{aligned}$$

$$\begin{aligned}\hat{\sigma}_{XY} &= \frac{1}{n} \sum_i X_i Y_i - \bar{X}\bar{Y}. \\ &= \frac{1}{n} \sum_i (X_i - \bar{X})(Y_i - \bar{Y}).\end{aligned}$$

*Given x_1, \dots, x_n and y_1, \dots, y_m , what is $\hat{\sigma}_{XY}$?
What is right (x_i, y_i) .*

$$E[X_i Y_i] = \text{Cov}(X, Y) + \mu_X \mu_Y$$

$$\begin{aligned} E[\overline{XY}] &= \text{Cov}(\overline{X}, \overline{Y}) + E[\overline{X}] E[\overline{Y}] \\ &= \text{Cov}\left(\frac{1}{n} \sum_i X_i, \frac{1}{n} \sum_i Y_i\right) + \mu_X \mu_Y \\ &= \frac{1}{n^2} \sum_i \sum_j \text{Cov}(X_i, Y_j) + \mu_X \mu_Y \\ &= \frac{1}{n} \text{Cov}(X, Y) + \mu_X \mu_Y \end{aligned}$$

$$\begin{aligned} E \sum_i (X_i - \overline{X})(Y_i - \overline{Y}) &= \\ &= E \left[\sum_i X_i Y_i - n \overline{XY} \right] \\ &= n E[XY] - n E[\overline{XY}] \\ &= n \sigma_{XY} + n \mu_X \mu_Y - \sigma_{XY} - n \mu_X \mu_Y \\ &= (n - 1) \sigma_{XY}. \end{aligned}$$

Therefore, $\frac{1}{n} \sum_i (X_i - \overline{X})(Y_i - \overline{Y})$ is biased for σ_{XY} .

Another proof for $E[\overline{XY}]$:

$$\begin{aligned}
 E[\overline{XY}] &= E\left[\left(\frac{1}{n}\sum_i X_i\right)\left(\frac{1}{n}\sum_i Y_i\right)\right] \\
 &= E\left[\frac{1}{n^2}\sum_i\sum_j X_i Y_j\right] \\
 &= \frac{1}{n^2}E\left[\sum_i X_i Y_i + \sum_{i\neq j}\sum_j X_i Y_j\right] \\
 &= \frac{1}{n^2}(nE[XY] + n(n-1)E[X_i Y_j]) \\
 &= \frac{1}{n}(E[XY] + (n-1)E[X_i Y_j]) \\
 &= \frac{1}{n}(\text{Cov}(X, Y) + \mu_X\mu_Y + (n-1)\mu_X\mu_Y) \\
 &= \frac{1}{n}\text{Cov}(X, Y) + \mu_X\mu_Y.
 \end{aligned}$$

8.3 The Method of Maximum Likelihood

Likelihood is a function of parameters:

$$\begin{aligned}lik(\theta) &= f_{X_1 \dots X_n}(x_1, \dots, x_n | \theta) \\ &= \prod_{i=1}^n f(x_i | \theta). \quad (\text{i.i.d.})\end{aligned}$$

- For given data x_1, \dots, x_n , what is the value of θ that maximizes $lik(\theta)$.
- Remember Example 15, Page 19 in Lecture Notes.
- Much easier, in many cases, to deal with the **log likelihood** :

$$l(\theta) = \sum_{i=1}^n \log f(x_i | \theta).$$

Example 22 (*Poisson* (λ))

$$p(x) = \frac{\lambda^x e^{-\lambda}}{x!}, \quad 0 \leq x.$$

$$lik(\lambda) = p(x_1, \dots, x_x) = \prod_{i=1}^n \left(\frac{\lambda^{x_i} e^{-\lambda}}{x_i!} \right),$$

$$\begin{aligned} l(\lambda) &= \sum_{i=1}^n \log \left(\frac{\lambda^{x_i} e^{-\lambda}}{x_i!} \right) \\ &= \sum_i [x_i \log \lambda - \lambda - \log(x_i!)] \\ &= \log(\lambda) \sum_i x_i - n\lambda - \sum_i \log(x_i!) \quad (8.1) \end{aligned}$$

$$l'(\lambda) = \frac{\sum_i x_i}{\lambda} - n, \quad (l'(\lambda) \stackrel{\text{set}}{=} 0)$$

$$\hat{\lambda} = \frac{1}{n} \sum x_i = \bar{X}, \quad (\text{MoM})$$

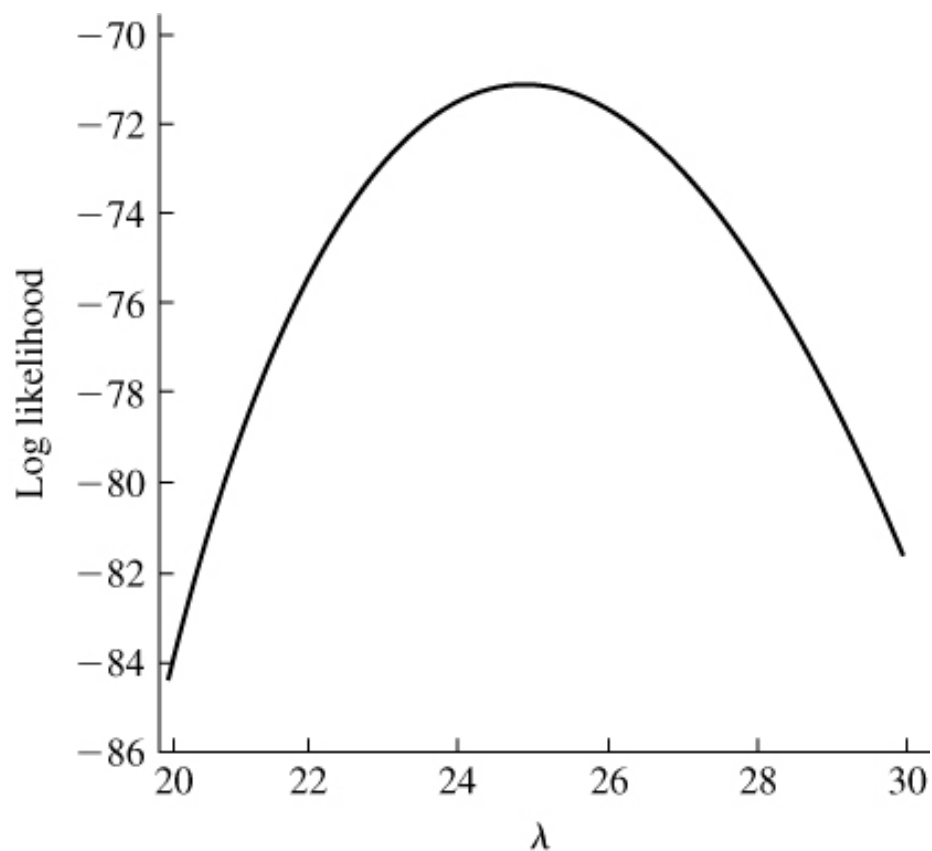
$$l''(\lambda) = \frac{-\sum_i x_i}{\lambda^2} \leq 0. \quad (x_i \geq 0)$$

Therefore, $\hat{\lambda} = \bar{X}$ is a point of local maxima; and

$$\lim_{\lambda \rightarrow \infty} l(\lambda) = -\infty,$$

then, $\hat{\lambda} = \bar{X}$ is a global maximum as well.

What does (8.1) mean for asbestos dataset?



Example 23 ($N(\mu, \sigma^2)$, both are unknown)

$$f(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[\frac{-1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2 \right]$$

$$\begin{aligned} l(\mu, \sigma) &= \sum_{i=1}^n \log f(x_i|\mu, \sigma) \\ &= \sum_i \left[-\log \sigma - \log \sqrt{2\pi} - \frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2 \right] \\ &= -n \log \sigma - n \log \sqrt{2\pi} - \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2 \end{aligned}$$

$$\frac{\partial l}{\partial \mu} = \frac{1}{\sigma^2} \sum_i (x_i - \mu) \quad \left(\frac{\partial l}{\partial \mu} \stackrel{\text{set}}{=} 0 \right)$$

$$0 = \sum_i x_i - n\hat{\mu},$$

$$\hat{\mu} = \frac{1}{n} \sum_i x_i = \bar{X}. \quad (\text{MoM})$$

$$\frac{\partial l}{\partial \sigma} = \frac{-n}{\sigma} + \frac{1}{\sigma^3} \sum_i (x_i - \mu)^2 \quad \left(\frac{\partial l}{\partial \sigma} \stackrel{\text{set}}{=} 0 \right)$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i \left(x_i - \bar{X} \right)^2. \quad (\text{MoM})$$

To verify that $(\hat{\mu}, \hat{\sigma})$ is a point of global maxima through calculus we have to satisfy:

First: it is a point of local maxima

- $\frac{\partial l}{\partial \mu}|_{\hat{\mu}} = \frac{\partial l}{\partial \sigma}|_{\hat{\sigma}} = 0$ (satisfied)
- $\frac{\partial^2 l}{\partial \mu^2}|_{\hat{\mu}} = 0$ or $\frac{\partial^2 l}{\partial \sigma^2}|_{\hat{\sigma}} = 0$ (satisfied)
- $\begin{vmatrix} \frac{\partial^2 l}{\partial \mu^2} & \frac{\partial^2 l}{\partial \mu \partial \sigma} \\ \frac{\partial^2 l}{\partial \mu \partial \sigma} & \frac{\partial^2 l}{\partial \sigma^2} \end{vmatrix}_{\hat{\mu}, \hat{\sigma}} > 0$ (needs work).

