

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



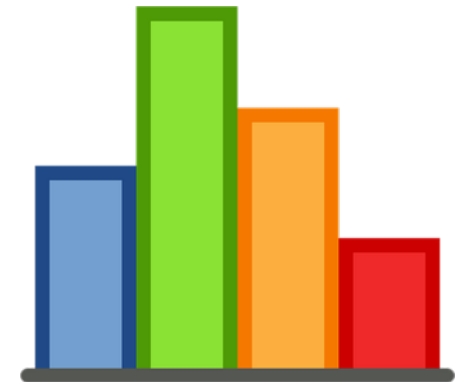
$Ber(p = 0.5)$



Probability
given model, predict data



$P(THHTHH)$



PROBABILITY VS STATISTICS



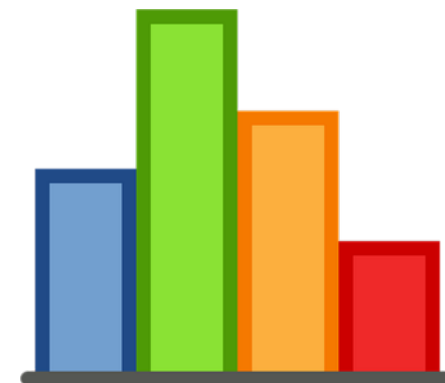
$Ber(p = 0.5)$



Probability
given model, predict data



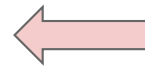
$P(THHTHH)$



$Ber(p = ???)$



Statistics
given data, predict model



$THHTHH$

RANDOM PICTURE



LIKELIHOOD (INTUITION)



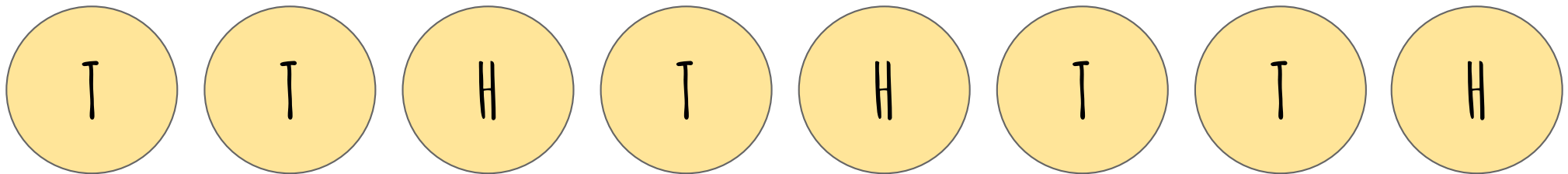
I give you and your classmates each 5 minutes with a coin with unknown probability of heads p . Whoever has the closest estimate will get an A+ in the class. What do you do in your precious 5 minutes, and what do you give as your estimate?

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?

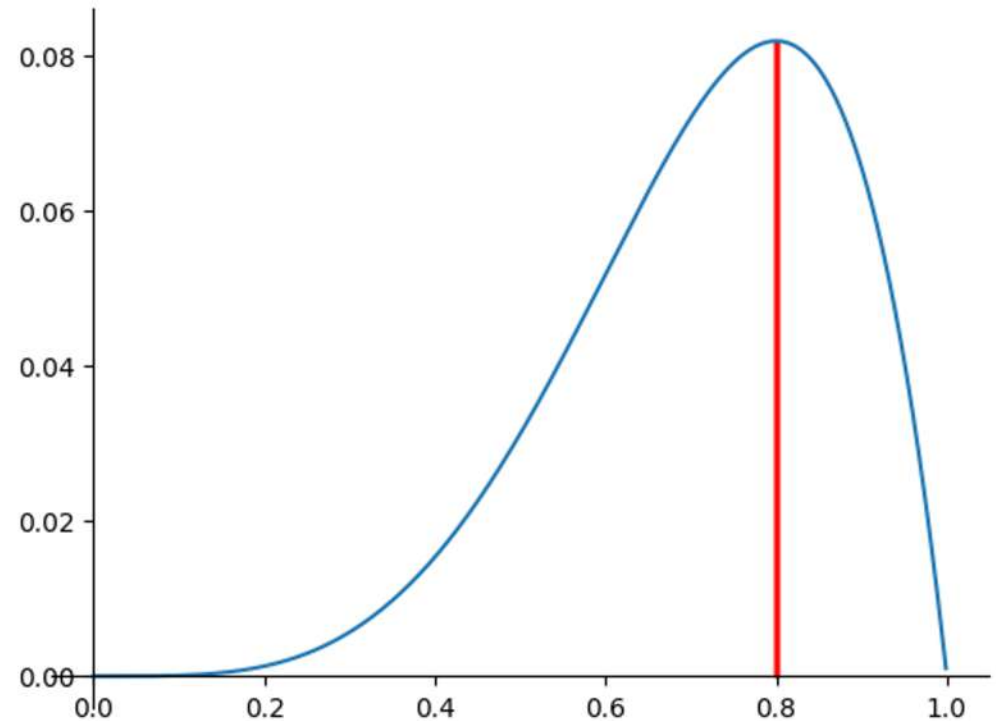
Flip it as many times as possible, and return $\# H / (\# H + \# T)$!



