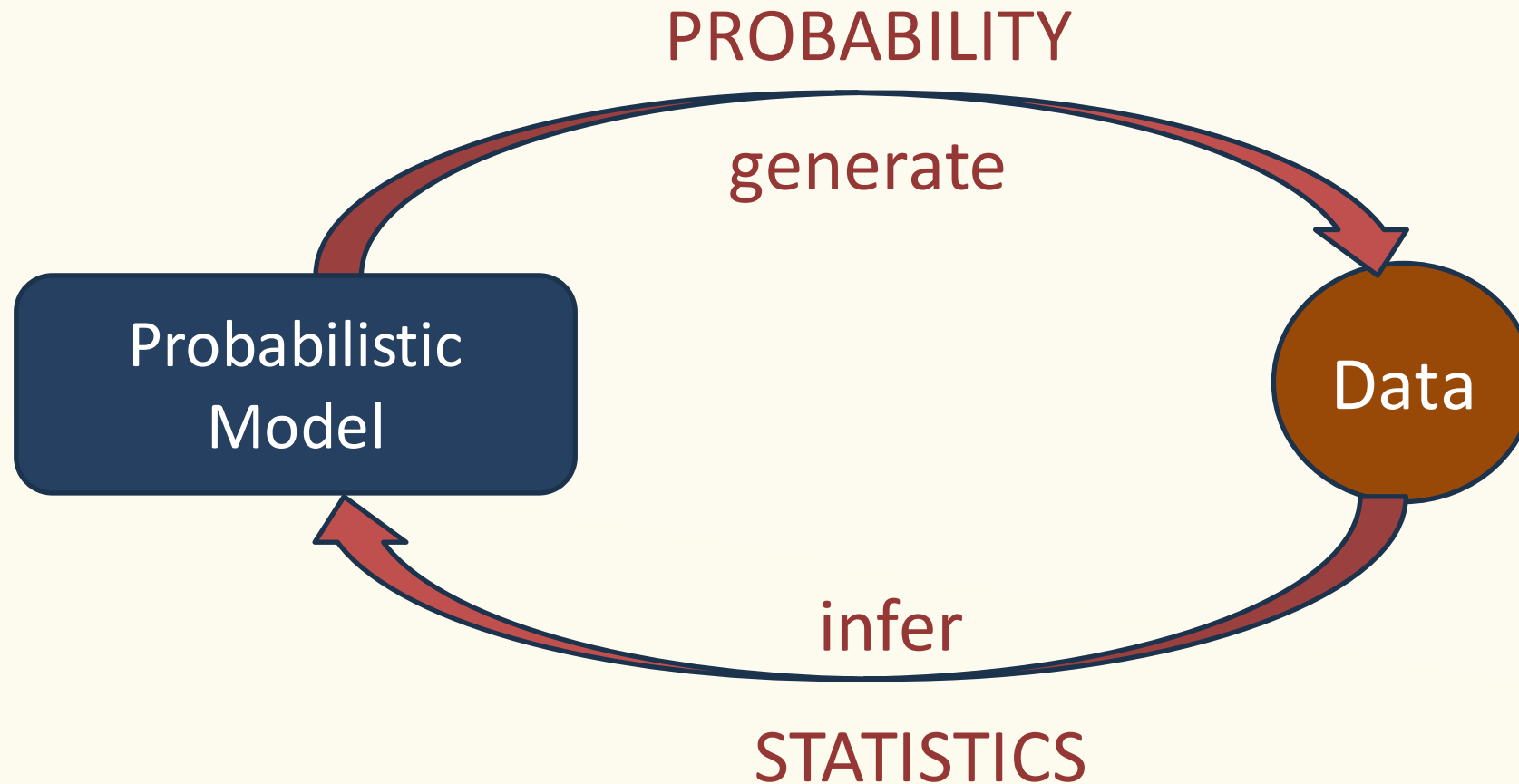


More Maximum Likelihood Estimation

CSE 312 Spring 26
Lecture 24

Several slides by Mor Harchol-Balter from
Book “Intro to Probability for Computing”

Probability vs. Statistics



Maximum Likelihood Estimation: Given some data drawn from a known distribution with unknown parameters, estimate the parameters of the distribution

Creating an maximum likelihood estimator

Goal: Estimate an unknown value θ , given sample data X from known distribution with unknown parameter(s) θ .