Second: there is no maximum at infinity (messy).

Instead, we can use a trick:

$$l(\mu, \sigma) = -n \log \sigma - n \log \sqrt{2\pi} - \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2$$

is maximized for

$$\sum_i (x_i - \mu)^2 = \sum_i (x_i - \bar{X})^2.$$

Then $l(\bar{X}, \sigma)$ is a function in single variable σ ,

$$\frac{\partial l}{\partial \sigma} = \frac{-n}{\sigma} + \frac{1}{\sigma^3} \sum_i (x_i - \bar{X})^2, \quad \left(\frac{\partial l}{\partial \sigma} \stackrel{set}{=} 0 \right)$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i (x_i - \bar{X})^2$$

$$\begin{aligned} \frac{\partial^2 l}{\partial \sigma^2} &= \frac{n}{\sigma^2} - \frac{3}{\sigma^4} \sum_i (x_i - \bar{X})^2 \\ &= \frac{n}{\sigma^2} \left(1 - \frac{3}{n\sigma^2} \sum_i (x_i - \bar{X})^2 \right), \end{aligned}$$

$$\left. \frac{\partial^2 l}{\partial \sigma^2} \right|_{\hat{\sigma}} = \frac{n}{\hat{\sigma}^2} (1 - 3) < 0,$$

which gives a local maximum for $l(\sigma)$. And

$$\lim_{\sigma \rightarrow \infty} l(\sigma) = -\infty.$$

Hence, $\hat{\sigma}$ attains a global maxima.

Example 24 ($\text{Gamma}(\alpha, \lambda)$) :

$$f(x) = \frac{1}{\Gamma(\alpha)} \lambda^\alpha x^{\alpha-1} e^{-\lambda x}, \quad 0 \leq x < \infty$$

$$l(\alpha, \lambda) = \sum_{i=1}^n (\alpha \log \lambda + (\alpha - 1) \log x_i - \lambda x_i - \log \Gamma(\alpha))$$

$$= n\alpha \log \lambda + (\alpha - 1) \sum_{i=1}^n \log x_i - \lambda \sum_{i=1}^n x_i$$

$$- n \log \Gamma(\alpha)$$

$$\frac{\partial l}{\partial \lambda} = \frac{n\alpha}{\lambda} - \sum_{i=1}^n x_i \quad \left(\frac{\partial l}{\partial \lambda} \stackrel{\text{set}}{=} 0 \right)$$

$$0 = \frac{n\hat{\alpha}}{\hat{\lambda}} - \sum_{i=1}^n x_i$$

$$\hat{\lambda} = \frac{\hat{\alpha}}{\bar{X}}.$$

$$\frac{\partial l}{\partial \alpha} = n \log \lambda + \sum_{i=1}^n \log x_i - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} \quad \left(\frac{\partial l}{\partial \alpha} \stackrel{\text{set}}{=} 0 \right)$$

$$0 = n \log \left(\frac{\hat{\alpha}}{\bar{X}} \right) + \sum_{i=1}^n \log x_i - n \frac{\Gamma'(\hat{\alpha})}{\Gamma(\hat{\alpha})}$$

$$0 = n \log \hat{\alpha} - n \log \bar{X} + \sum_{i=1}^n \log x_i - n \frac{\Gamma'(\hat{\alpha})}{\Gamma(\hat{\alpha})},$$

- no closed-form solution.
- solution has to be found either by numerical methods or bootstrap (later)
- more complications for checking the second derivatives.

Example 25

$$f(x) = \frac{1}{\theta}, \quad 0 \leq x \leq \theta$$

$$= \frac{1}{\theta} I_{(0 \leq x \leq \theta)}$$

$$l(\theta) = \sum_{i=1}^n -\log \theta, \quad x_i \leq \theta$$

$$= -n \log \theta, \quad x^{(n)} \leq \theta$$

$$\hat{\theta} = x^{(n)}.$$

- *Intuitively, this is clear.*
- *We know $f_{X^{(n)}}(x)$ for $X \sim \text{Uniform}(0, \theta)$.*
- *Compare to MoM:*

$$\mu_1 = \frac{\theta}{2}$$

$$\hat{\theta} = 2\bar{X}.$$

Example 26 (*Multinomial*(p_1, \dots, p_m)) :

$$\sum_{i=1}^m p_i = 1, \quad \sum_{i=1}^m x_i = n$$

$$f(x_1, \dots, x_m) = \frac{n!}{x_1! \dots x_m!} p_1^{x_1} \dots p_m^{x_m}$$

$$l(p_1, \dots, p_m) = \log n! - \sum_{i=1}^m \log x_i! + \sum_{i=1}^m x_i \log p_i$$

Using Lagrange multiplier

$$L(p_1, \dots, p_m, \lambda) = \log n! - \sum_{i=1}^m \log x_i! + \sum_{i=1}^m x_i \log p_i \\ + \lambda \left(\sum_{i=1}^m p_i - 1 \right)$$

$$\frac{\partial L}{\partial p_i} = \frac{x_i}{p_i} + \lambda \quad \left(\frac{\partial L}{\partial p_i} \stackrel{set}{=} 0 \right)$$

$$\hat{p}_i = \frac{-x_i}{\lambda},$$

$$1 = \sum_i \hat{p}_i = \sum_{i=1}^m \frac{-x_i}{\lambda} = \frac{-n}{\lambda},$$

$$\lambda = -n,$$

$$\hat{p}_i = \frac{x_i}{n} \quad (\text{intuitive})$$

- A special case is *Binomial* (n, p) , where $m = 2$, $p_1 = p$, $x_1 = x$, n is known

$$\hat{p} = \frac{x}{n},$$

- n above is a parameter; the number of observations is l , which is the vector (x_1, \dots, x_m)

For K observations: $(x_{11}, \dots, x_{1m}), \dots, (x_{K1}, \dots, x_{Km})$.

$$f(x_1, \dots, x_K) = \prod_{k=1}^K \frac{n!}{x_{k1}! \dots x_{km}!} p_1^{x_{k1}} \dots p_m^{x_{km}}$$

$$L(p_1, \dots, p_m, \lambda) = \log(n!)^K - \sum_{i=1}^m \sum_{k=1}^K \log x_{ki}! \\ + \sum_{i=1}^m \sum_{k=1}^K x_{ki} \log p_i + \lambda \left(\sum_{i=1}^m p_i - 1 \right)$$

$$\frac{\partial L}{\partial p_i} = \frac{\sum_{k=1}^K x_{ki}}{p_i} + \lambda,$$

$$\hat{p}_i = \frac{-\sum_{k=1}^K x_{ki}}{\lambda}$$

$$1 = \frac{-\sum_{i=1}^m \sum_{k=1}^K x_{ki}}{\lambda} = \frac{-nK}{\lambda}$$

$$\hat{p}_i = \frac{\sum_{k=1}^K x_{ki}}{nK} = \frac{\overline{X}_i}{n},$$

which for Binomial (n, p) will be

$$\hat{p} = \frac{\overline{X}}{n},$$

which is very intuitive.

8.3.1 Large Sample Theory for MLE

Reminder:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i \quad (\bar{X})$$

$$\hat{\mu} \xrightarrow{p} E[X] \quad (\text{WLLN})$$

$$\sqrt{n} \frac{\hat{\mu} - \mu}{\sigma} \xrightarrow{d} N(0, 1) \quad (\text{CLT})$$

$$\lim_{n \rightarrow \infty} \Pr \left(\sqrt{n} \frac{\hat{\mu} - \mu}{\sigma} \leq x \right) = \Pr(N(0, 1) \leq x)$$

$$\lim_{n \rightarrow \infty} \Pr(\sqrt{n}(\hat{\mu} - \mu) \leq \sigma x) = \Pr(\sigma N(0, 1) \leq \sigma x)$$

$$= \Pr(N(0, \sigma^2) \leq \sigma x)$$

$$\sqrt{n}(\hat{\mu} - \mu) \xrightarrow{d} N(0, \sigma^2) \quad (\text{CLT}')$$

Definition 27 (Asymptotic Mean and Variance)

: For any statistic (or estimator) T_n , if

$$k_n \frac{T_n - \mu}{\sigma} \xrightarrow{d} N(0, 1), \quad (k_n \text{ can be } \sqrt{n})$$

we call μ and σ^2 the asymptotic mean and variance (even if $E[T_n] \neq \mu$ and $\text{Var}[T_n] \neq \sigma^2$).