LIKELIHOOD (INTUITION)

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

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

LIKELIHOOD (INTUITION)



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

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

LIKELIHOOD (INTUITION)



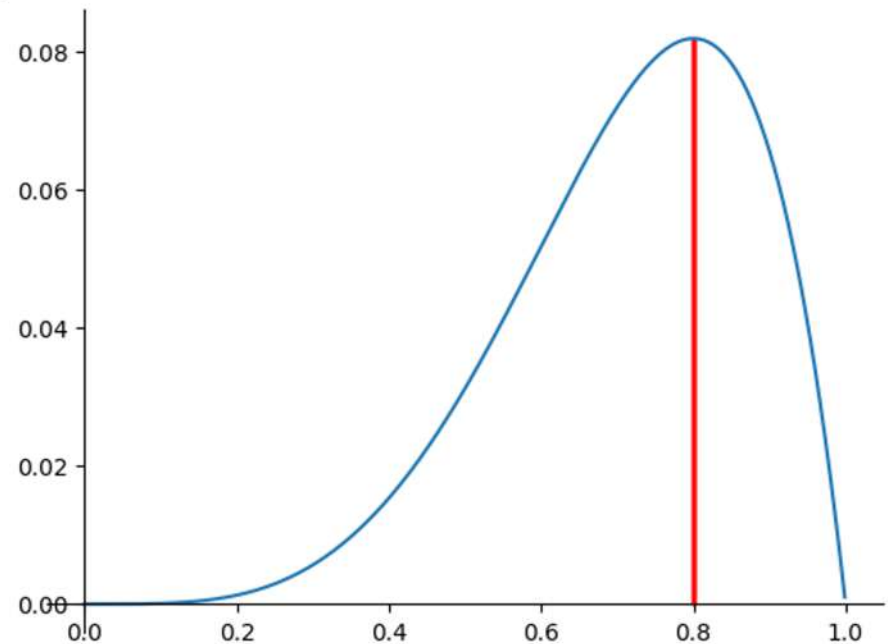
You assume a model (Bernoulli in our case) with unknown parameter θ , and receive iid samples $\mathbf{x} = (x_1, \dots, x_n) \sim \text{Ber}(\theta)$. The **likelihood** of the data given a parameter θ is

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

MAXIMUM LIKELIHOOD ESTIMATION (BERNOULLI)



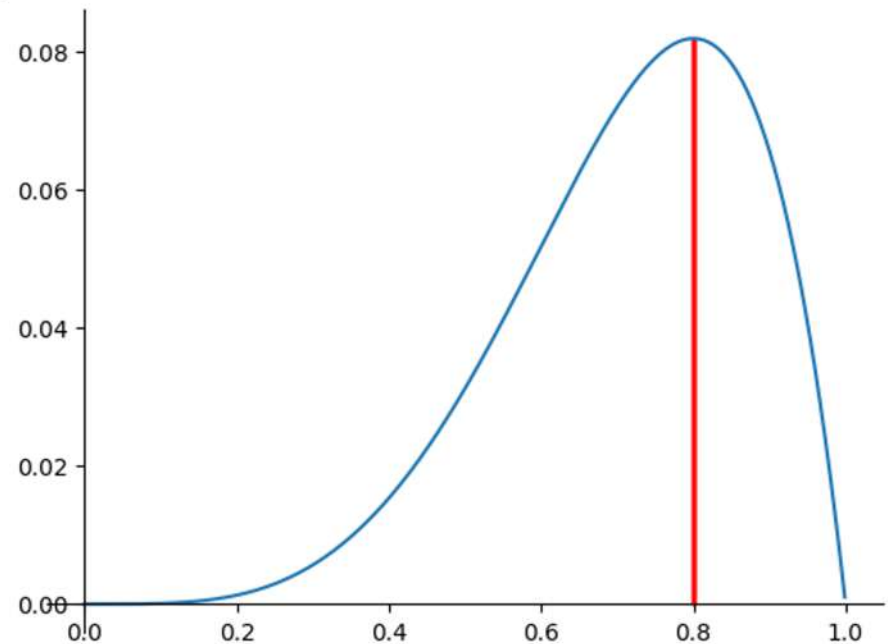
$$L(HHHHT \mid \theta) = \theta^4(1 - \theta) = \theta^4 - \theta^5$$