Step 1: Define

$$\hat{\theta}_{ML}(X = x) = \arg \max_{\theta} \mathbf{P}(X = x; \theta)$$

Our estimator is the value of θ that maximizes the likelihood function

This is the “likelihood function”

Example of MLE

Goal: Estimate θ = Number of pink jelly beans

Experiment: Randomly sample $n = 20$ beans w/ replacement

X = # pink jelly beans in sample



Example of MLE

Goal: Estimate θ = Number of pink jelly beans

Experiment: Randomly sample $n = 20$ beans w/ replacement

X = # pink jelly beans in sample is Binomial ($n, \theta / 1000$)



Suppose we observe $X = 3$ pinks in our sample with $n = 20$.

$$P(X = 3; \theta) = \binom{20}{3} \left(\frac{\theta}{1000} \right)^3 \cdot \left(1 - \frac{\theta}{1000} \right)^{17}$$

What value
of θ
maximizes
this?

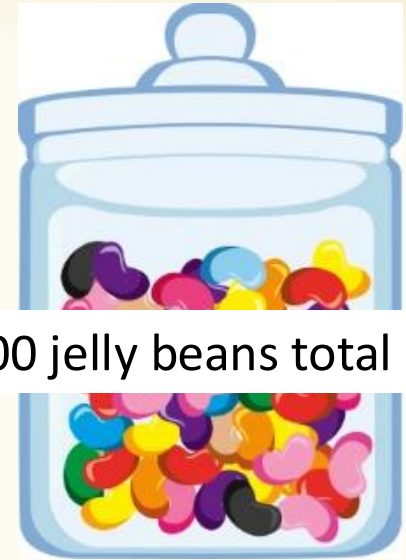
Example of MLE

Goal: Estimate θ = Number of pink jelly beans

Experiment: Randomly sample $n = 20$ beans w/ replacement

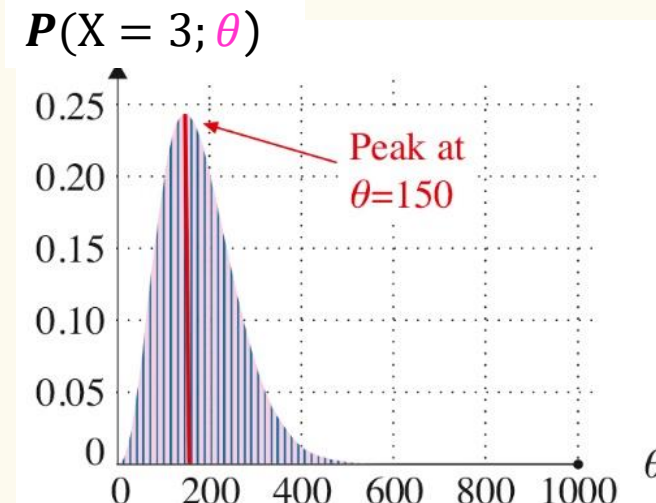
X = # pink jelly beans in sample

Q: Suppose we observe $X = 3$ pinks in our sample. What is $P(X = 3; \theta)$?



What value of θ maximizes this?

$\theta = 150$



$$\hat{\theta}_{ML}(X = 3)$$

$$= \arg \max_{\theta} P(X = 3; \theta)$$

$$= 150$$

Example of MLE

Q: What is the likelihood function $\mathbf{P}(X = x; \theta)$?

$$\mathbf{P}(X = x; \theta) = \binom{n}{x} \left(\frac{\theta}{1000}\right)^x \cdot \left(1 - \frac{\theta}{1000}\right)^{n-x}$$

Q: What is $\hat{\theta}_{ML}(X = x) = \arg \max_{\theta} \mathbf{P}(X = x; \theta)$?

$$\begin{aligned} 0 &= \frac{d}{d\theta} \mathbf{P}(X = x; \theta) = \frac{d}{d\theta} \binom{n}{x} \left(\frac{\theta}{1000}\right)^x \cdot \left(1 - \frac{\theta}{1000}\right)^{n-x} \\ &= \binom{n}{x} \left(\frac{\theta}{1000}\right)^x \cdot (n-x) \cdot \left(1 - \frac{\theta}{1000}\right)^{n-x-1} \cdot \frac{-1}{1000} + \binom{n}{x} \cdot x \left(\frac{\theta}{1000}\right)^{x-1} \cdot \frac{1}{1000} \cdot \left(1 - \frac{\theta}{1000}\right)^{n-x} \\ &= -\frac{n-x}{1000} \cdot \frac{\theta}{1000} + \frac{x}{1000} \cdot \left(1 - \frac{\theta}{1000}\right) \end{aligned}$$

Solving, we get:

$$\theta = \frac{1000x}{n}$$

2nd derivative is negative, so this is a max ✓



Maximizing Log likelihood Simpler

Goal: Estimate an unknown value θ , given sample data X

Define

$$\hat{\theta}(X = x) = \arg \max_{\theta} \underbrace{P(X = x; \theta)}_{\text{This is the likelihood function}} = \arg \max_{\theta} \underbrace{\ln P(X = x; \theta)}_{\text{This is the log likelihood function}}$$

equivalent
Why??