MoM:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i \quad (\bar{X})$$

$$\hat{\mu} \xrightarrow{p} E[X] \quad (\text{WLLN})$$

$$\sqrt{n} \frac{\hat{\mu} - E[X]}{\sqrt{\text{Var}[X]}} \xrightarrow{d} N(0, 1) \quad (\text{CLT})$$

$$\hat{\mu}_r = \frac{1}{n} \sum_{i=1}^n X_i^r, \quad (\text{MoM})$$

$$\hat{\mu}_r \xrightarrow{p} E[X^r] \quad (E[\hat{\mu}_r] \overset{\text{always}}{=} E[X^r])$$

$$\sqrt{n} \frac{\hat{\mu}_r - E[X^r]}{\sqrt{\text{Var}[X^r]}} \xrightarrow{d} N(0, 1)$$

Notice that:

- $E[\hat{\mu}_r] = E[X^r]$ (always unbiased $\forall n$)
- the estimated parameters, e.g., $\hat{\sigma}^2$, may be biased for finite n .

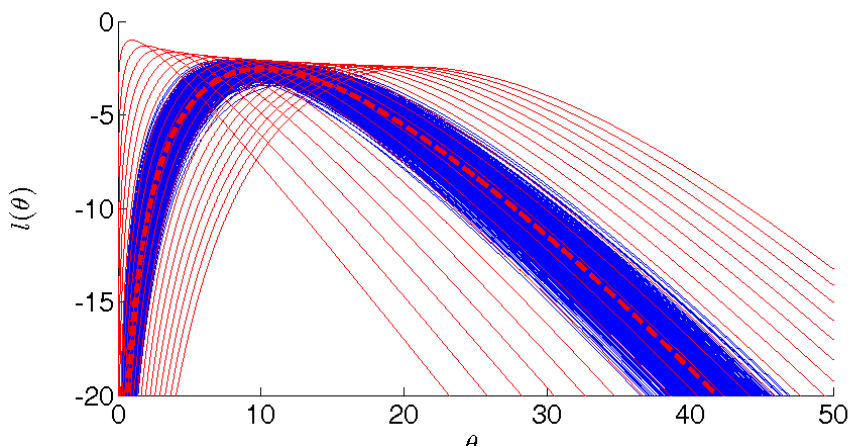
Some Intuition First:

$$l(\theta|X) = X \log \theta - \theta - \log(X!)$$

$$E[l(\theta|X)] = E[X] \log \theta - \theta - E[\log(X!)]$$

$$l(\theta|X_1, \dots, X_n) = \sum_i X_i \log \theta - n\theta - \sum_i \log(X_i!)$$

$$\frac{1}{n} l(\theta) \xrightarrow{p} E[\log f(X|\theta)]$$



- **Take care:** $E[X]$ above is $E_{X|\theta_0}[X]$.
- Why curves are less than zero?
- We simulated 1000 curves, why few are there

Matlab Code 8.2:

```
theta0=10; theta = (0:.01:50) ' ;  
C = 1000;  
ltheta = zeros(length(theta) , C);  
  
figure1 = figure; fs=20;  
set(gcf, 'Units', 'inches');  
haxes=axes('Parent',figure1,'YLim'  
    ,[-20 0], 'XLim',[0 50], 'FontSize',  
    fs);  
xlabel(' $\theta$ ', 'Interpreter', 'latex'  
    , 'FontSize', fs, 'Units', '  
    normalized');  
ylabel(' $l(\theta)$ ', 'Interpreter', '  
    latex', 'FontSize', fs, 'Units', '  
    normalized');  
  
hold all;
```

```

n=10;
for c=1:C
    x=random( 'Poisson' ,theta0 ,[n,1]);
    ltheta (:, c)=mean(x)*log(theta)-
        theta-sum(log(factorial(x)))/n;
    plot(theta, ltheta (:, c), 'b');
end;

```

```

n=1;
for c=1:C
    x=random( 'Poisson' ,theta0 ,[n,1]);
    ltheta (:, c)=x*log(theta)-theta-
        sum(log(factorial(x)))/n;
    plot(theta, ltheta (:, c), 'r');
end;
plot(theta, mean(ltheta, 2), 'r—', '
    LineWidth', 4);

```

Theorem 28 *Under regularity conditions on f , the MLE estimator is consistent*

Semi-Proof. :Under regularity conditions

$$l(\theta) = \sum_{i=1}^n \log f(X_i|\theta),$$

$$\frac{1}{n} l(\theta) \xrightarrow{p} E[\log f(X|\theta)], \quad (E_{X|\theta_0})$$

$$\operatorname{argmax} l(\theta) = \operatorname{argmax} \frac{1}{n} l(\theta) \quad (\text{of course})$$

$$\stackrel{I \text{ hope}}{=} \operatorname{argmax} E[\log f(X|\theta)]$$

$$\frac{\partial}{\partial \theta} E[\log f(X|\theta)] = \frac{\partial}{\partial \theta} \int \log f(x|\theta) f(x|\theta_0) dx$$

$$= \int \frac{\partial}{\partial \theta} \log f(x|\theta) f(x|\theta_0) dx$$

$$= \int \frac{\frac{\partial}{\partial \theta} f(x|\theta)}{f(x|\theta)} f(x|\theta_0) dx$$

$$\left. \frac{\partial}{\partial \theta} E[\log f(X|\theta)] \right|_{\theta_0} = \left. \int \frac{\partial}{\partial \theta} f(x|\theta) dx \right|_{\theta_0}$$

$$= \left. \frac{\partial}{\partial \theta} \int f(x|\theta) dx \right|_{\theta_0}$$

$$= \left. \frac{\partial}{\partial \theta} 1 \right|_{\theta_0} = 0$$

■

Lemma 29 *Under regularity conditions:*

$$\mathbb{E} \left[\frac{\partial}{\partial \theta} \log f(X|\theta) \right] = 0 \quad (\mathbb{E}_{X|\theta})$$

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

which is called $I(\theta)$, the Fisher information (information number) of one observation.

- What is the meaning of “Information” here?
Let’s see on the figure.
- Meaning of both equations.

Proof.

$$f(x|\theta) \frac{\partial}{\partial \theta} \log f(x|\theta) = f(x|\theta) \frac{\frac{\partial}{\partial \theta} f(x|\theta)}{f(x|\theta)} = \frac{\partial}{\partial \theta} f(x|\theta)$$

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} (1) = \frac{\partial}{\partial \theta} \int f(x|\theta) dx = \int \frac{\partial}{\partial \theta} f(x|\theta) dx \\ &= \int f(x|\theta) \frac{\partial}{\partial \theta} \log f(x|\theta) dx && (E_{X|\theta_0}) \\ &= \frac{\partial}{\partial \theta} \int f(x|\theta) \frac{\partial}{\partial \theta} \log f(x|\theta) dx \\ &= \int \frac{\partial}{\partial \theta} f(x|\theta) \frac{\partial}{\partial \theta} \log f(x|\theta) dx + \\ &\quad \int f(x|\theta) \frac{\partial^2}{\partial \theta^2} \log f(x|\theta) dx \\ &= \int f(x|\theta) \left(\frac{\partial}{\partial \theta} \log f(x|\theta) \right)^2 dx + \\ &\quad \int f(x|\theta) \left(\frac{\partial^2}{\partial \theta^2} \log f(x|\theta) \right) dx \\ &= E \left[\left(\frac{\partial}{\partial \theta} \log f(x|\theta) \right)^2 \right] + E \left[\frac{\partial^2}{\partial \theta^2} \log f(x|\theta) \right] \end{aligned}$$

■

Theorem 30 Let $X_1, \dots, X_n \stackrel{iid}{\sim} f(X|\theta)$, $\hat{\theta}$ is the MLE of θ . Then, under regularity conditions

$$\sqrt{n} \frac{\hat{\theta} - \theta}{1/\sqrt{I(\theta)}} \xrightarrow{d} N(0, 1),$$

$$\sqrt{n} \frac{\tau(\hat{\theta}) - \tau(\theta)}{1/\sqrt{I(\theta)}} \xrightarrow{d} N(0, 1).$$

That is, any estimator $\tau(\hat{\theta})$ (or $\hat{\theta}$) is asymptotically unbiased for $\tau(\theta)$ (or θ) with asymptotic variance of $1/I(\theta)$. So, we have $\xrightarrow{d} N(0, 1)$ in addition to $\xrightarrow{p} \theta$.