MAXIMUM LIKELIHOOD ESTIMATION (BERNOULLI)



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

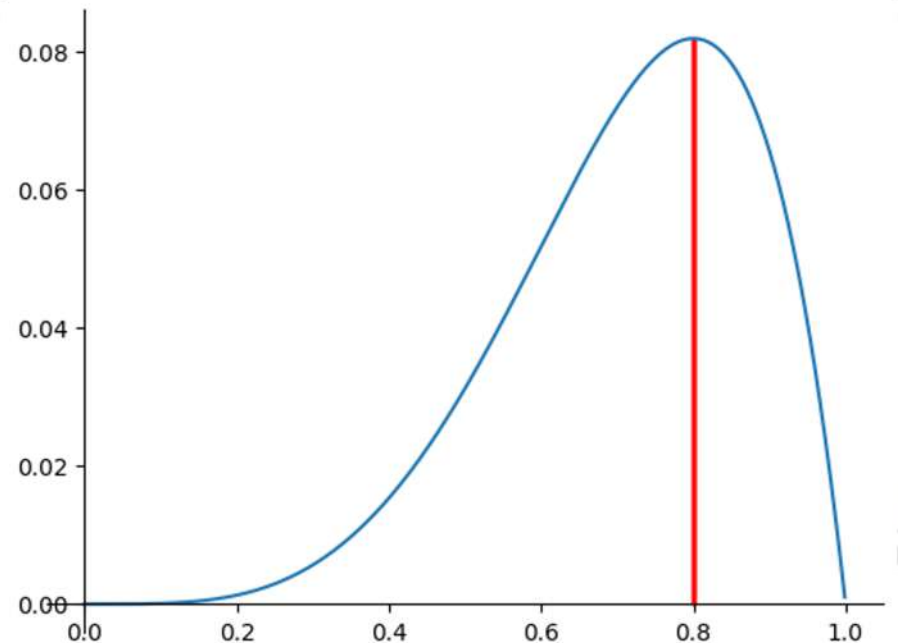


MAXIMUM LIKELIHOOD ESTIMATION (BERNOULLI)



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

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



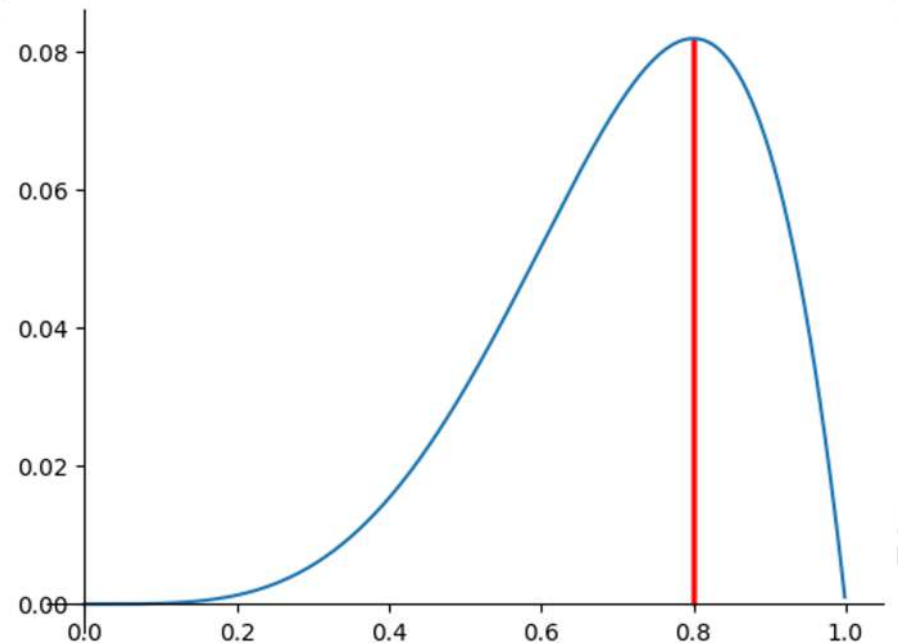
MAXIMUM LIKELIHOOD ESTIMATION (BERNOULLI)



