

7.1, 7.2 MAXIMUM LIKELIHOOD ESTIMATION

ANNA KARLIN
MOST SLIDES BY ALEX TSUN

AGENDA

- PROBABILITY VS STATISTICS
- LIKELIHOOD
- MAXIMUM LIKELIHOOD ESTIMATION (MLE)
- MLE EXAMPLE (POISSON)
- MLE EXAMPLE (NORMAL)

PROBABILITY VS STATISTICS

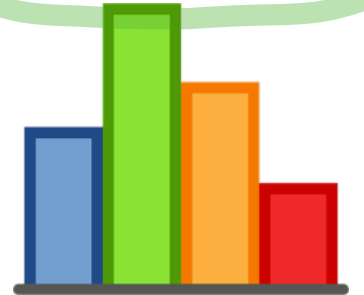
$Ber(p = 0.5)$



Probability
given model, predict data



$P(THHTHH)$



PROBABILITY VS STATISTICS



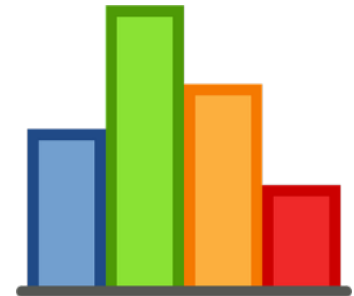
$Ber(p = 0.5)$



Probability
given model, predict data



$P(\text{THHTHH})$



$Ber(p = ???)$



Statistics
given data, predict model



THHTHH

using parametric model of data.

Bin(n, p), Exp(λ)
 Θ parameter unknown.

RANDOM PICTURE



LIKELIHOOD (INTUITION)



I give you and your classmates each 5 minutes with a coin with unknown probability of heads p . Whoever has the closest estimate will get an A+ in the class. What do you do in your precious 5 minutes, and what do you give as your estimate?

T T H T H T T H

What is your estimate?

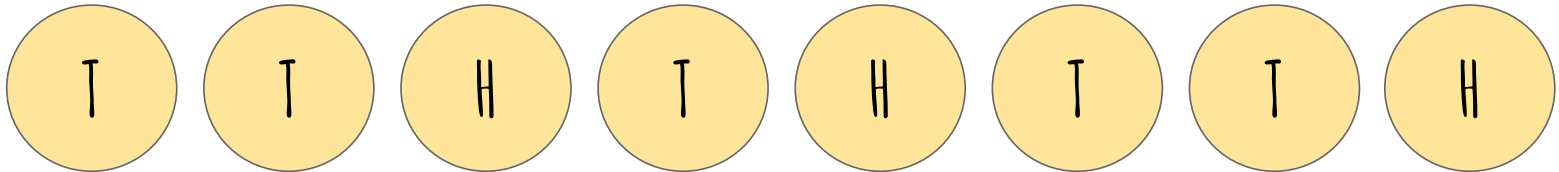
- a) $\frac{1}{2}$
- b) $\frac{1}{4}$
- c) $\frac{1}{2}$
- d) $\frac{1}{2}$

LIKELIHOOD (INTUITION)



I give you and your classmates each 5 minutes with a coin with unknown probability of heads p . Whoever has the closest estimate will get an A+ in the class. What do you do in your precious 5 minutes, and what do you give as your estimate?

Flip it as many times as possible, and return $\# H / (\# H + \# T)$!