Proof. Suppose that the true value of θ is θ_0

$$l(\theta) = \sum_{i=1}^n \log f(X_i|\theta)$$

$$l'(\theta) = l'(\theta_0) + (\theta - \theta_0) l''(\theta_0) + \dots$$

$$l'(\hat{\theta}) = l'(\theta_0) + (\hat{\theta} - \theta_0) l''(\theta_0) + \dots$$

$$(\hat{\theta} - \theta_0) \approx -l'(\theta_0) / l''(\theta_0) \quad (\text{MLE def.})$$

$$\sqrt{n} \frac{(\hat{\theta} - \theta_0)}{\sqrt{1/I(\theta_0)}} \approx \frac{\sqrt{n} \frac{1}{n} l'(\theta_0) / \sqrt{I(\theta_0)}}{\frac{-1}{n} l''(\theta_0) / I(\theta_0)}.$$

$$\frac{1}{n} l'(\theta_0) = \frac{1}{n} \sum_i \frac{\partial}{\partial \theta} \log f(X_i | \theta) \Big|_{\theta_0}$$

$$E \left[\frac{\partial}{\partial \theta} \log f(X_i | \theta) \Big|_{\theta_0} \right] = 0 \quad (E_{X|\theta_0})$$

$$\begin{aligned} \text{Var} \left[\frac{\partial}{\partial \theta} \log f(X_i | \theta) \Big|_{\theta_0} \right] &= E \left[\left(\frac{\partial}{\partial \theta} \log f(X | \theta) \right)^2 \Big|_{\theta_0} \right] \\ &= I(\theta_0) \end{aligned}$$

$$\sqrt{n} \frac{\frac{1}{n} l'(\theta_0) - 0}{\sqrt{I(\theta_0)}} \xrightarrow{d} N(0, 1) \quad (\text{CLT})$$

$$-\frac{1}{n} l''(\theta_0) = -\frac{1}{n} \sum_i \frac{\partial^2}{\partial \theta^2} \log f(X_i | \theta)$$

$$E \left[\frac{\partial^2}{\partial \theta^2} \log f(X | \theta) \Big|_{\theta_0} \right] = -I(\theta_0)$$

$$-\frac{1}{n} l''(\theta_0) \xrightarrow{p} I(\theta_0)$$

$$-\frac{1}{n} l''(\theta_0) / I(\theta_0) \xrightarrow{p} 1$$

$$\sqrt{n} \frac{(\hat{\theta} - \theta_0)}{\sqrt{1/I(\theta_0)}} \xrightarrow{d} N(0, 1).$$

■

Said differently

$$\sqrt{n} \frac{\hat{\theta} - \theta_0}{\sqrt{1/I(\theta_0)}} \xrightarrow{d} N(0, 1),$$
$$\sqrt{n} (\hat{\theta} - \theta_0) \xrightarrow{d} N(0, 1/I(\theta_0)),$$

which means that the MLE $\hat{\theta}$

- Asymptotically unbiased
- Asymptotic variance = $1/I(\theta_0)$
- Asymptotically normally distributed.

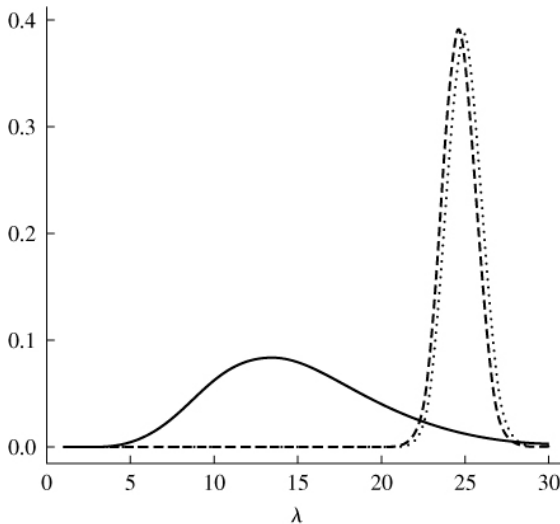
Why variance decreases with $I(\theta_0)$?

$$I(\theta_0) = -E \left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta) \Big|_{\theta_0} \right]$$

High $I(\theta_0)$ means very sharp curve at θ_0 , which means very probable θ_0 , which means less likely that the next dataset will not support that inference; and hence less variable the next estimator is.

8.4 The Bayesian Approach to Parameter Estimation

- We treat θ as r.v. with **subjective** prior knowledge f_{θ} ; as opposed to “Frequentist (or Classical) Approach”
- Data $\mathbf{x} = x_1, \dots, x_n$ for $\mathbf{X} = X_1, \dots, X_n$ modifies our belief and produces the posterior $f_{\theta|\mathbf{X}}$?
- We estimate θ by many criteria; e.g.,:



1. Posterior Mode/Max. A Posteriori (MAP):

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x})$$

2. Posterior Mean:

$$\hat{\theta} = \underset{\Theta}{\operatorname{E}}[\theta] = \int \theta f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta$$

3. Posterior loss function optimization:

$$\begin{aligned}\hat{\theta} &= \underset{\eta}{\operatorname{argmin}} \underset{\Theta}{\operatorname{E}}[L(\eta, \theta)] \\ &= \underset{\eta}{\operatorname{argmin}} \int L(\eta, \theta) f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta\end{aligned}$$

General Framework:

$$\begin{aligned}f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) &= \frac{f_{\mathbf{X},\Theta}(\mathbf{x}, \theta)}{f_{\mathbf{X}}(\mathbf{x})} \\ &= \frac{f_{\mathbf{X}|\Theta}(\mathbf{x}|\theta) f_{\Theta}(\theta)}{\int f_{\mathbf{X},\Theta}(\mathbf{x}, \theta) d\theta} \\ &= \frac{f_{\mathbf{X}|\Theta}(\mathbf{x}|\theta) f_{\Theta}(\theta)}{\int f_{\mathbf{X}|\Theta}(\mathbf{x}|\theta) f_{\Theta}(\theta) d\theta} \\ &= \operatorname{Const}(\mathbf{x}) f_{\mathbf{X}|\Theta}(\mathbf{x}|\theta) f_{\Theta}(\theta)\end{aligned}$$

Posterior \propto Likelihood \times Prior.

Connection to MLE:

$$f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) = \text{Const}(\mathbf{x}) f_{\mathbf{X}|\Theta}(\mathbf{x}|\theta) f_{\Theta}(\theta) \\ \propto \text{Likelihood} \times \text{Prior}$$

if we choose an uninformative prior $\Theta \sim U$ to let data speak for themselves:

$$f_{\Theta|X}(\theta|x) = \text{Const}(x) f_{X|\Theta}(x|\theta) \\ \propto \text{Likelihood}$$

Then, if we choose MAP criterion

$$\hat{\theta} = \operatorname{argmax} l(\theta), \quad (\text{MLE})$$

Example 31 (Poisson) \mathbf{X} denotes X_1, \dots, X_n :

$$f_{\mathbf{X}|\Lambda} = \prod_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}, \quad 0 \leq x_i,$$

$$= \frac{\lambda^{\sum_i x_i} e^{-n\lambda}}{\prod_i x_i!}$$

$$f_{\Lambda|\mathbf{X}} = \frac{f_{\mathbf{X}|\Lambda}(\mathbf{x}|\lambda) f_{\Lambda}(\lambda)}{\int f_{\mathbf{X}|\Lambda}(\mathbf{x}|\lambda) f_{\Lambda}(\lambda) d\lambda}$$

$$= \frac{\lambda^{\sum_i x_i} e^{-n\lambda} f_{\Lambda}(\lambda) / \prod_i x_i!}{\int \lambda^{\sum_i x_i} e^{-n\lambda} f_{\Lambda}(\lambda) / \prod_i x_i! d\lambda}$$

$$= \frac{\lambda^{\sum_i x_i} e^{-n\lambda} \frac{1}{100}}{\int \lambda^{\sum_i x_i} e^{-n\lambda} \frac{1}{100} d\lambda} \quad (\Lambda \sim U(0, 100))$$

$$= \frac{v^{\alpha}}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-v\lambda} \quad (Gamma(\alpha, v))$$

$$\sim Gamma(S_n + 1, n)$$

$$\hat{\lambda} = E[\Lambda] = \frac{S_n + 1}{n} = \bar{X} + \frac{1}{n} \quad (\text{Post. Mean})$$

$$\frac{\partial f_{\Lambda|\mathbf{X}}}{\partial \lambda} = \frac{v^{\alpha}}{\Gamma(\alpha)} ((\alpha - 1) \lambda^{\alpha-2} e^{-v\lambda} - v \lambda^{\alpha-1} e^{-v\lambda})$$

$$\hat{\lambda} = \frac{\alpha - 1}{v} = \frac{S_n}{n} = \bar{X} \quad (\text{MAP} \equiv \text{MLE})$$

$$\frac{S_n}{n} = \frac{573}{23} = 24.9, \quad \frac{S_n + 1}{n} = 25$$

On the other hand, if we have the prior knowledge that Λ has $\mu = 15$ and $\sigma = 5$ then, we can assume that $\Lambda \sim \text{Gamma}(\alpha, \nu)$ with

$$\mu = \alpha / \nu,$$

$$\sigma^2 = \alpha / \nu^2,$$

$$\nu = \frac{\mu}{\sigma^2} = 0.6 \ll n \quad (n = 23)$$

$$\alpha = \nu \mu = 9 \ll S_n, \quad (S_n = 573)$$

$$\begin{aligned} f_{\Lambda|\mathbf{X}} &= \frac{\lambda^{\sum_i x_i} e^{-n\lambda} f_{\Lambda}(\lambda)}{\int \lambda^{\sum_i x_i} e^{-n\lambda} f_{\Lambda}(\lambda) d\lambda} \\ &= \frac{\lambda^{\sum_i x_i} e^{-n\lambda} \frac{\nu^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\nu\lambda}}{\int \lambda^{\sum_i x_i} e^{-n\lambda} \frac{\nu^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\nu\lambda} d\lambda} \\ &= \frac{\lambda^{(S_n+\alpha-1)} e^{-(n+\nu)\lambda}}{\int \lambda^{(S_n+\alpha-1)} e^{-(n+\nu)\lambda} d\lambda} \\ &\sim \text{Gamma}(S_n + \alpha, n + \nu) \\ \hat{\lambda} &= \frac{S_n + \alpha}{n + \nu} = \frac{573 + 9}{23 + .6} = 24.7 \quad (\text{Post. Mean}) \\ \hat{\lambda} &= \frac{S_n + \alpha - 1}{n + \nu} = \frac{573 + 9 - 1}{23 + .6} = 24.6 \quad (\text{MAP}) \end{aligned}$$

Example 32 ($Ber(p)$) : n obs., then

$$\mu_1 = p,$$

$$\hat{p} = \overline{X} = \frac{\sum_i x_i}{n} = \frac{\#Heads}{n}, \quad (\text{MoM})$$

$$p_X(x) = p^x (1-p)^{1-x}, \quad x = 0, 1$$

$$l(p) = \sum_i x_i \log p + \sum_i (1-x_i) \log(1-p)$$

$$l'(p) = \frac{\sum_i x_i}{p} - \frac{\sum_i (1-x_i)}{1-p} \quad (l'(p) \stackrel{set}{=} 0)$$

$$\hat{p} = \overline{X} = \frac{\sum_i x_i}{n} = \frac{\#Heads}{n}. \quad (\text{MLE})$$

Now, if we get 5 heads in 5 trials \hat{p} will be 1 !!!!