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

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

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



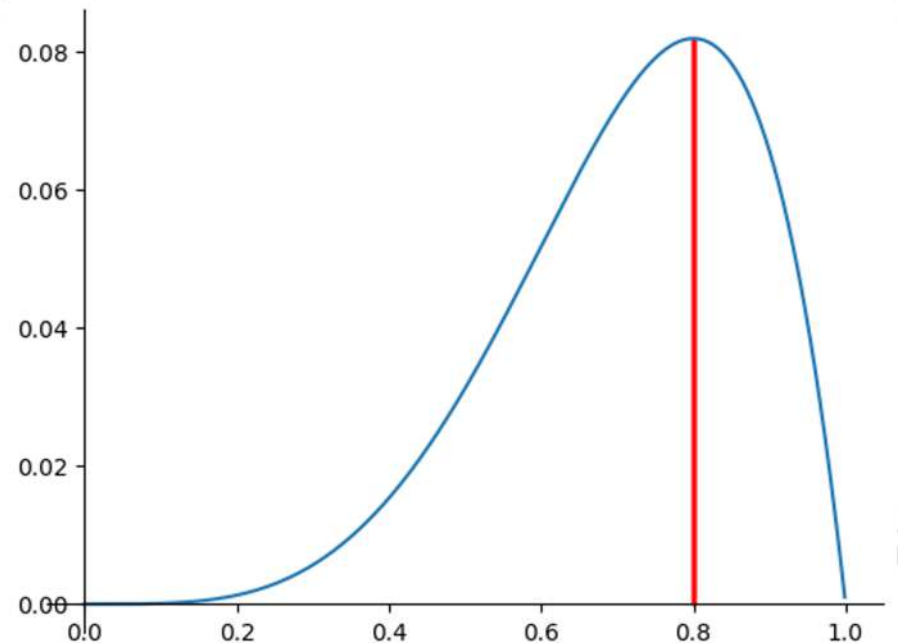
MAXIMUM LIKELIHOOD ESTIMATION (BERNOULLI)



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

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

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



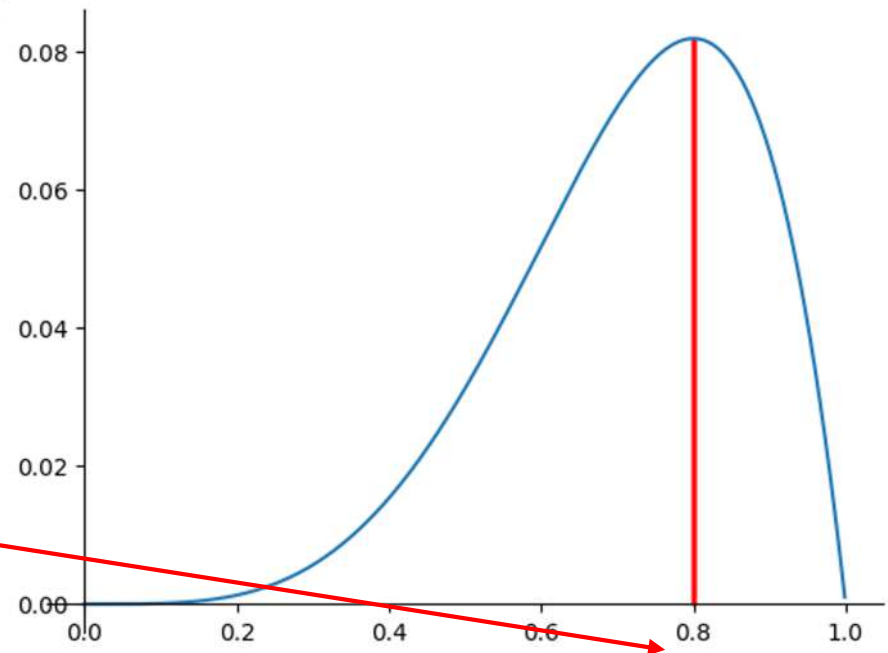
MAXIMUM LIKELIHOOD ESTIMATION (BERNOULLI)