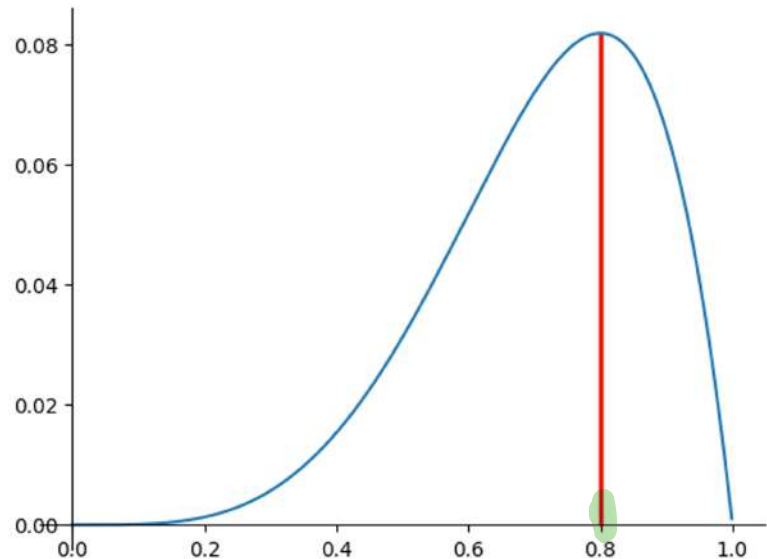


LIKELIHOOD (INTUITION)

Let's say you saw 4 heads
and 1 tail. You tell me $\hat{p} = \frac{4}{5} = 0.8$

How can you argue, *objectively*,
that this is the "best" estimate?

Is there some objective
function it maximizes?



LIKELIHOOD (INTUITION)



You assume a model (Bernoulli in our case) with unknown parameter θ , and receive iid samples $x = (x_1, \dots, x_n) \sim \text{Ber}(\theta)$. The likelihood of the data given a parameter θ is

$$L(x|\theta) = P(\text{seeing data} | \theta)$$

$$x_1 = H \quad x_2 = H \\ \dots$$

$$\left[\begin{array}{l} H H H H T \\ \theta \cdot \theta \cdot \theta \cdot \theta \cdot (1-\theta) \\ = \theta^4 (1-\theta) \end{array} \right]$$

LIKELIHOOD (INTUITION)



You assume a model (Bernoulli in our case) with unknown parameter θ , and receive iid samples $x = (x_1, \dots, x_n) \sim \text{Ber}(\theta)$. The **likelihood** of the data given a parameter θ is

$$\begin{aligned} L(x|\theta) &= P(\text{seeing data} \mid \theta) \\ &= P(x_1, \dots, x_n \mid \theta) \end{aligned}$$

LIKELIHOOD (INTUITION)



You assume a model (Bernoulli in our case) with unknown parameter θ , and receive iid samples $x = (x_1, \dots, x_n) \sim \text{Ber}(\theta)$. The **likelihood** of the data given a parameter θ is

$$\begin{aligned} L(x|\theta) &= P(\text{seeing data} \mid \theta) \\ &= P(x_1, \dots, x_n \mid \theta) \\ &= \prod_{i=1}^n p_X(x_i; \theta) \end{aligned}$$

$$p_X(H; \theta) = \theta$$

$$p_X(T; \theta) = 1 - \theta$$

MAXIMUM LIKELIHOOD ESTIMATION (BERNOULLI)

$\text{Ber}(\theta)$



$$L(\text{HHHHT} | \theta) = \theta^4(1 - \theta) = \theta^4 - \theta^5$$

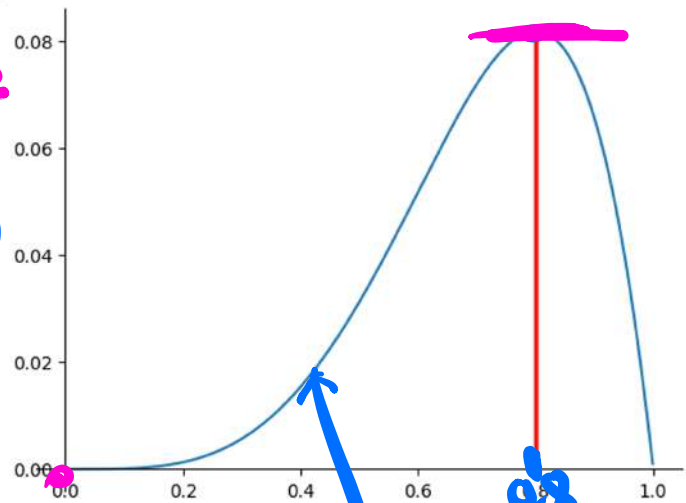
find param θ
that maximizes $L(\text{HHHHT} | \theta)$