Example of MLE

$$\begin{aligned}\ln(ab) &= \ln(a) + \ln(b) \\ \ln(a/b) &= \ln(a) - \ln(b) \\ \ln(a^b) &= b \cdot \ln(a)\end{aligned}$$

The likelihood function $\mathbf{P}(X = \mathbf{x}; \boldsymbol{\theta})$?

$$\mathbf{P}(X = \mathbf{x}; \boldsymbol{\theta}) = \binom{n}{x} \left(\frac{\theta}{1000}\right)^x \cdot \left(1 - \frac{\theta}{1000}\right)^{n-x}$$

$$\ln \mathbf{P}(X = \mathbf{x}; \boldsymbol{\theta}) = \ln \binom{n}{x} + x \ln \left(\frac{\theta}{1000}\right) + (n - x) \ln \left(1 - \frac{\theta}{1000}\right)$$

$$0 = \frac{d}{d\theta} \ln \mathbf{P}(X = \mathbf{x}; \boldsymbol{\theta})$$



1000 jelly beans total

X = # pink jelly beans
in sample



Example of MLE

What is the likelihood function $\mathbf{P}\{X = x \mid \theta\}$?

$$\mathbf{P}(X = x; \theta) = \binom{n}{x} \left(\frac{\theta}{1000}\right)^x \cdot \left(1 - \frac{\theta}{1000}\right)^{n-x}$$



1000 jelly beans total


X = # pink jelly beans
in sample

What is $\hat{\theta}(X = x) = \arg \max_{\theta} \mathbf{P}(X = x; \theta)$ = $\arg \max_{\theta} \ln \mathbf{P}(X = x; \theta)$?

$$\hat{\theta}(X = x) = \arg \max_{\theta} \mathbf{P}(X = x; \theta) = \frac{1000x}{n}$$

This holds
 $\forall x$

Agenda

- MLE for Continuous Distributions 
- MLE for Normal Distribution
- Unbiased and Consistent Estimators

The Continuous Case

Given n (independent) samples $X_1 = x_1, \dots, X_n = x_n$ from (continuous) parametric model $f(x_i; \theta)$ which is now a family of densities

Definition. The **likelihood** of independent observations x_1, \dots, x_n is

$$\mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

Replace pmf with pdf!

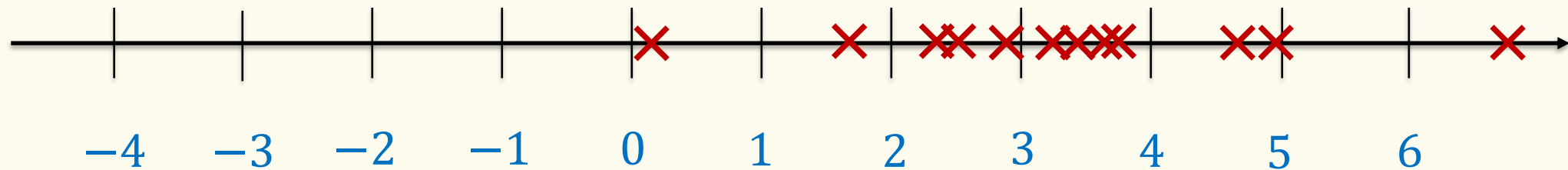
Why density?

- Density \neq probability, but:
 - For maximizing likelihood, **we really only care about relative likelihoods**, and density captures that
 - has desired property that likelihood increases with better fit to the model

Agenda

- MLE for for Continuous Distributions
- **MLE for Normal Distribution** 
- Unbiased and Consistent Estimators

n samples $X_1 = x_1, \dots, X_n = x_n \in \mathbb{R}$ from Gaussian $\mathcal{N}(\mu, 1)$. Most likely μ ?
[i.e., we are given the promise that the variance is 1]



Example – Gaussian Parameters

$$\begin{aligned}\ln(ab) &= \ln(a) + \ln(b) \\ \ln(a/b) &= \ln(a) - \ln(b) \\ \ln(a^b) &= b \cdot \ln(a)\end{aligned}$$

Normal outcomes $X_1 = x_1, \dots, X_n = x_n$, known variance $\sigma^2 = 1$

Goal: estimate θ , the unknown expectation

$$\begin{aligned}\mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta) &= \prod_{i=1}^n f(x_i; \theta) = \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i - \theta)^2}{2}} \right) \\ &= \left(\frac{1}{\sqrt{2\pi}} \right)^n \prod_{i=1}^n e^{-\frac{(x_i - \theta)^2}{2}}\end{aligned}$$