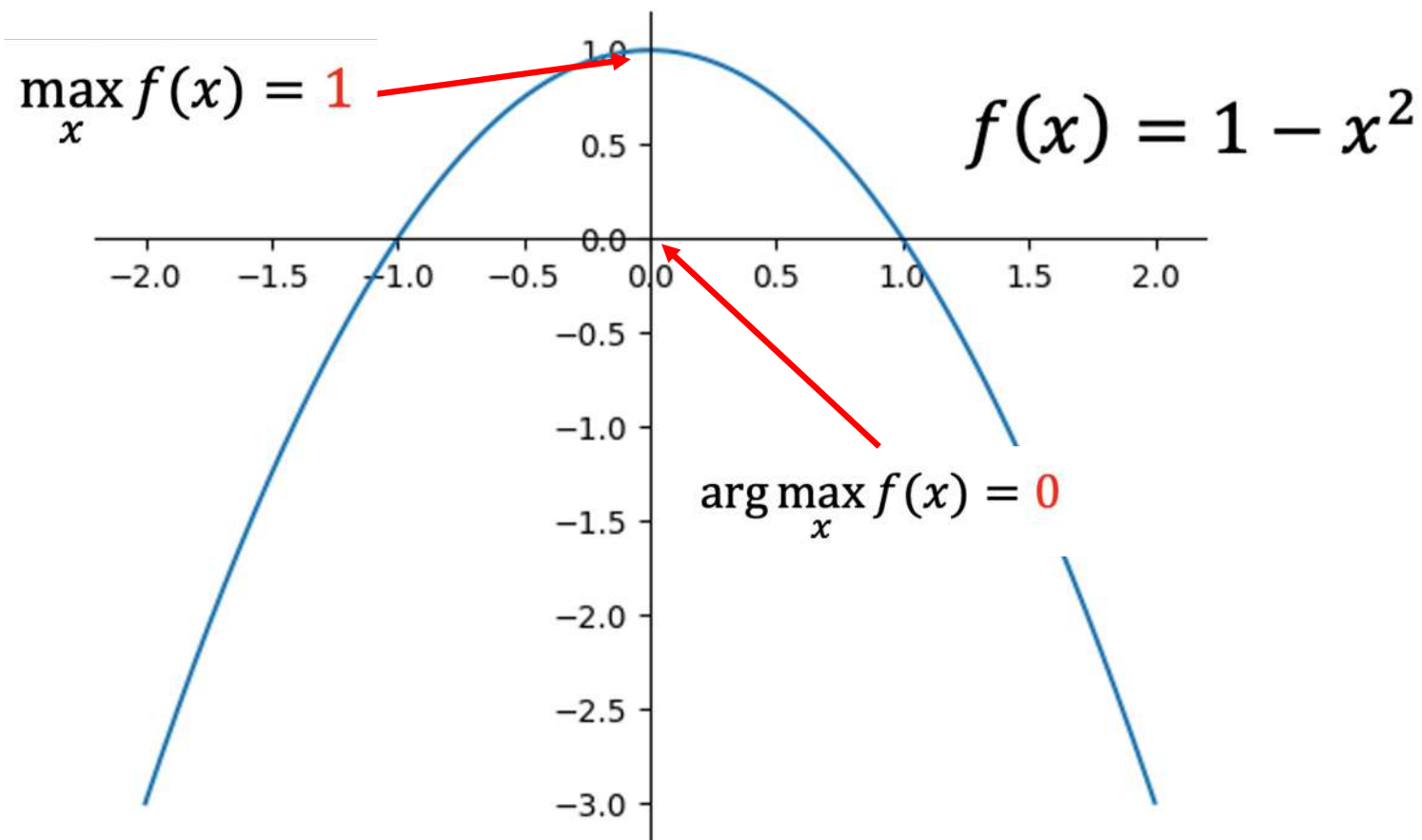
$$L(\text{HHHHT} \mid \theta) = \theta^4(1 - \theta) = \theta^4 - \theta^5$$

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

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



MAX VS 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 ?