$$\frac{d}{d\theta} [\theta^4 - \theta^5] = 4\theta^3 - 5\theta^4$$

$$\theta^3(4 - 5\theta) = 0$$

$$\theta = 0$$

$$\theta = \frac{4}{5}$$



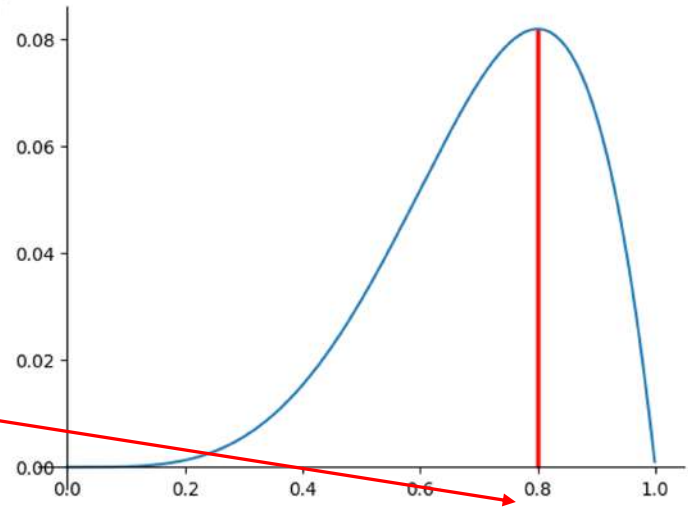
MAXIMUM LIKELIHOOD ESTIMATION (BERNOULLI)



$$L(\text{HHHHT} \mid \theta) = \theta^4(1 - \theta) = \theta^4 - \theta^5$$

$$\frac{\partial}{\partial \theta} L(x \mid \theta) = 4\theta^3 - 5\theta^4 = \theta^3(4 - 5\theta)$$

$$\hat{\theta}^3(4 - 5\hat{\theta}) = 0 \rightarrow \hat{\theta} = \frac{4}{5} \text{ or } 0$$



Likelihood

vs Probability. $\text{Ber}(\theta)$

$L(x; \theta)$
as fn of θ
(fixed x)

$P(x; \theta)$ = prob event x
given model $\text{Ber}(\theta)$
as fn of x (for fixed θ)

$\sum_{\theta} L(x; \theta)$ anything.

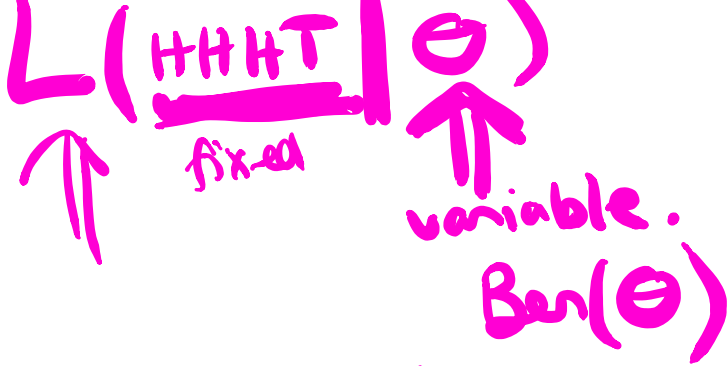
$$\sum_x P(x; \theta) = 1$$

$$L(\text{HHHT} | 0.6) > L(\text{HHHT} | 0.5)$$

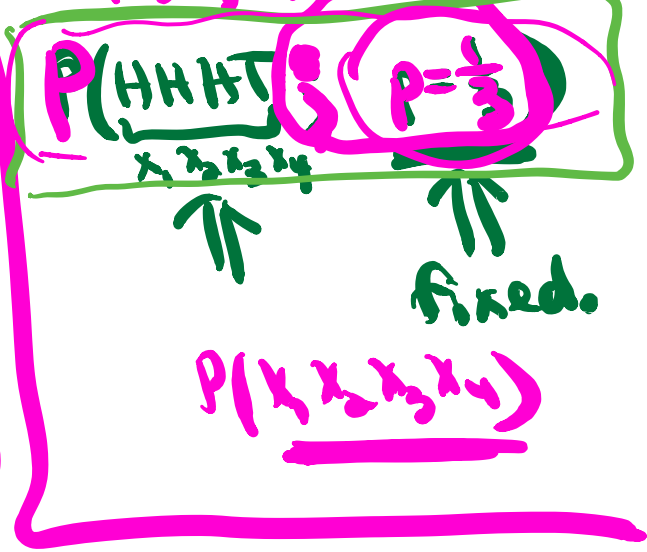
Finds parameters θ
that maximize likelihood



x_1, x_2, x_3, x_4 iid. $\text{Ber}(\frac{1}{3})$



Prob (seeing data for param θ)
 ↑ var.