Goal: estimate θ = expectation

Example – Gaussian Parameters

Normal outcomes $X_1 = x_1, \dots, X_n = x_n$, known variance $\sigma^2 = 1$

$$\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta) = -n \frac{\ln 2\pi}{2} - \sum_{i=1}^n \frac{(x_i - \theta)^2}{2}$$

$$\frac{d}{d\theta} \ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta) =$$

Note: $\frac{d}{d\theta} \frac{(x_i - \theta)^2}{2} = \frac{1}{2} \cdot 2 \cdot (x_i - \theta) \cdot (-1) = \theta - x_i$

Goal: estimate θ = expectation

Example – Gaussian Parameters

Normal outcomes $X_1 = x_1, \dots, X_n = x_n$, known variance $\sigma^2 = 1$

$$\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta) = -n \frac{\ln 2\pi}{2} - \sum_{i=1}^n \frac{(x_i - \theta)^2}{2}$$

Note: $\frac{\partial}{\partial \theta} \frac{(x_i - \theta)^2}{2} = \frac{1}{2} \cdot 2 \cdot (x_i - \theta) \cdot (-1) = \theta - x_i$

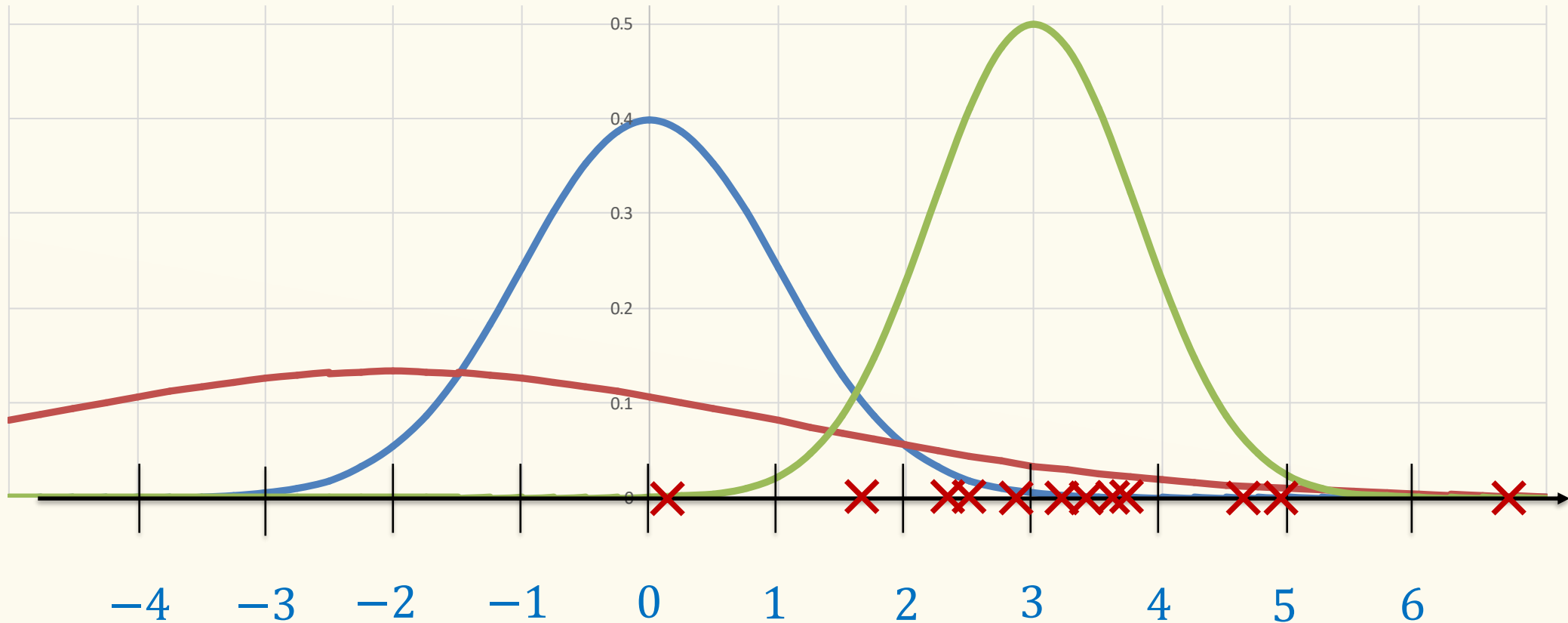
$$\frac{\partial}{\partial \theta} \ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta) = \sum_{i=1}^n (x_i - \theta) = \sum_{i=1}^n x_i - n\theta$$

So... solve $\sum_{i=1}^n x_i - n\hat{\theta} = 0$ for $\hat{\theta}$

$$\hat{\theta}(X_1 = x_1, \dots, X_n = x_n) = \frac{\sum_{i=1}^n x_i}{n}$$

In other words, MLE is the *sample mean* of the data.