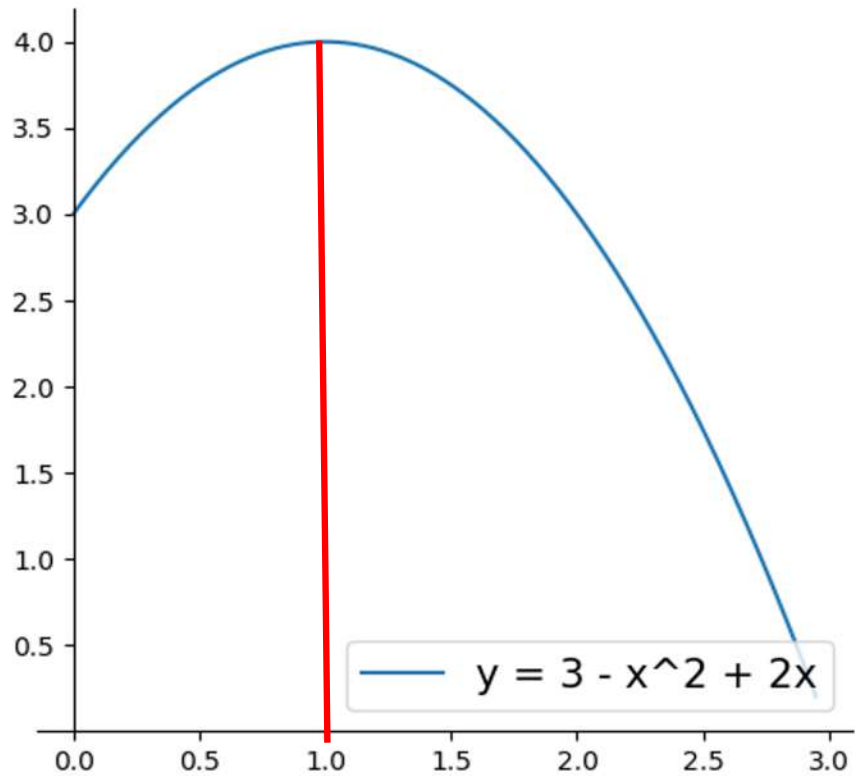
MAXIMUM LIKELIHOOD ESTIMATION (POISSON)



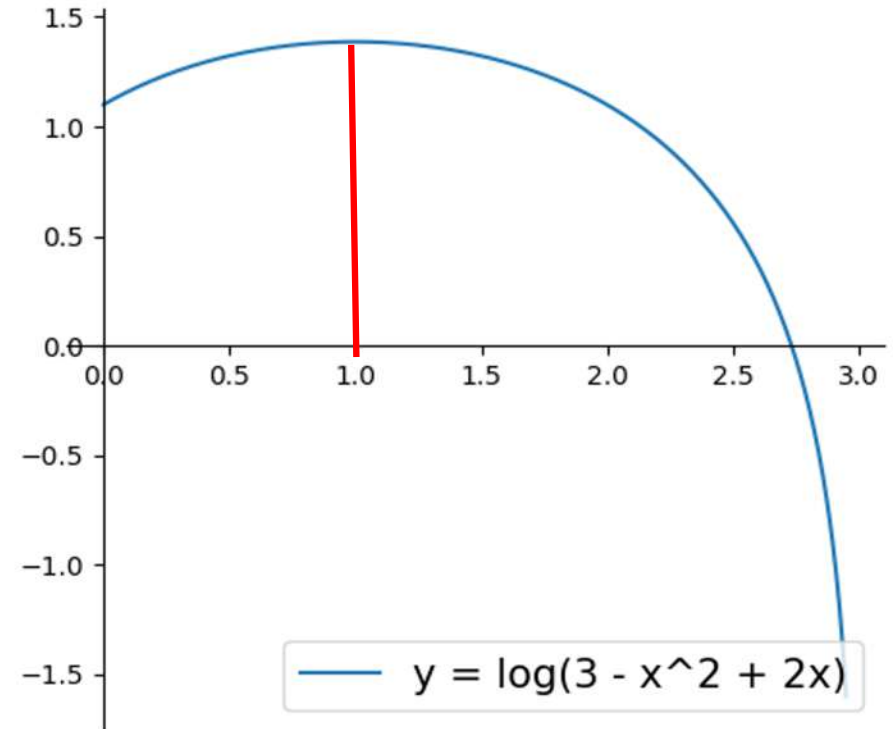
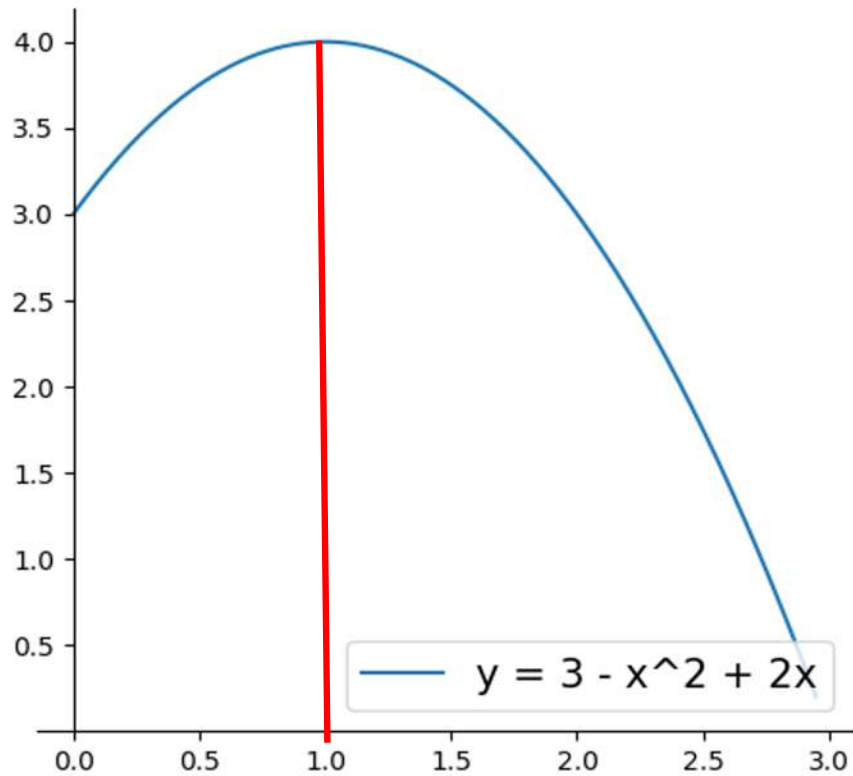
Let's say x_1, x_2, \dots, x_n are iid samples from $Poi(\theta)$. (might look like $x_1 = 3, x_2 = 5, x_3 = 4$, etc.) What is the MLE of θ ?