$L(\underline{x_1 \dots x_n} | \theta)$ ← function of θ
 ↑

[assuming the param of distn is θ
 [not cond prob]



Assuming data i.i.d Poisson (θ)

x_1, \dots, x_n Samples

MAXIMUM LIKELIHOOD ESTIMATION (POISSON)

Let's say x_1, x_2, \dots, x_n are iid samples from $Poi(\theta)$. (might look like $x_1 = 3, x_2 = 5, x_3 = 4,$ etc.) What is the MLE of θ ?

↑
unknown.



MAXIMUM LIKELIHOOD ESTIMATION (POISSON)



Let's say x_1, x_2, \dots, x_n are iid samples from $Poi(\theta)$. (might look like $x_1 = 3, x_2 = 5, x_3 = 4$, etc.) What is the MLE of θ ?

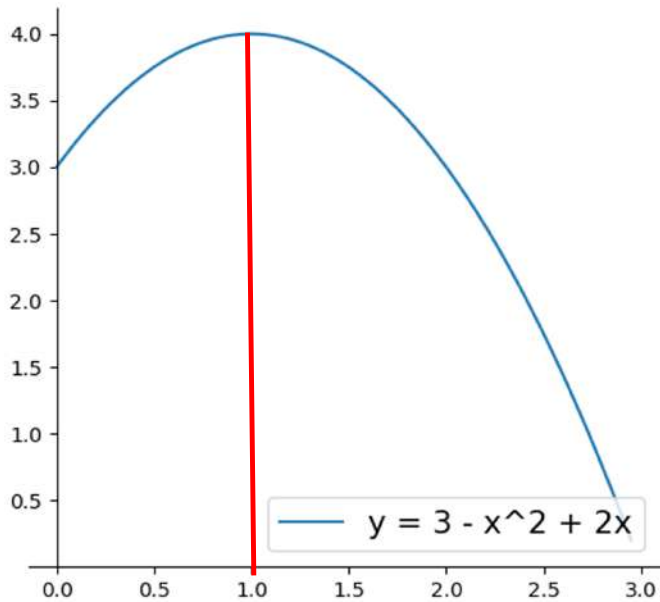
$$L(x | \theta) = \prod_{i=1}^n p_X(x_i; \theta) = \prod_{i=1}^n e^{-\theta} \frac{\theta^{x_i}}{x_i!}$$