Next: n samples $X_1 = x_1, \dots, X_n = x_n \in \mathbb{R}$ from Gaussian $\mathcal{N}(\mu, \sigma^2)$. Most likely μ and σ^2 ?



Two-parameter optimization

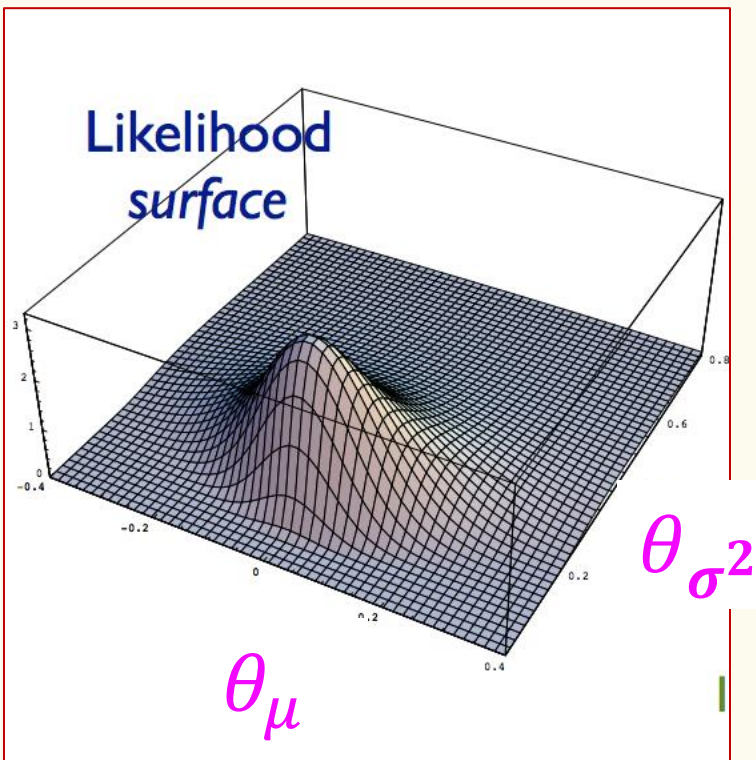
$$\begin{aligned}\ln(ab) &= \ln(a) + \ln(b) \\ \ln(a/b) &= \ln(a) - \ln(b) \\ \ln(a^b) &= b \cdot \ln(a)\end{aligned}$$

Normal outcomes $X_1 = x_1, \dots, X_n = x_n$

Goal: estimate θ_μ = expectation and θ_{σ^2} = variance

$$\mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) = \prod_{i=1}^n f(x_i; \theta_\mu, \theta_{\sigma^2})$$

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



Two-parameter estimation

$$\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) = -\frac{\ln(2\pi \theta_{\sigma^2})}{2} - \sum_{i=1}^n \frac{(x_i - \theta_\mu)^2}{2\theta_{\sigma^2}}$$

Find pair $\hat{\theta}_\mu, \hat{\theta}_{\sigma^2}$ that maximizes $\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2})$

Two-parameter estimation

$$\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) = -\frac{\ln(2\pi \theta_{\sigma^2})}{2} - \sum_{i=1}^n \frac{(x_i - \theta_\mu)^2}{2\theta_{\sigma^2}}$$

We need to find a solution $\hat{\theta}_\mu, \hat{\theta}_{\sigma^2}$ to

$$\frac{\partial}{\partial \theta_\mu} \ln \mathcal{L}(x_1, \dots, x_n; \theta_\mu, \theta_{\sigma^2}) = 0$$
$$\frac{\partial}{\partial \theta_{\sigma^2}} \ln \mathcal{L}(x_1, \dots, x_n; \theta_\mu, \theta_{\sigma^2}) = 0$$

And then check second order conditions.

MLE for Expectation

$$\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) = -n \frac{\ln(2\pi \theta_{\sigma^2})}{2} - \sum_{i=1}^n \frac{(x_i - \theta_\mu)^2}{2\theta_{\sigma^2}}$$

MLE for Expectation

$$\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) = -n \frac{\ln(2\pi \theta_{\sigma^2})}{2} - \sum_{i=1}^n \frac{(x_i - \theta_\mu)^2}{2\theta_{\sigma^2}}$$

$$\frac{\partial}{\partial \theta_\mu} \ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) = \frac{1}{\theta_{\sigma^2}} \sum_{i=1}^n (x_i - \theta_\mu) = 0$$

$$\frac{\partial}{\partial \theta_{\sigma^2}} \ln \mathcal{L}(x_1, \dots, x_n; \hat{\theta}_\mu, \theta_{\sigma^2}) = -\frac{n}{2\theta_{\sigma^2}} + \frac{1}{2\theta_{\sigma^2}^2} \sum_{i=1}^n (x_i - \theta_\mu)^2 = 0$$