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

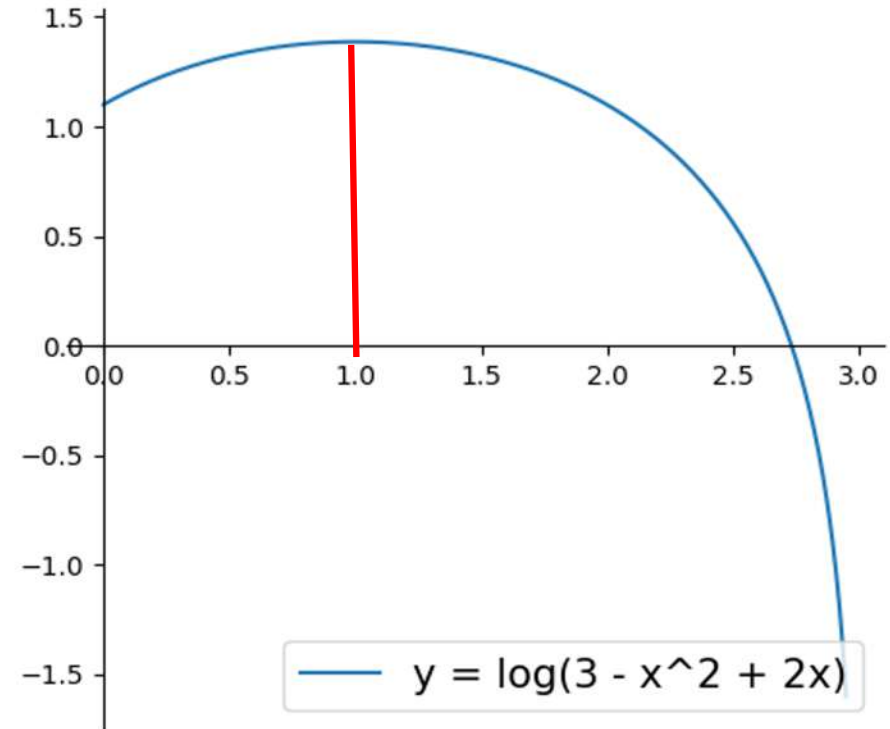
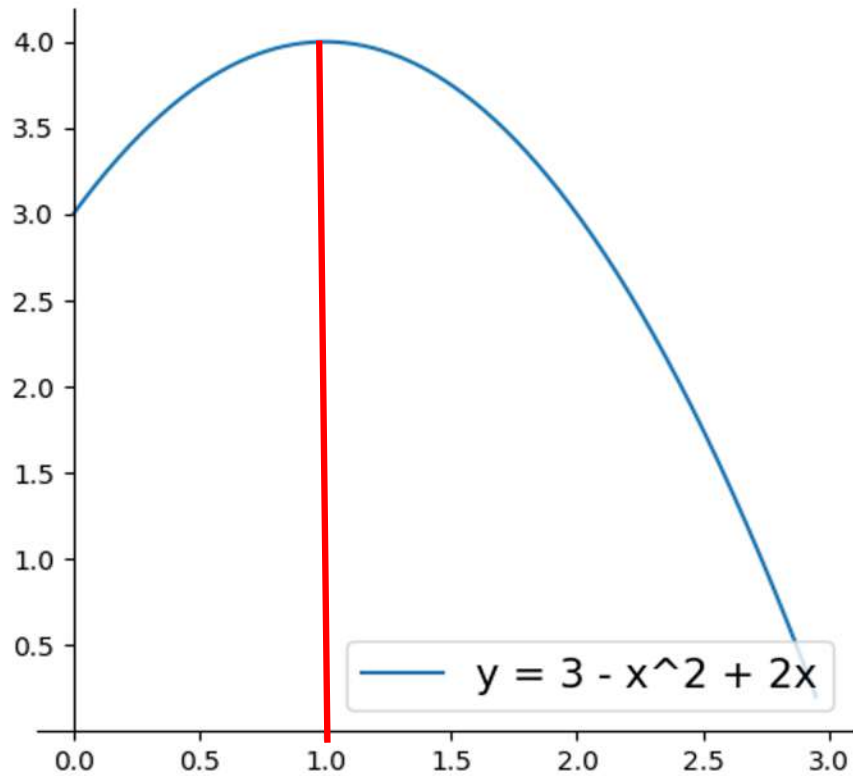
OPTIMIZING FUNCTION VS LOG(FUNCTION)