Prob of seeing this data
as fn of parameter θ

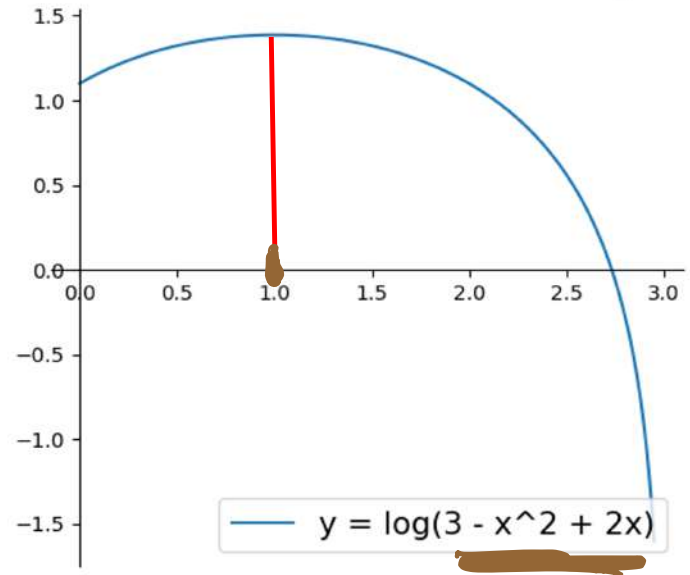
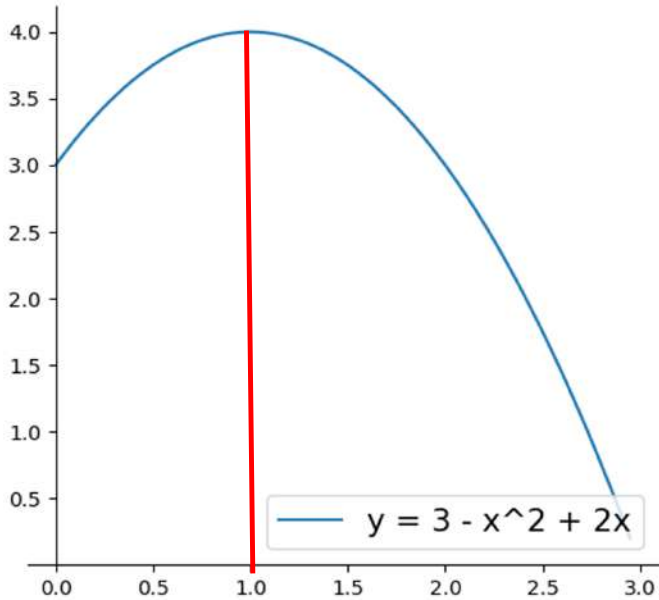
Find θ that maximizes $L(x | \theta)$

$$e^{-\theta} \frac{\theta^3}{3!} e^{-\theta} \frac{\theta^5}{5!} e^{-\theta} \frac{\theta^4}{4!}$$

OPTIMIZING FUNCTION VS LOG(FUNCTION)

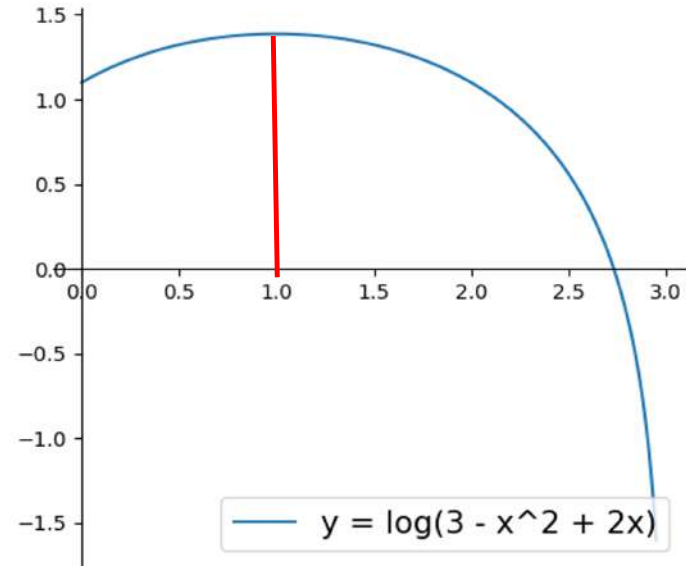
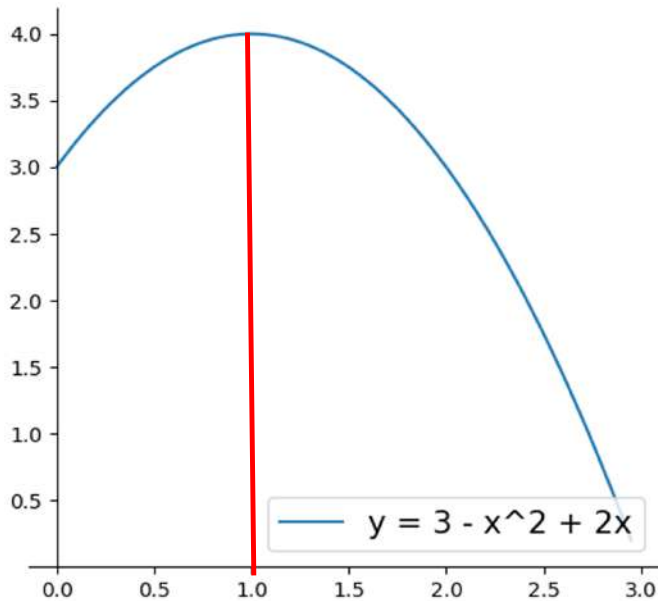


OPTIMIZING FUNCTION VS LOG(FUNCTION)



$f(x_1) > f(x_2)$ $\ln f(x_1) > \ln f(x_2)$
x that maximizes
 $f(x) \equiv x$ that maximizes $\ln f(x)$

OPTIMIZING FUNCTION VS LOG(FUNCTION)



Since $g(x) = \log x$ is **strictly increasing**, it preserves order, and in particular, the argmax.