Let's see the Bayesian approach.

$$f_{\mathbf{X}|P} = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum_i x_i} (1-p)^{\sum_i (1-x_i)}$$

$$f_P(p) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} p^{a-1} (1-p)^{b-1} \quad (\sim \text{Beta}(a, b))$$

$$\begin{aligned} f_{P|\mathbf{X}} &= \frac{f_{\mathbf{X}|P}(\mathbf{x}|p) f_P(p)}{\int f_{\mathbf{X}|P}(\mathbf{x}|p) f_P(p) dp} \\ &\propto p^{a-1+S} (1-p)^{b-1+(n-S)} \\ &\sim \text{Beta}(a+S, b+n-S). \end{aligned}$$

$$\hat{p} = \frac{A-1}{A+B-2} = \frac{a+S-1}{a+b+n-2} \quad (\text{MAP})$$

$$= \frac{a+S-1}{2a+n-2} \quad (\text{Symmetric Prior})$$

$$a=1: U(0,1), \hat{p} = \frac{S}{n} \equiv \text{MLE}.$$

$$a=2: \text{not uniform but spread. } \hat{p} = (S+1)/(n+2).$$

- $S=n: \hat{p} = (n+1)/(n+2) \rightarrow 1.$
- $S=n/2: \hat{p} = 1/2$ (of course).

$$a \gg: \text{insisting on fair coin, } \hat{p} \approx a/(2a) = \frac{1}{2}$$

$$f_{P|X} \sim \text{Beta}(a + S, b + n - S)$$

$$\begin{aligned}\hat{p} &= \frac{A}{A + B} \\ &= \frac{a + S}{a + b + n}\end{aligned}\quad (\text{Posterior Mean})$$

8.4.1 Large Sample Theory of Bayesian Inference

\mathbf{X} and \mathbf{x} denote X_1, \dots, X_n and x_1, \dots, x_n , respectively, to simplify notation.

$$f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) \propto f_{\Theta}(\theta) f_{\mathbf{X}|\Theta}(\mathbf{x}|\theta),$$

which is dominated by $f_{\mathbf{X}|\Theta}$ as $n \rightarrow \infty$.

$$\begin{aligned} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) &\propto f_{\mathbf{X}|\Theta}(\mathbf{x}|\theta) && (\text{as } n \rightarrow \infty) \\ &= \exp [\log f_{\mathbf{X}|\Theta}(\mathbf{x}|\theta)] \\ &= \exp [l(\theta)] \\ &= \exp [l(\hat{\theta}) + (\theta - \hat{\theta}) l'(\hat{\theta}) \\ &\quad + \frac{1}{2} (\theta - \hat{\theta})^2 l''(\hat{\theta}) + \dots] \\ &\propto \exp \left[-\frac{1}{2} \frac{(\theta - \hat{\theta})^2}{1/l''(\hat{\theta})} \right] && (l'(\hat{\theta}) = 0) \\ &\sim N(\hat{\theta}, -1/l''(\hat{\theta})). \end{aligned}$$

Do not confuse it with the MLE asymptotic normality.

8.5 Assessing Estimators, Efficiency, and the Cramér-Rao Lower Bound

8.5.1 Mean Squared Error (MSE) Criterion

$$\begin{aligned}MSE(\hat{\theta}) &= \mathbb{E}_{\mathbf{X}} \left[(\hat{\theta} - \theta)^2 \right] \\&= \text{Var}_{\mathbf{X}} [\hat{\theta}] + \left(\mathbb{E}_{\mathbf{X}} \hat{\theta} - \theta \right)^2 \\&= \text{Variance}(\hat{\theta}) + \left(\text{Bias}(\hat{\theta}) \right)^2.\end{aligned}$$

- Since $MSE = MSE(\theta)$ no best estimator; e.g., $\hat{\theta} = 12.3$ is the best when $\theta = 12.3$ but terrible otherwise.
- If $\text{Bias}(\hat{\theta}) = 0$, $\hat{\theta}$ is unbiased for θ .
- Tradeoff exists between Bias and Variance.
- A biased estimator may have lower MSE.

Example 33 ($\hat{\sigma}^2$ vs. S^2 for $N(\mu, \sigma^2)$) :

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i \left(X_i - \bar{X} \right)^2,$$

$$S^2 = \frac{1}{n-1} \sum_i \left(X_i - \bar{X} \right)^2$$

$$E[S^2] = \sigma^2 \quad (\text{unbiased})$$

$$\text{Var}[S^2] = \frac{2\sigma^4}{n-1} \quad (\text{see Extra Materials})$$

$$MSE(S^2) = \frac{2\sigma^4}{n-1} + (\sigma^2 - \sigma^2)^2 = \frac{2\sigma^4}{n-1}$$

$$E[\hat{\sigma}^2] = \frac{n-1}{n} \sigma^2 \quad (\text{biased})$$

$$\text{Var}[\hat{\sigma}^2] = \text{Var}\left[\frac{n-1}{n} S^2\right] = \left(\frac{n-1}{n}\right)^2 \text{Var}[S^2]$$

$$= \left(\frac{n-1}{n}\right)^2 \left(\frac{2\sigma^4}{n-1}\right) = \frac{2(n-1)\sigma^4}{n^2}$$

$$MSE(\hat{\sigma}^2) = \frac{2(n-1)\sigma^4}{n^2} + \left(\frac{n-1}{n}\sigma^2 - \sigma^2\right)^2$$

$$= \frac{2n-1}{n^2} \sigma^4 < \frac{2\sigma^4}{n-1} \quad \forall \sigma, n.$$

Remarks:

- Although S^2 is unbiased, $\hat{\sigma}^2$ has less MSE.
- MSE, for scale parameter, may not be reasonable since $\sigma^2 > 0$.
- $\hat{\theta}_1$ may be better than $\hat{\theta}_2$ under some criterion and the other way around and another criterion.

Example 34 (\hat{p} of $Ber(p)$) :

$$\hat{p}_M = \bar{X} \quad (\text{MLE})$$

$$E[\hat{p}_M] = p$$

$$\text{Var}[\hat{p}_M] = \frac{1}{n}p(1-p)$$

$$MSE(\hat{p}_M) = \frac{1}{n}p(1-p)$$

$$\hat{p}_B = \frac{S + a}{a + b + n} \quad (\text{Posterior Mean})$$

$$E[\hat{p}_B] = \frac{np + a}{a + b + n}$$

$$\text{Var}[\hat{p}_B] = \frac{np(1-p)}{(a + b + n)^2}$$

$$MSE(\hat{p}_B) = \frac{np(1-p)}{(a + b + n)^2} + \left(\frac{np + a}{a + b + n} - p \right)^2$$

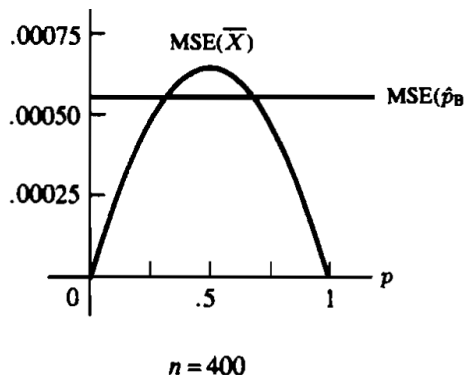
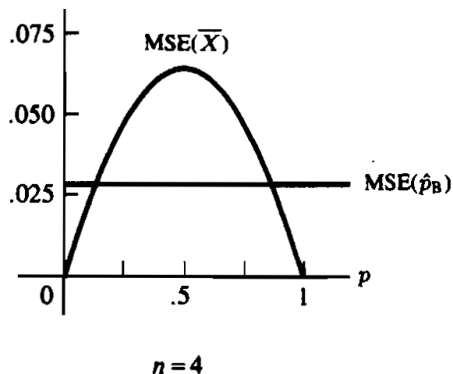
Choosing $a = b = \sqrt{n}/2$ relaxes dependence on p :

$$\hat{p}_B = \frac{S + \sqrt{n}/2}{n + \sqrt{n}},$$

$$MSE(\hat{p}_B) = \frac{n}{4(n + \sqrt{n})^2}.$$

$$MSE(\hat{p}_M) = \frac{1}{n}p(1-p)$$

$$MSE(\hat{p}_B) = \frac{n}{4(n + \sqrt{n})^2}$$



- For small n , \hat{p}_B is better unless p is on the boundary.
- For large n , \hat{p}_M is better unless p is in the middle.
- Having knowledge about the problem allows choosing the right estimator.

