



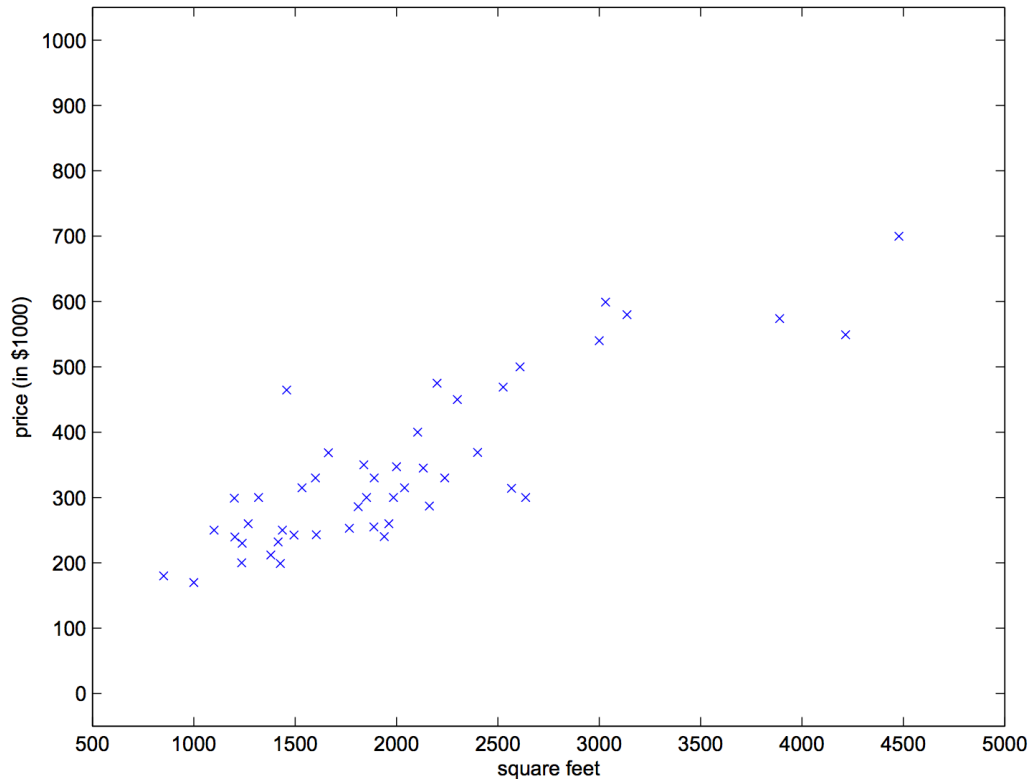
GUIDE TO THE GALAXY

Maximum Likelihood

CS 109
Lecture 21
May 13th, 2016

Predict Housing Prices

Living area (feet ²)	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
⋮	⋮



Review

Our Path

Neural Networks

Linear
Regression

Naive
Bayes

Logistic
Regression

Parameter Estimation

What are Parameters?

- Consider some probability distributions:
 - Ber(p) $\theta = p$
 - Poi(λ) $\theta = \lambda$
 - Uni(α, β) $\theta = (\alpha, \beta)$
 - Normal(μ, σ^2) $\theta = (\mu, \sigma^2)$
 - $Y = mX + b$ $\theta = (m, b)$
 - etc...
- Call these “parametric models”
- Given model, parameters yield actual distribution
 - Usually refer to parameters of distribution as θ
 - Note that θ that can be a vector of parameters

Our Path

Neural Networks

Linear
Regression

Naive
Bayes

Logistic
Regression

Unbiased
estimators

Method of
moments

Maximizing
likelihood

Bayesian
estimation

Parameter Estimation

Neural Networks

Linear
Regression

Naive
Bayes

Logistic
Regression

Unbiased
estimators

Method of
moments

Maximizing
likelihood

Bayesian
estimation

Parameter Estimation

Neural Networks

Linear
Regression

Naive
Bayes

Logistic
Regression

Unbiased
estimators

Method of
moments

Maximizing
likelihood

Bayesian
estimation

Recall Sample Mean + Variance?