MAXIMUM LIKELIHOOD ESTIMATION (POISSON)



Let's say x_1, x_2, \dots, x_n are iid samples from $Poi(\theta)$. (might look like $x_1 = 3, x_2 = 5, x_3 = 4$, etc.) What is the MLE of θ ?

$\log(ab) = \log(a) + \log(b)$

$\log(a/b) = \log(a) - \log(b)$

$\log(a^b) = b \log a$

$$L(x | \theta) = \prod_{i=1}^n p_X(x_i; \theta) = \prod_{i=1}^n e^{-\theta} \frac{\theta^{x_i}}{x_i!}$$

log of a:
= sum log a:

log likelihood.
 $LL(x|\theta) = \log L(x|\theta) = \sum_{i=1}^n \log \left(e^{-\theta} \frac{\theta^{x_i}}{x_i!} \right)$

$\underbrace{\ln e^{-\theta}}_{= -\theta} + x_i \ln \theta - \ln(x_i!)$

$-\theta \ln e$

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

MAXIMUM LIKELIHOOD ESTIMATION (POISSON)



Let's say x_1, x_2, \dots, x_n are iid samples from $Poi(\theta)$. (might look like $x_1 = 3, x_2 = 5, x_3 = 4$, etc.) What is the MLE of θ ?

$\log(ab) = \log(a) + \log(b)$ $L(x|\theta) = \prod_{i=1}^n p_X(x_i; \theta) = \prod_{i=1}^n e^{-\theta} \frac{\theta^{x_i}}{x_i!}$

$\log(a/b) = \log(a) - \log(b)$

$\log(a^b) = b \log a$

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

$$\frac{d}{d\theta} LL(x|\theta) = -n + \sum_{i=1}^n x_i \frac{1}{\theta}$$

$$= -n + \sum_{i=1}^n \frac{x_i}{\theta}$$

$$0 = -n + \sum_{i=1}^n \frac{x_i}{\theta}$$

$$\sum_{i=1}^n \frac{x_i}{\theta} = n$$

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i$$

MAXIMUM LIKELIHOOD ESTIMATION (POISSON)



Let's say x_1, x_2, \dots, x_n are iid samples from $Poi(\theta)$. (might look like $x_1 = 3, x_2 = 5, x_3 = 4$, etc.) What is the MLE of θ ?

$$L(\mathbf{x} | \theta) = \prod_{i=1}^n p_X(x_i; \theta) = \prod_{i=1}^n e^{-\theta} \frac{\theta^{x_i}}{x_i!}$$

$$\ln L(\mathbf{x} | \theta) = \sum_{i=1}^n [-\theta + x_i \ln \theta - \ln(x_i!)]$$

$$\frac{\partial}{\partial \theta} \ln L(\mathbf{x} | \theta) = \sum_{i=1}^n \left[-1 + \frac{x_i}{\theta}\right]$$

$$\sum_{i=1}^n \left[-1 + \frac{x_i}{\hat{\theta}}\right] = 0 \rightarrow -n + \frac{1}{\hat{\theta}} \sum_{i=1}^n x_i = 0 \rightarrow \hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i$$



MAXIMUM LIKELIHOOD ESTIMATION (POISSON)

Let's say x_1, x_2, \dots, x_n are iid samples from $Poi(\theta)$. (might look like $x_1 = 3, x_2 = 5, x_3 = 4$, etc.) What is the MLE of θ ?