OPTIMIZING FUNCTION VS LOG(FUNCTION)



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 ?

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

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

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

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 ?

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

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

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

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



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 ?

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

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

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

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



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 ?

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

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

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

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

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

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,

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

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,

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

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).

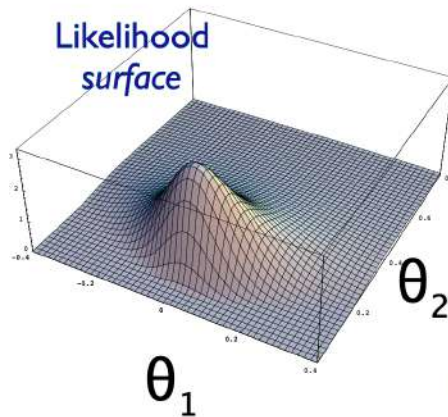
$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(\mathbf{x} | \theta) = \arg \max_{\theta} \ln L(\mathbf{x} | \theta)$$

RANDOM PICTURE



MAXIMUM LIKELIHOOD ESTIMATION (NORMAL)

$x_i \sim N(\mu, \sigma^2)$, μ, σ^2 both unknown



MAXIMUM LIKELIHOOD ESTIMATION (NORMAL)



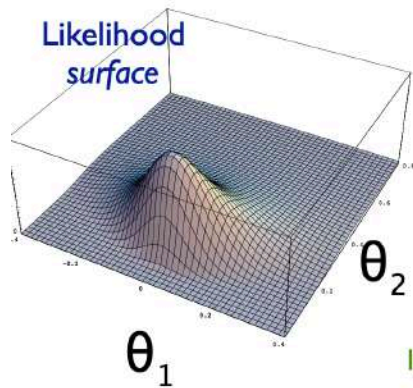
$x_i \sim N(\mu, \sigma^2)$, μ, σ^2 both unknown

$$\ln L(x_1, x_2, \dots, x_n | \theta_1, \theta_2) = \sum_{i=1}^n -\frac{1}{2} \ln(2\pi\theta_2) - \frac{(x_i - \theta_1)^2}{2\theta_2}$$

$$\frac{\partial}{\partial \theta_1} \ln L(x_1, x_2, \dots, x_n | \theta_1, \theta_2) = \sum_{i=1}^n \frac{(x_i - \theta_1)}{\theta_2} = 0$$

$$\hat{\theta}_1 = \left(\sum_{i=1}^n x_i \right) / n = \bar{x}$$

Sample mean is MLE of population mean, again



In general, a problem like this results in 2 equations in 2 unknowns.
Easy in this case, since θ_2 drops out of the $\partial/\partial\theta_1 = 0$ equation



$x_i \sim N(\mu, \sigma^2)$, μ, σ^2 both unknown

$$\ln L(x_1, x_2, \dots, x_n | \theta_1, \theta_2) = \sum_{i=1}^n -\frac{1}{2} \ln(2\pi\theta_2) - \frac{(x_i - \theta_1)^2}{2\theta_2}$$

$$\frac{\partial}{\partial \theta_2} \ln L(x_1, x_2, \dots, x_n | \theta_1, \theta_2) = \sum_{i=1}^n -\frac{1}{2} \frac{2\pi}{2\pi\theta_2} + \frac{(x_i - \theta_1)^2}{2\theta_2^2} = 0$$

$$\hat{\theta}_2 = \left(\sum_{i=1}^n (x_i - \hat{\theta}_1)^2 \right) / n = \bar{s}^2$$

**Sample variance is MLE of
population variance**