MLE for Expectation

$$\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) = -n \frac{\ln(2\pi \theta_{\sigma^2})}{2} - \sum_{i=1}^n \frac{(x_i - \theta_\mu)^2}{2\theta_{\sigma^2}}$$

$$\frac{\partial}{\partial \theta_\mu} \ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) = \frac{1}{\theta_{\sigma^2}} \sum_{i=1}^n (x_i - \theta_\mu) = 0$$

$$\hat{\theta}_\mu(X_1 = x_1, \dots, X_n = x_n) = \frac{\sum_{i=1}^n x_i}{n}$$

In other words, MLE of expectation is (again) the *sample mean* of the data, regardless of θ_2

What about the variance?

MLE for Variance

$$\begin{aligned}\ln \mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta_\mu, \theta_{\sigma^2}) &= -n \frac{\ln(2\pi \theta_{\sigma^2})}{2} - \sum_{i=1}^n \frac{(x_i - \theta_\mu)^2}{2\theta_{\sigma^2}} \\ &= -n \frac{\ln 2\pi}{2} - n \frac{\ln \theta_{\sigma^2}}{2} - \frac{1}{2\theta_{\sigma^2}} \sum_{i=1}^n (x_i - \hat{\theta}_\mu)^2\end{aligned}$$

$$\frac{\partial}{\partial \theta_{\sigma^2}} \ln \mathcal{L}(x_1, \dots, x_n; \hat{\theta}_\mu, \theta_{\sigma^2}) = -\frac{n}{2\theta_{\sigma^2}} + \frac{1}{2\theta_{\sigma^2}^2} \sum_{i=1}^n (x_i - \hat{\theta}_\mu)^2 \quad \text{Set } = 0$$

$$\hat{\theta}_{\sigma^2} (X_1 = x_1, \dots, X_n = x_n) = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\theta}_\mu)^2$$

In other words, MLE of variance is the *population variance* of the data.

Likelihood – Continuous Case

Definition. The **likelihood** of independent observations x_1, \dots, x_n is

$$\mathcal{L}(X_1 = x_1, \dots, X_n = x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

Normal outcomes $(X_1 = x_1, \dots, X_n = x_n)$

$$\hat{\theta}_\mu (X_1 = x_1, \dots, X_n = x_n) = \frac{\sum_{i=1}^n x_i}{n}$$

MLE estimator for
expectation

$$\hat{\theta}_{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\theta}_\mu)^2$$

MLE estimator for
variance

General Recipe (single parameter)

1. **Input** Given n i.i.d. samples $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ from parametric model with parameter θ .
2. **Likelihood** Define your likelihood $\mathcal{L}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n; \theta)$.
 - For discrete $\mathcal{L}(x_1, \dots, x_n; \theta) = \prod_{i=1}^n P(x_i; \theta)$
 - For continuous $\mathcal{L}(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$
3. **Log** Compute $\ln \mathcal{L}(x_1, \dots, x_n; \theta)$
4. **Differentiate** Compute $\frac{d}{d\theta} \ln \mathcal{L}(x_1, \dots, x_n; \theta)$
5. **Solve for $\hat{\theta}$** by setting derivative to 0 and solving for max.

Do a second derivative test to verify it is a maximum, but we won't ask you to do that in CSE 312.

General Recipe (multiple parameters)

1. **Input** Given n i.i.d. samples $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ from parametric model with parameters $\vec{\theta} = (\theta_1, \dots, \theta_k)$.

2. **Likelihood** Define your likelihood $\mathcal{L}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n; \vec{\theta})$.

– For discrete $\mathcal{L}(x_1, \dots, x_n; \vec{\theta}) = \prod_{i=1}^n P(x_i; \vec{\theta})$

– For continuous $\mathcal{L}(x_1, \dots, x_n; \vec{\theta}) = \prod_{i=1}^n f(x_i; \vec{\theta})$

3. **Log** Compute $\ln \mathcal{L}(x_1, \dots, x_n; \vec{\theta})$


4. **Differentiate** Compute $\frac{\partial}{\partial \theta_i} \ln \mathcal{L}(x_1, \dots, x_n; \vec{\theta})$

5. **Solve for $\hat{\theta}$** by setting all the derivatives to 0 and solving.

Check second order conditions, but we won't ask you to do that in CSE 312.