8.5.2 Best Unbiased Estimator

Definition 35 (UMVUE) : An estimator $\hat{\theta}^*$, for θ , is a best unbiased estimator or uniform minimum variance unbiased estimator (UMVUE) if it satisfies $E[\hat{\theta}^*] = \theta \forall \theta$ and for any other estimator $\hat{\theta}$ we have $\text{Var}[\hat{\theta}^*] \leq \text{Var}[\hat{\theta}]$.

Theorem 36 (Cramér-Rao Inequality) : Let $X_1, \dots, X_n \stackrel{i.i.d}{\sim} f(x|\theta)$ with regularity condition. Then for any estimator $T = T(X_1, \dots, X_n) = T(\mathbf{X})$

$$\text{Var}(T) \geq \frac{\left(\frac{d}{d\theta} E[T]\right)^2}{nI(\theta)},$$

$$\text{Var}(T) \geq \frac{1}{nI(\theta)}. \quad (\text{if } T \text{ is unbiased})$$

- For all estimators with particular bias: the higher the *information number* the lower the *lower bound*.
- An estimator *attains (attainment)* the lower bound is called *efficient*.

Proof. : Since $1 \leq \rho = \text{Cov}(T, Z) / \sqrt{\text{Var}(T) \text{Var}(Z)}$

$$\text{Var}[T] \geq (\text{Cov}(T, Z))^2 / \text{Var}(Z)$$

$$Z = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(X_i | \theta)$$

$$\begin{aligned} \text{Var}[Z] &= n \text{Var} \left[\frac{\partial}{\partial \theta} \log f(X_i | \theta) \right] \\ &= n I(\theta) \quad (\text{Proof of Th. 30}) \end{aligned}$$

$$\begin{aligned} \sigma_{TZ} &= E(Z - E[Z])(T - E[T]) = E[T(Z - E[Z])] \\ &= E[ZT] \quad (E[Z] = 0) \end{aligned}$$

$$\begin{aligned} &= E \left[T \frac{\partial}{\partial \theta} \log \prod_i f(X_i | \theta) \right] \\ &= E \left[T \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right] \quad (\mathbf{X} = X_1, \dots, X_n) \end{aligned}$$

$$= \int T(\mathbf{x}) \frac{\frac{\partial}{\partial \theta} f(\mathbf{x} | \theta)}{f(\mathbf{x} | \theta)} f(\mathbf{x} | \theta) d\mathbf{x}$$

$$= \frac{\partial}{\partial \theta} \int T(\mathbf{x}) f(\mathbf{x} | \theta) d\mathbf{x}$$

$$= \frac{\partial}{\partial \theta} E_{\mathbf{X}}[T(\mathbf{X})]$$

■

Example 37 (Poisson) :

$$\begin{aligned} I(\lambda) &= \mathbb{E} \left[\left(\frac{\partial}{\partial \lambda} \log \frac{\lambda^X e^{-\lambda}}{X!} \right)^2 \right] \\ &= \mathbb{E} \left[\left(\frac{\partial}{\partial \lambda} (X \log \lambda - \lambda - \log X!) \right)^2 \right] \\ &= \mathbb{E} \left[\left(\frac{X}{\lambda} - 1 \right)^2 \right] \\ &= -\mathbb{E} \left[\frac{\partial^2}{\partial \lambda^2} \log \frac{\lambda^X e^{-\lambda}}{X!} \right] \quad (\text{easier}) \\ &= -\mathbb{E} \left[\frac{-X}{\lambda^2} \right] = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda}, \end{aligned}$$

$$\begin{aligned} \text{Var}[T] &\geq \frac{\left(\frac{\partial}{\partial \lambda} \mathbb{E}[T] \right)^2}{n I(\lambda)} \\ &= \frac{\lambda}{n} \quad (\text{for unbiased estimators}) \end{aligned}$$

$$\hat{\lambda} = \bar{X} \quad (\text{MLE})$$

$$\mathbb{E}[\hat{\lambda}] = \lambda \quad (\text{unbiased})$$

$$\text{Var}[\hat{\lambda}] = \text{Var}[\bar{X}] = \frac{1}{n} \text{Var}[X] = \frac{\lambda}{n}, \quad (\text{attainment})$$

Example 38 ($U(0, \theta)$) : $f(x|\theta) = 1/\theta$, then

$$\begin{aligned} I(\theta) &= E \left[\left(\frac{\partial}{\partial \theta} \log(1/\theta) \right)^2 \right] \\ &= E \left[\left(-\frac{\partial}{\partial \theta} \log \theta \right)^2 \right] = 1/\theta^2, \end{aligned}$$

$$\begin{aligned} \text{Var} [\hat{\theta}] &\geq \frac{\left(\frac{\partial}{\partial \theta} E[T] \right)^2}{n I(\theta)} \\ &= \frac{\theta^2}{n}, \end{aligned} \quad (\text{for unbiased estimators})$$

$$\hat{\theta} = 2\bar{X}, \quad (\text{MoM})$$

$$E[\hat{\theta}] = \theta \quad (\text{unbiased})$$

$$\begin{aligned} \text{Var} [\hat{\theta}] &= \frac{4}{n} \text{Var} [X] = \frac{4}{n} \frac{\theta^2}{12} \\ &= \frac{\theta^2}{3n} < \frac{\theta^2}{n}. \quad (!!!\text{where is the problem?}) \end{aligned}$$

The regularity condition assumes ($n = 1$):

$$\begin{aligned}\frac{\partial}{\partial \theta} \mathbb{E}[T] &= \frac{\partial}{\partial \theta} \int T f(x|\theta) dx & (\mathbf{x} = x) \\ &= \int T \frac{\partial}{\partial \theta} f(x|\theta) dx\end{aligned}$$

Let's see

$$\begin{aligned}\frac{\partial}{\partial \theta} \mathbb{E}[T] &= \frac{\partial}{\partial \theta} \int_0^\theta T \frac{1}{\theta} dx \\ &= \frac{\partial}{\partial \theta} \left(\frac{1}{\theta} \int_0^\theta T dx \right) \\ &= \left(\frac{\partial}{\partial \theta} \frac{1}{\theta} \right) \int_0^\theta T dx + \frac{1}{\theta} \frac{\partial}{\partial \theta} \int_0^\theta T dx \\ &= \left(\frac{\partial}{\partial \theta} \frac{1}{\theta} \right) \int_0^\theta T dx + \frac{T(\theta)}{\theta} \\ \int_0^\theta T \frac{\partial}{\partial \theta} f(x|\theta) dx &= \left(\frac{\partial}{\partial \theta} \frac{1}{\theta} \right) \int_0^\theta T dx, \\ &\neq \frac{\partial}{\partial \theta} \mathbb{E}[T],\end{aligned}$$

unless $T(\theta) = 0 \forall \theta$.

Homework: repeat with the MLE estimator, scale it to be unbiased, then find its variance.

Loss Function

- Not only for assessment and comparison,
- but also for designing and optimization!

The loss function:

$$L(\theta, T(\mathbf{X})) = |\theta - T(\mathbf{X})| \quad (\text{absolute error (AE)})$$

$$L(\theta, T(\mathbf{X})) = (\theta - T(\mathbf{X}))^2 \quad (\text{squared error (SE)})$$

\vdots

expresses how the estimate $T(\mathbf{X})$ deviates from θ .

The risk:

$$R(\theta, T) = \mathbb{E}_{\mathbf{X}} L(\theta, T(\mathbf{X}))$$

is a function of θ . $R(\theta, T_1)$ may cross with $R(\theta, T_2)$.