- Consider n I.I.D. random variables X_1, X_2, \dots, X_n
 - X_j have distribution F with $E[X_j] = \mu$ and $\text{Var}(X_j) = \sigma^2$
 - We call sequence of X_j a **sample** from distribution F

- Recall sample mean: $\bar{X} = \sum_{i=1}^n \frac{X_i}{n}$ where $E[\bar{X}] = \mu$

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \text{ as } n \rightarrow \infty$$

- Recall sample variance:

$$S^2 = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{n-1} = \text{undefined}$$

Estimate parameters for
Bernoulli Poisson and
Normal

Method of Moments

- Recall: n -th moment of distribution for variable X :

$$m_n = E[X^n]$$

- Consider I.I.D. random variables X_1, X_2, \dots, X_n

$$\hat{m}_1 = \frac{1}{n} \sum_{i=1}^n X_i \quad \hat{m}_2 = \frac{1}{n} \sum_{i=1}^n X_i^2 \quad \dots \quad \hat{m}_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

are called the “sample moments”

- Estimates of the moments of distribution based on data
- Method of moments estimators
 - Estimate model parameters by equating “true” moments to sample moments:

$$m_i \approx \hat{m}_i$$

Estimate parameters for
Bernoulli Poisson and
Normal

Method of Moments with Uniform

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Uni}(\alpha, \beta)$

- Estimate mean:

$$\mu \approx \hat{m}_1 = \frac{1}{n} \sum_{i=1}^n X_i = \hat{\mu}$$

- Estimate variance:

$$\sigma^2 \approx \hat{m}_2 - (\hat{m}_1)^2 = \frac{\sum_{i=1}^n (X_i^2 - \bar{X}^2)}{n} = \hat{\sigma}^2$$

- For $\text{Uni}(\alpha, \beta)$, know that: $\mu = \frac{\alpha + \beta}{2}$ and $\sigma^2 = \frac{(\beta - \alpha)^2}{12}$

- Solve (two equations, two unknowns):

- Set $\beta = 2\mu - \alpha$, substitute into formula for σ^2 and solve:

$$\hat{\alpha} = \bar{X} - \sqrt{3}\hat{\sigma} \quad \text{and} \quad \hat{\beta} = \bar{X} + \sqrt{3}\hat{\sigma}$$

Method of Moments with Uniform



End Review

Parameter Estimation

Neural Networks

Linear
Regression

Naive
Bayes

Logistic
Regression

Unbiased
estimators

Method of
moments

Maximizing
likelihood

Bayesian
estimation

Parameter Estimation

Neural Networks

Linear
Regression

Naive
Bayes

Logistic
Regression

Unbiased
estimators

Method of
moments

Maximizing
likelihood

Bayesian
estimation

Great idea in Machine Learning

Likelihood of Data

- Consider n I.I.D. random variables X_1, X_2, \dots, X_n
 - X_i is a sample from density function $f(X_i | \theta)$
 - Note: now explicitly specify parameter θ of distribution
 - We want to determine how “likely” the observed data (x_1, x_2, \dots, x_n) is based on density $f(X_i | \theta)$
 - Define the **Likelihood function**, $L(\theta)$:

$$L(\theta) = \prod_{i=1}^n f(X_i | \theta)$$

- This is just a product since X_i are I.I.D.
 - Intuitively: what is probability of observed data using density function $f(X_i | \theta)$, for some choice of θ

[Demo](#)

Maximum Likelihood Estimator

- The **Maximum Likelihood Estimator** (MLE) of θ , is the value of θ that maximizes $L(\theta)$
 - More formally: $\theta_{MLE} = \arg \max_{\theta} L(\theta)$