$$L(\mathbf{x} | \theta) = \prod_{i=1}^n p_X(x_i; \theta) = \prod_{i=1}^n e^{-\theta} \frac{\theta^{x_i}}{x_i!}$$

$$\ln L(\mathbf{x} | \theta) = \sum_{i=1}^n [-\theta + x_i \ln \theta - \ln(x_i!)]$$

$$\frac{\partial}{\partial \theta} \ln L(\mathbf{x} | \theta) = \sum_{i=1}^n \left[-1 + \frac{x_i}{\theta}\right]$$

$$\sum_{i=1}^n \left[-1 + \frac{x_i}{\hat{\theta}}\right] = 0 \rightarrow -n + \frac{1}{\hat{\theta}} \sum_{i=1}^n x_i = 0 \rightarrow \hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial^2}{\partial \theta^2} \ln L(\mathbf{x} | \theta) = \sum_{i=1}^n \left[-\frac{x_i}{\theta^2}\right] < 0 \rightarrow \text{concave down everywhere}$$

General recipe

Given indep samples x_1, \dots, x_n
from parametric model

→ Bern(θ)

→ Poisson(θ)

Exp(θ)

Find $\hat{\theta}$ that
maximizes

$$L(x_1, \dots, x_n | \theta)$$

discrete

$$= \prod_{i=1}^n P_X(x_i; \theta) \leftarrow$$

cont.

$$\triangleq \prod_{i=1}^n f_X(x_i; \theta)$$

find $\hat{\theta}$ that maximizes
log likelihood

$$LL(x_1, \dots, x_n | \theta)$$

discrete

$$= \sum_{i=1}^n \ln [P_X(x_i; \theta)] \leftarrow$$

cont

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

Find $\hat{\theta}$ to max $LL(\vec{x} | \theta)$
 x_1, \dots, x_n

distn has 1 parameter

compute $\frac{dLL(\theta)}{d\theta}$

set $\frac{dLL(\theta)}{d\theta} = 0$

solve for $\hat{\theta}$

[verify soln is max (2nd deriv < 0)]

don't
need
to do
this

multiple params $\theta_1, \dots, \theta_k$

$$\frac{\partial LL}{\partial \theta_1} = 0$$

$$\frac{\partial LL}{\partial \theta_2} = 0$$

...

$$\frac{\partial LL}{\partial \theta_k} = 0$$

find
 $\hat{\theta}_1, \dots, \hat{\theta}_k$
that are
soln
to this
system

check max

check Hessian
-ve definite

LIKELIHOOD

Realization/Sample: A realization/sample x of a random variable X is the value that is actually observed.

LIKELIHOOD

Realization/Sample: A realization/sample x of a random variable X is the value that is actually observed.

Likelihood: Let $x = (x_1, \dots, x_n)$ be iid realizations from probability mass function $p_X(t; \theta)$ (if X discrete), or from density $f_X(t; \theta)$ (if X continuous), where θ is a parameter (or vector of parameters). We define the likelihood of x as the probability of seeing the data.

LIKELIHOOD

Realization/Sample: A realization/sample x of a random variable X is the value that is actually observed.

Likelihood: Let $x = (x_1, \dots, x_n)$ be iid realizations from probability mass function $p_X(t; \theta)$ (if X discrete), or from density $f_X(t; \theta)$ (if X continuous), where θ is a parameter (or vector of parameters). We define the likelihood of x as the probability of seeing the data.

If X is discrete,

$$L(x | \theta) = \prod_{i=1}^n p_X(x_i; \theta)$$

X cont.

$$L(\vec{x} | \theta) = \prod_{i=1}^n f_X(x_i; \theta)$$

MAXIMUM LIKELIHOOD ESTIMATION (MLE)

Maximum Likelihood Estimation (MLE): Let $\mathbf{x} = (x_1, \dots, x_n)$ be iid realizations from probability mass function $p_X(t; \theta)$ (if X discrete), or from density $f_X(t; \theta)$ (if X continuous), where θ is a parameter (or vector of parameters). We define the maximum likelihood estimator $\hat{\theta}_{MLE}$ of θ to be the parameter which maximizes the likelihood (or equivalently, the log-likelihood).

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(\mathbf{x} | \theta) = \arg \max_{\theta} \ln L(\mathbf{x} | \theta)$$

RANDOM PICTURE