MSE (special case):

$$MSE(\theta) = R(\theta, T)$$

$$= \mathbb{E}_{\mathbf{X}} [L(\theta, T(\mathbf{X}))],$$

$$L(\theta, T(\mathbf{X})) = (\theta - T(\mathbf{X}))^2.$$

Example 39 (Risk of σ^2 Est.) :

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n \left(X_i - \bar{X} \right)^2, \quad (R(\sigma^2, S^2) = \frac{2\sigma^4}{n-1})$$

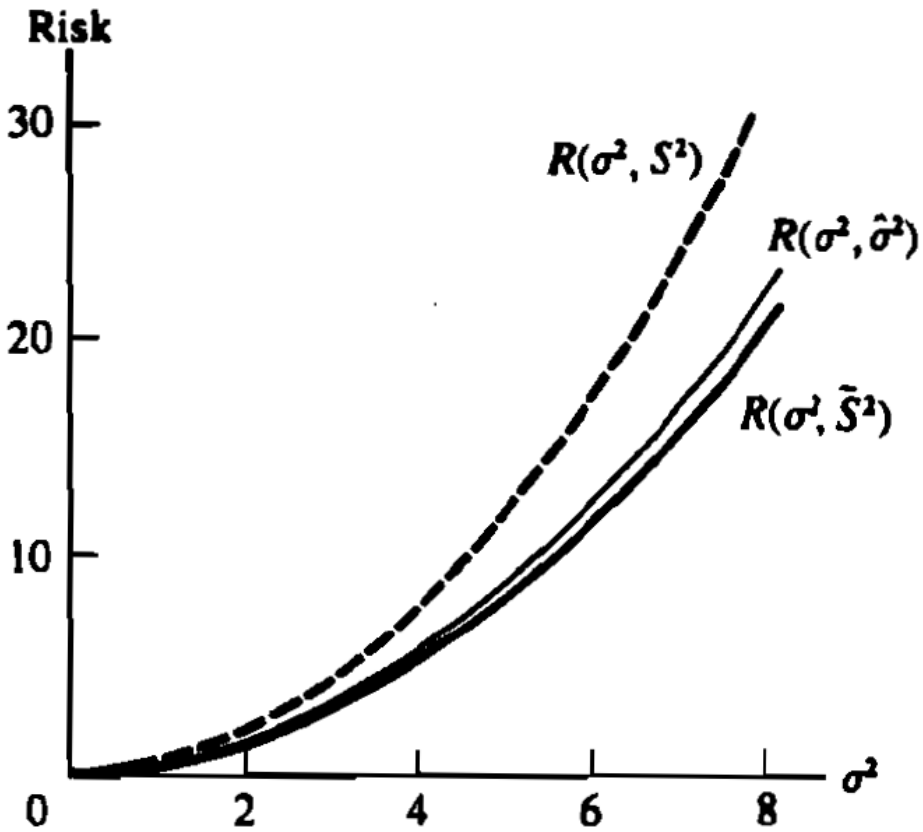
$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left(X_i - \bar{X} \right)^2, \quad (R(\sigma^2, \hat{\sigma}^2) = \frac{2n-1}{n^2} \sigma^4)$$

$$\tilde{S}^2 = b \sum_{i=1}^n \left(X_i - \bar{X} \right)^2 \quad (R(\sigma^2, \tilde{S}^2)?)$$

$$\begin{aligned} R(\sigma^2, \tilde{S}^2) &= \text{Var} [b(n-1) S^2] \\ &\quad + (\text{E} [b(n-1) S^2] - \sigma^2)^2 \\ &= b^2 (n-1)^2 \frac{2\sigma^4}{n-1} + (b(n-1) - 1)^2 \sigma^4 \\ &= (2b^2 (n-1) + (b(n-1) - 1)^2) \sigma^4, \\ &= c\sigma^4, \end{aligned}$$

$$c_{\min} = \frac{2}{n+1} \quad (\text{at } b = \frac{1}{n+1})$$

$$\begin{aligned} \tilde{S}^2 &= \frac{1}{n+1} \sum_{i=1}^n \left(X_i - \bar{X} \right)^2 \\ (R(\sigma^2, \tilde{S}^2) &= \frac{2}{n+1} \sigma^4) \end{aligned}$$



Connection to Cramér-Rao Inequality

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right)$$

$$l(\theta) = -\log \sqrt{2\pi} - \frac{1}{2} \log \theta - \frac{1}{2\theta} (x-\mu)^2$$

$(\theta = \sigma^2)$

$$l'(\theta) = \frac{-1}{2\theta} + \frac{(x-\mu)^2}{2\theta^2}$$

$$l''(\theta) = \frac{1}{2\theta^2} - \frac{(x-\mu)^2}{\theta^3}$$

$$E[l''(\theta)] = \frac{1}{2\theta^2} - \frac{\theta}{\theta^3} = \frac{-1}{2\theta^2}$$

$$I(\theta) = -E\left[\frac{\partial^2 l(\theta)}{\partial \theta^2}\right] = \frac{1}{2\sigma^4}$$

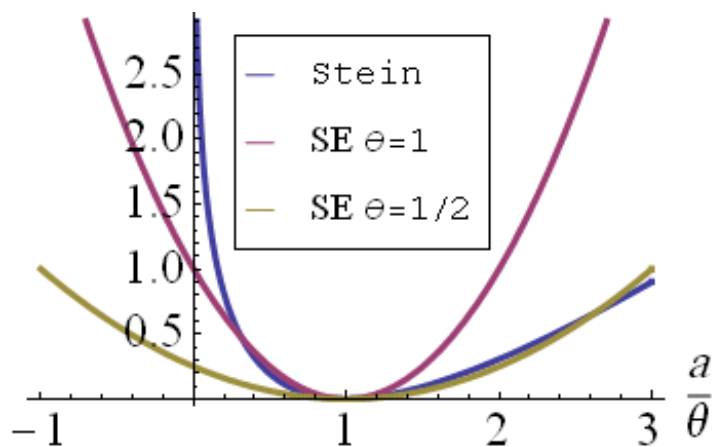
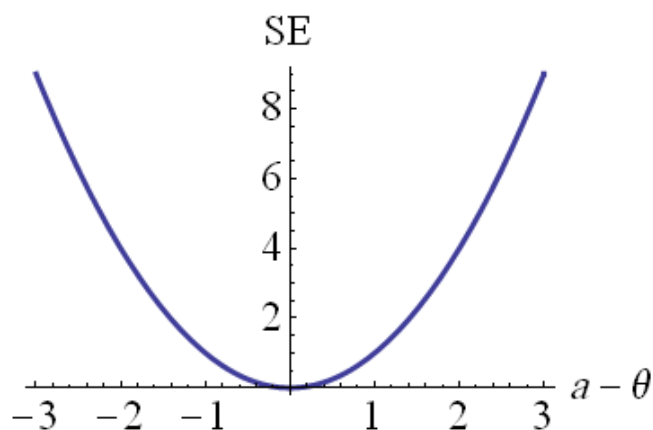
$$\text{Var}[T] \geq \frac{1}{nI(\theta)} = \frac{2\sigma^4}{n},$$

- lower bound of any unbiased estimator of σ^2 .
- not attainable by the unbiased version above

Assessing with different Loss Function:

$$L(\theta, a) = (a - \theta)^2 = \theta \left(\frac{a}{\theta} - 1 \right)^2 \quad (\text{SE loss})$$

$$L(\theta, a) = \frac{a}{\theta} - 1 - \log \left(\frac{a}{\theta} \right) \quad (\text{Stien's loss})$$



$$\tilde{S}^2 = b \sum_{i=1}^n \left(X_i - \bar{X} \right)^2$$

$$L(\theta, a) = \frac{a}{\theta} - 1 - \log\left(\frac{a}{\theta}\right)$$

$$\begin{aligned} R(\sigma^2, \tilde{S}^2) &= \mathbb{E} \left[b(n-1) \frac{S^2}{\sigma^2} - 1 - \log \frac{b(n-1) S^2}{\sigma^2} \right] \\ &= b \mathbb{E} [\chi_{n-1}^2] - 1 - \log b - \mathbb{E} \log \chi_{n-1}^2 \end{aligned}$$

$$\frac{\partial R}{\partial b} = \mathbb{E} [\chi_{n-1}^2] - \frac{1}{b} \quad (\stackrel{set}{=} 0)$$

$$b = \frac{1}{\mathbb{E} [\chi_{n-1}^2]} = \frac{1}{n-1}$$

$$\tilde{S}^2 = \frac{1}{n-1} \sum_{i=1}^n \left(X_i - \bar{X} \right)^2 = S^2.$$

“Better” in which sense?

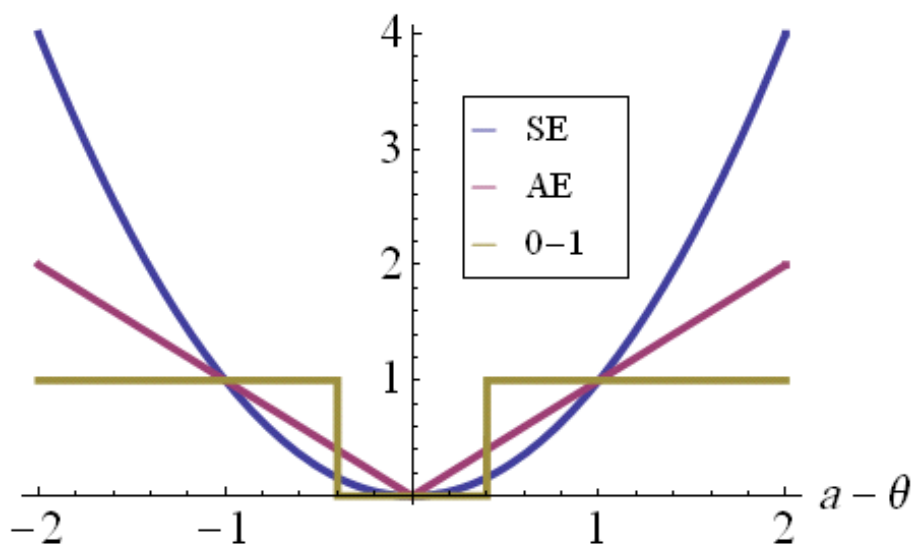
Obtaining Bayesian's Estimator by Loss Function Optimization!