Agenda

- MLE for for Continuous Distributions
- Normal Distribution
- Unbiased and Consistent Estimators 

Definition of Estimator

e.g. $\theta =$
parameter of a Bernoulli distri

θ : quantity we're trying to estimate

think of these as
i. i. d. instances of X
 $\sim \text{Ber}(\theta)$

This is a
constant

X_1, X_2, \dots, X_n : i.i.d. data

$\hat{\theta}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$: estimation of θ based on
specific instantiation of the data

This is a r.v.
because it's a
function of r.v.s

$\hat{\theta}(X_1, X_2, \dots, X_n)$: estimator of the unknown θ

Sometimes just
write $\hat{\theta}$

MLE for pink jelly beans

Q: What is the likelihood function $\mathbf{P}\{X = x \mid \theta\}$?

$$\mathbf{P}(X = x; \theta) = \binom{n}{x} \left(\frac{\theta}{1000}\right)^x \cdot \left(1 - \frac{\theta}{1000}\right)^{n-x}$$

Q: What is $\hat{\theta}_{ML}(X = x) = \arg \max_{\theta} \mathbf{P}(X = x; \theta)$?

$$\hat{\theta}_{ML}(X = x) = \arg \max_{\theta} \mathbf{P}(X = x; \theta) = \frac{1000x}{n}$$

$$\rightarrow \hat{\theta}_{ML}(X) = \frac{1000X}{n}$$



1000 jelly beans total



X = # pink jelly beans
in sample



This holds
 $\forall x$

MLE estimation for normal distribution

$$\hat{\theta}_{\mu}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = \frac{x_1 + x_2 + \dots + x_n}{n}$$

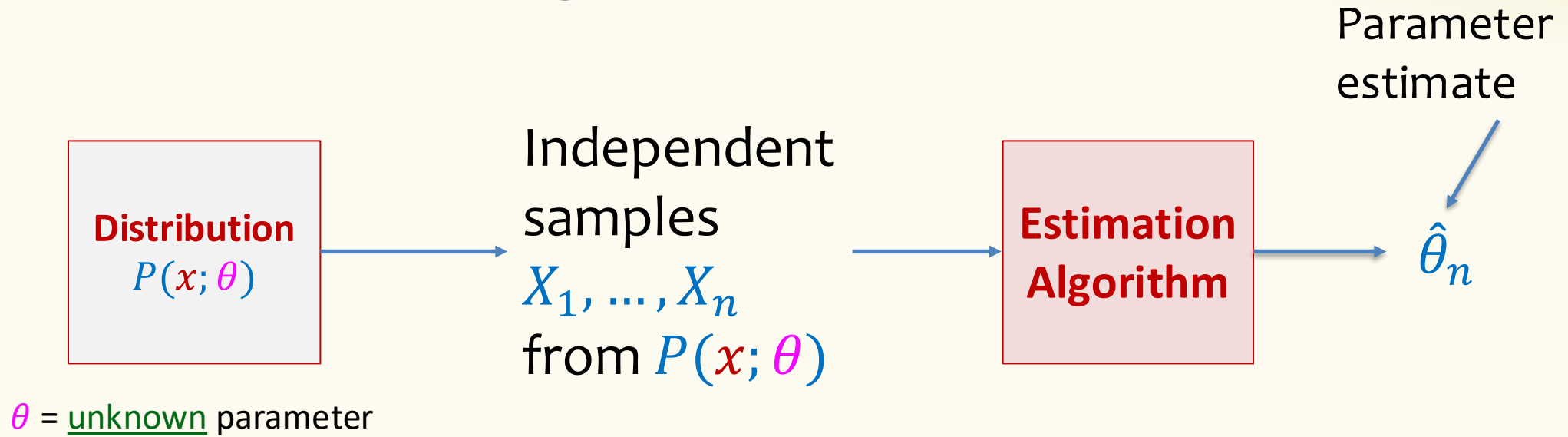
These hold for all
 x_1, \dots, x_n

$$\hat{\theta}_{\sigma^2}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = \frac{1}{n} \sum_{i=1}^n \left(x_i - \frac{x_1 + x_2 + \dots + x_n}{n} \right)^2$$

→
$$\hat{\theta}_{\mu}(X_1, X_2, \dots, X_n) = \frac{X_1 + X_2 + \dots + X_n}{n}$$

→
$$\hat{\theta}_{\sigma^2}(X_1, X_2, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \left(X_i - \frac{X_1 + X_2 + \dots + X_n}{n} \right)^2$$

When is an estimator good?