$$\begin{aligned} R(\theta, T) &= \mathbb{E}_{\mathbf{X}} L(\theta, T(\mathbf{X})) \\ &= \int L(\theta, T(\mathbf{x})) f_{\mathbf{X}}(\mathbf{x}|\theta) d\mathbf{x} \end{aligned}$$

- no uniformly “best” estimator.
- $R(\theta, T_1)$ may cross with $R(\theta, T_2)$.

$$\begin{aligned} \mathbb{E}_{\Theta} R(\theta, T) &= \int_{\theta} R(\theta, T) f_{\Theta}(\theta) d\theta \\ &= \int_{\theta} \left[\int_{\mathbf{x}} L(\theta, T(\mathbf{x})) f_{\mathbf{X}}(\mathbf{x}|\theta) d\mathbf{x} \right] f_{\Theta}(\theta) d\theta \\ &= \int_{\mathbf{x}} \left[\int_{\theta} L(\theta, T(\mathbf{x})) f_{\mathbf{X}}(\mathbf{x}|\theta) f_{\Theta}(\theta) d\theta \right] d\mathbf{x} \\ &= \int_{\mathbf{x}} \left[\int_{\theta} L(\theta, T(\mathbf{x})) f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \right] f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \\ T &= \underset{T}{\operatorname{argmin}} \int_{\theta} L(\theta, T(\mathbf{x})) f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \end{aligned}$$

Solutions under different loss functions:



$$T_1 = \arg \min_T \int_{\theta} (T - \theta)^2 f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \quad (\text{SE loss})$$

$$= \int_{\theta} \theta f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \quad (\text{Posterior mean})$$

$$T_2 = \underset{T}{\operatorname{argmin}} \int_{\theta} |T - \theta| f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \quad (\text{AE loss})$$

$$\begin{aligned} R &= \int_{\theta} |T - \theta| f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \\ &= \int_{-\infty}^T (T - \theta) f(\theta) d\theta + \int_T^{\infty} - (T - \theta) f(\theta) d\theta \\ &= T \int_{-\infty}^T f(\theta) d\theta - \int_{-\infty}^T \theta f(\theta) d\theta - \\ &\quad T \int_T^{\infty} f(\theta) d\theta + \int_T^{\infty} \theta f(\theta) d\theta \end{aligned}$$

$$\begin{aligned} \frac{\partial R}{\partial T} &= \left(\int_{-\infty}^T f(\theta) d\theta + T f(T) \right) - T f(T) - \\ &\quad \left(\int_T^{\infty} f(\theta) d\theta - T f(T) \right) - T f(T) \\ &= \int_{-\infty}^T f(\theta) d\theta - \int_T^{\infty} f(\theta) d\theta \quad (\stackrel{set}{=} 0) \end{aligned}$$

$$0 = F_{\Theta|\mathbf{X}}^{-1}(T) - (1 - F_{\Theta|\mathbf{X}}^{-1}(T))$$

$$0.5 = F_{\Theta|\mathbf{X}}^{-1}(T)$$

$$T_2 = F_{\Theta|\mathbf{X}}^{-1}(0.5) \quad (\text{Posterior median})$$

$$T_3 = \arg \min_T \int_{\theta} I_{0 \leq |T-\theta|} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \quad (0-1 \text{ loss})$$

$$\begin{aligned} R &= \int_{\theta} I_{a \leq |T-\theta|} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \\ &= \int_{a \leq |T-\theta|} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \\ &= 1 - \int_{|T-\theta| < a} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \\ &= 1 - \int_{T-a}^{T+a} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \\ &= 1 - \Pr_{\Theta|\mathbf{X}}[|\theta - T| < a] \end{aligned}$$

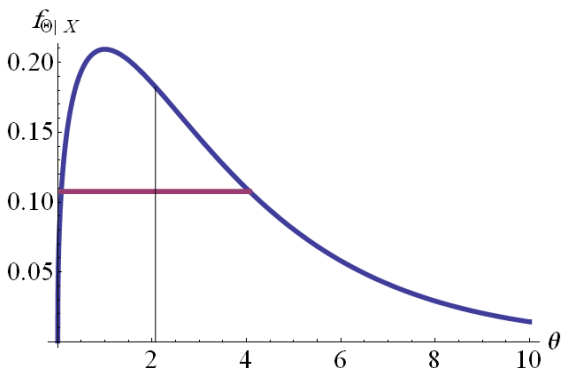
Notice that: we have to maximize the probability $\int_{T-a}^{T+a} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta$. The period $[T-a, T+a]$ has

- a length of $(T+a) - (T-a) = 2a$
- mid point of $\frac{1}{2} [(T+a) + (T-a)] = T$.
- T and mode do not necessarily coincide.,

which means that T_3 is mid-point of $2a$ modal interval.

$$\frac{\partial R}{\partial T} = f_{\Theta|\mathbf{X}}(T + a|\mathbf{X}) - f_{\Theta|\mathbf{X}}(T - a|\mathbf{X}), \quad (\stackrel{set}{=} 0)$$

$$f_{\Theta|\mathbf{X}}(T + a|\mathbf{X}) = f_{\Theta|\mathbf{X}}(T - a|\mathbf{X}).$$



For unimodal symmetric $f_{\Theta|\mathbf{X}}$:

$f_{\Theta|\mathbf{X}}(\theta - M) = f_{\Theta|\mathbf{X}}(\theta + M)$. Therefore,

$$T_3 = Mode. \quad (MAP)$$

For $a \rightarrow 0$

$$\begin{aligned} R &\approx 1 - f_{\Theta|\mathbf{X}}(T|\mathbf{x}) \cdot 2a, \\ T_3 &= \arg\max_T f_{\Theta|\mathbf{X}}(T|\mathbf{x}) = \textit{Mode} \quad (\text{MAP}) \end{aligned}$$

Of course T_3 could have been any point if we started minimizing the risk from beginning not by obtaining the limit:

$$\begin{aligned} R &= 1 - \int_{T-a}^{T+a} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \\ &= 1 - \int_T^T f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta \\ &= 1, \end{aligned}$$

unless Θ is discrete or categorical as in Pattern Recognition.

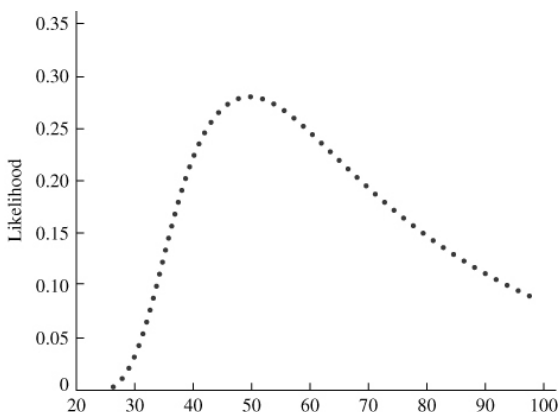
Estimation for Discrete Θ

MLE, Bayesian, Loss Functions have same treatment. However, maximization, expectation,..etc are taken over discrete space. Also, Cramér-Rao Lower Bound is derived for continuous case!

Example 40 (Capture Recapture Method) : *as in Example 15, page 19, first course. x captured animal in a population of θ animals. x was found to be 4 (we renamed variables):*

$$L(\theta) = P(x|\theta) = \frac{\binom{10}{4} \binom{\theta-10}{20-4}}{\binom{\theta}{20}}, \quad (\text{Likelihood})$$

$$\hat{\theta}_{MLE} = 50$$



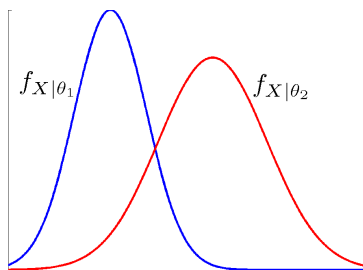
- maximization is obtained by $L_{\theta} / L_{\theta+1}$ not by $\frac{\partial L}{\partial \theta}$.
- Bayesian estimation is exactly the same through defining $f_{\Theta}(\theta)$.
- However, $f_{\Theta|\mathbf{X}}(\theta|\mathbf{x})$ will be discrete.

Estimation for Categorical Θ (basis for Pattern Recognition)

- $\Theta = \{\theta_1, \dots, \theta_K\}$, with K categories (classes).
- E.g., $\Theta = \{Male, Female\}$
- MoM is not applicable here (Θ is not numeric).

$$X|\theta_1 \sim N(1.5, .08),$$

$$X|\theta_2 \sim N(1.7, .1).$$



Suppose we got 1.77, 1.58, 1.77, 1.86, 1.75, 1.80, 1.77, 1.67, 1.73, 1.62. Are these readings obtained from Male or Female population?

Bayesian Estimation and MLE

8.5.3 Asymptotic Relative Efficiency (ARE)

Definition 41 *The (sequence of) estimator T_n is said to be asymptotically efficient for θ if*

$$\sqrt{n}(T_n - \theta) \xrightarrow{d} N(0, \sigma^2),$$
$$\sigma^2 = \frac{1}{I(\theta)},$$

which is Cramér-Rao Lower Bound.

It is clear that MLE is asymptotically efficient.