Definition. An estimator of parameter θ is an **unbiased estimator** if

$$\mathbb{E}[\hat{\theta}_n(X_1, \dots, X_n)] = \theta.$$

Note: This expectation is over the samples X_1, \dots, X_n

MLE for pink jelly beans

Q: What is the likelihood function $\mathbf{P}\{X = x \mid \theta\}$?

$$\mathbf{P}(X = x; \theta) = \binom{n}{x} \left(\frac{\theta}{1000}\right)^x \cdot \left(1 - \frac{\theta}{1000}\right)^{n-x}$$

Q: What is $\hat{\theta}_{ML}(X = x) = \arg \max_{\theta} \mathbf{P}(X = x; \theta)$?

$$\hat{\theta}_{ML}(X = x) = \arg \max_{\theta} \mathbf{P}(X = x; \theta) = \frac{1000x}{n}$$

→ $\hat{\theta}_{ML}(X) = \frac{1000X}{n}$

This holds
 $\forall x$

Fact. $\hat{\theta}_{ML}(X)$ is unbiased

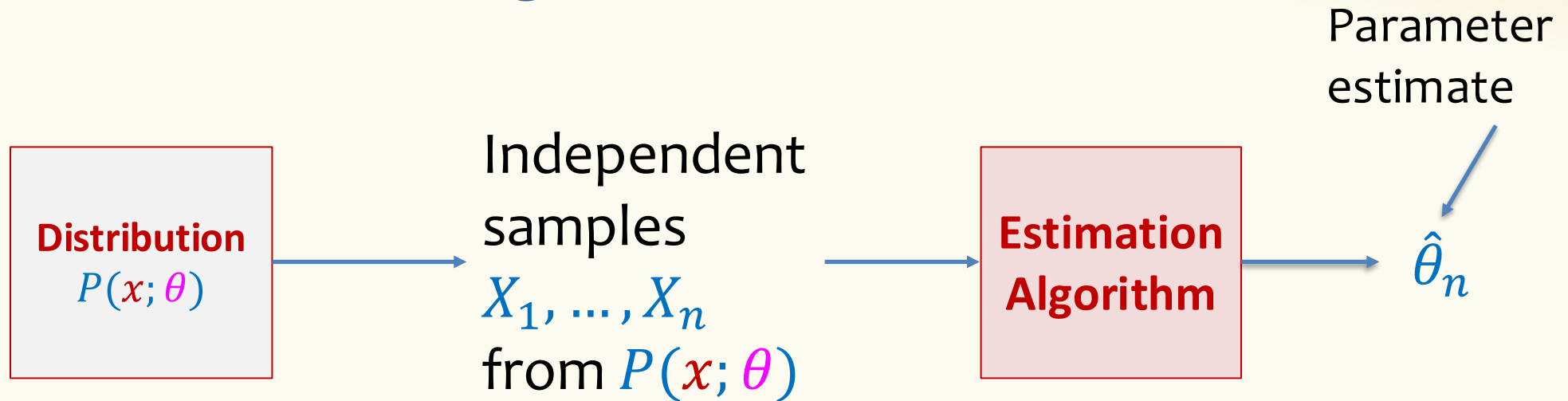


1000 jelly beans total

X = # pink jelly beans
in sample

Three samples from $U(0, \theta)$

When is an estimator good?



θ = unknown parameter

Definition. An estimator is **unbiased** if $\mathbb{E}[\hat{\theta}_n] = \theta$ for all $n \geq 1$.

Definition. An estimator is **consistent** if $\lim_{n \rightarrow \infty} \mathbb{E}[\hat{\theta}_n] = \theta$.

Theorem. MLE estimators are consistent.

(But not necessarily unbiased)

Example – Consistency

Normal outcomes X_1, \dots, X_n i.i.d. according to $\mathcal{N}(\mu, \sigma^2)$ Assume: $\sigma^2 > 0$

$$\hat{\Theta}_{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\Theta}_{\mu})^2$$

Population variance – Biased!

$\hat{\Theta}_{\sigma^2}$ is “consistent”

Example – Consistency

Normal outcomes X_1, \dots, X_n i.i.d. according to $\mathcal{N}(\mu, \sigma^2)$ **Assume:** $\sigma^2 > 0$

$$\hat{\Theta}_{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\Theta}_{\mu})^2$$

Population variance – Biased!

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\Theta}_{\mu})^2$$

Sample variance – Unbiased!

$\hat{\Theta}_{\sigma^2}$ converges to same value as S_n^2 , i.e., σ^2 , as $n \rightarrow \infty$.

$\hat{\Theta}_{\sigma^2}$ is “consistent”

Why does it matter?

- When statisticians are estimating a variance from a sample, they usually divide by $n-1$ instead of n .
- They and we not only want good estimators (unbiased, consistent)
 - They/we also want **confidence bounds**
 - Upper bounds on the probability that these estimators are far the truth about the underlying distributions
 - Confidence bounds are just like what we wanted for our polling problems, but CLT is usually not the best thing to use to get them (unless the variance is known)