Argmax

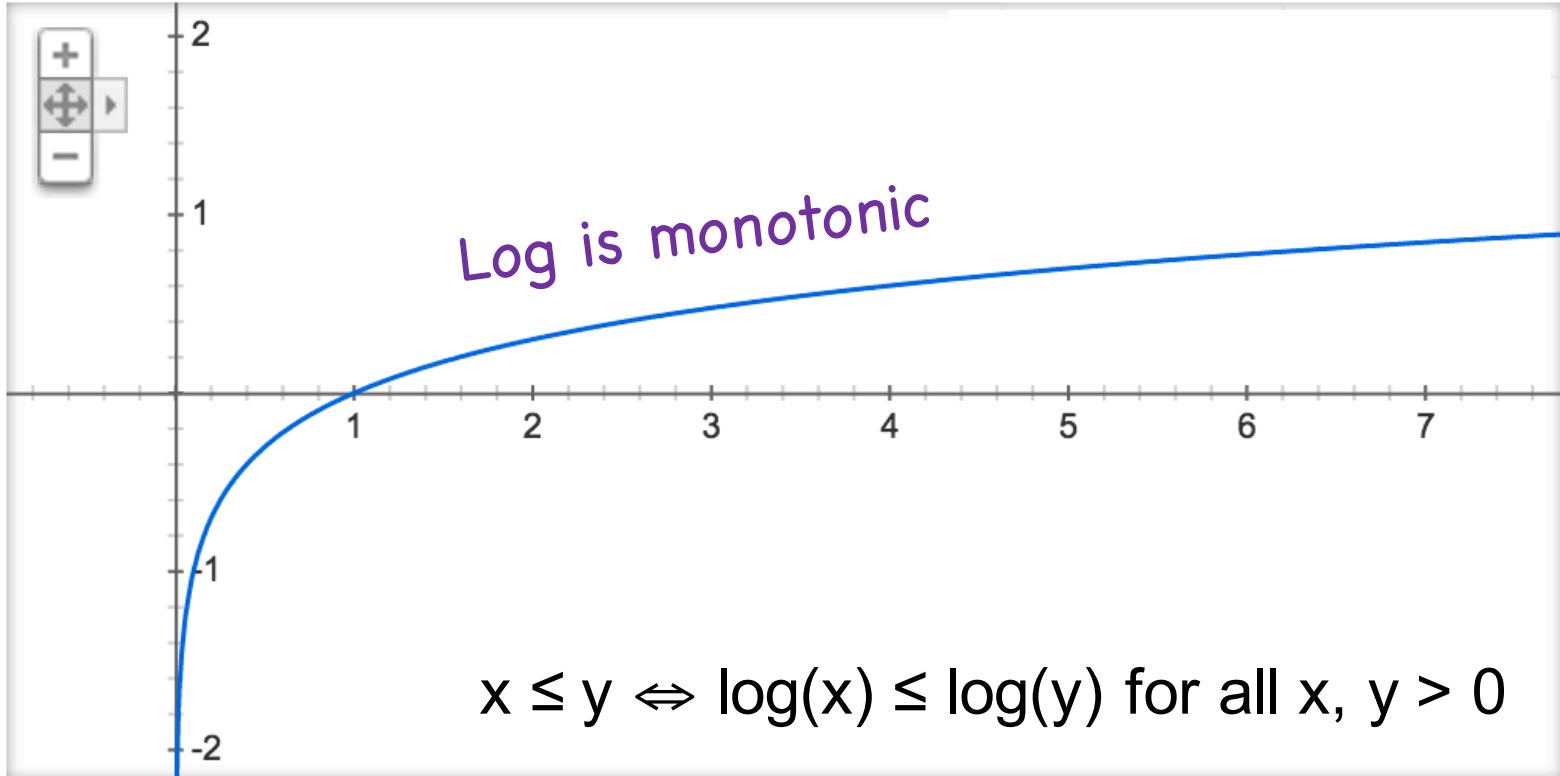
$$f(x) = -x^2 + 5$$

$$\max_x -x^2 + 5 = 5$$

$$\operatorname{argmax}_x -x^2 + 5 = 0$$

Argmax of Log

Graph for $\log(x)$

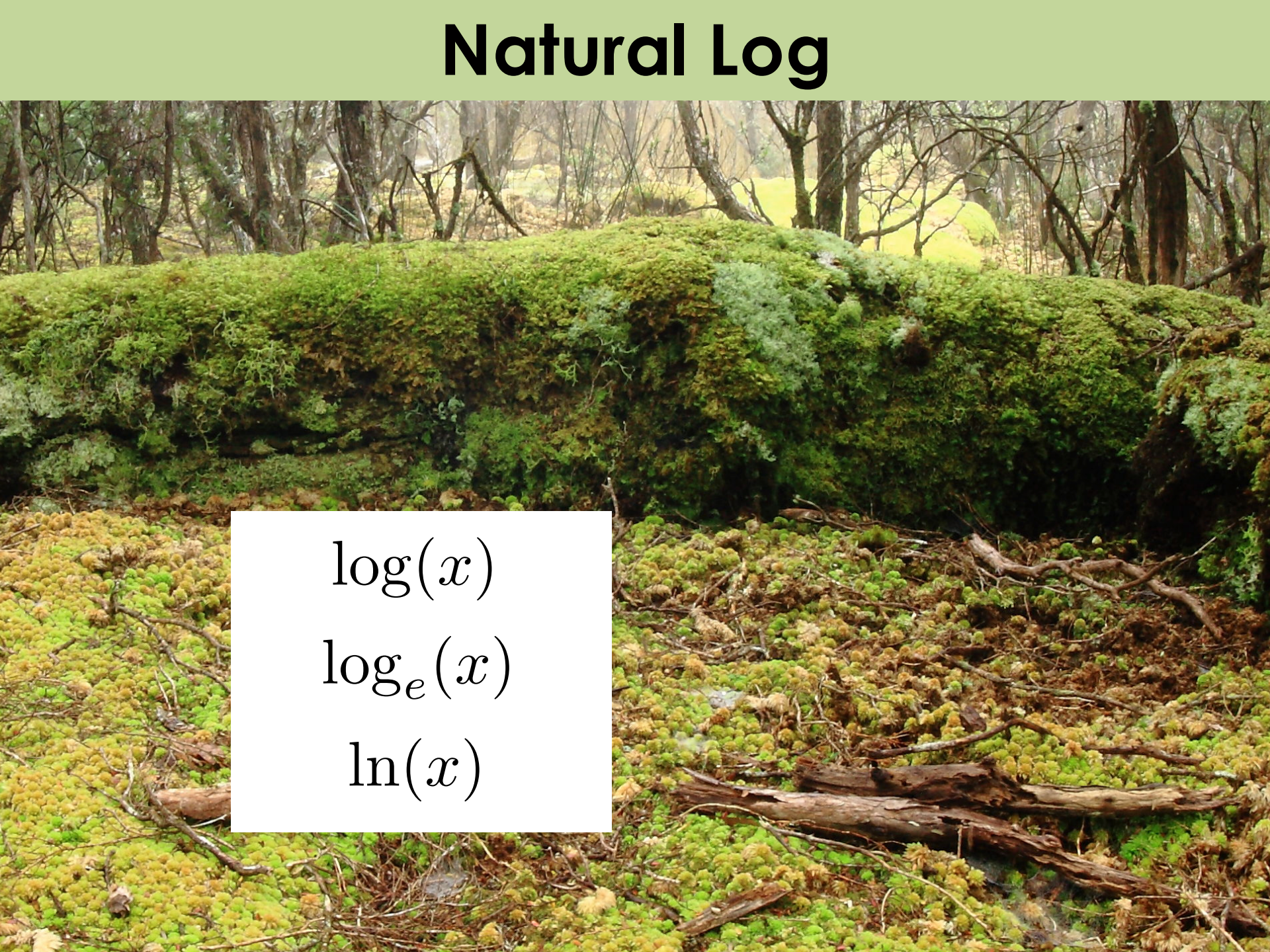


Claim: $\operatorname{argmax}_x f(x) = \operatorname{argmax}_x \log f(x)$

Log I Love You

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

Natural Log

A photograph of a forest floor covered in moss and fallen branches. The moss is a vibrant green, and the branches are brown and weathered. The background shows a dense forest of trees with bare branches, suggesting a late autumn or winter setting. A white text box is overlaid in the center of the image, containing the mathematical notations for the natural logarithm.

$\log(x)$
 $\log_e(x)$
 $\ln(x)$

Maximum Likelihood Estimator

- The **Maximum Likelihood Estimator** (MLE) of θ , is the value of θ that maximizes $L(\theta)$
 - More formally: $\theta_{MLE} = \arg \max_{\theta} L(\theta)$
 - More convenient to use **log-likelihood function**, $LL(\theta)$:

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

- Note that *log* function is “monotone” for positive values
 - Formally: $x \leq y \Leftrightarrow \log(x) \leq \log(y)$ for all $x, y > 0$
- So, θ that maximizes $LL(\theta)$ also maximizes $L(\theta)$
 - Formally: $\arg \max_{\theta} LL(\theta) = \arg \max_{\theta} L(\theta)$
 - Similarly, for any positive constant c (not dependent on θ):
$$\arg \max_{\theta} (c \cdot LL(\theta)) = \arg \max_{\theta} LL(\theta) = \arg \max_{\theta} L(\theta)$$

Computing the MLE

- General approach for finding MLE of θ
 - Determine formula for $LL(\theta)$
 - Differentiate $LL(\theta)$ w.r.t. (each) θ : $\frac{\partial LL(\theta)}{\partial \theta}$
 - To maximize, set $\frac{\partial LL(\theta)}{\partial \theta} = 0$
 - Solve resulting (simultaneous) equations to get θ_{MLE}
 - Make sure that derived $\hat{\theta}_{MLE}$ is actually a maximum (and not a minimum or saddle point). E.g., check $LL(\theta_{MLE} \pm \varepsilon) < LL(\theta_{MLE})$
 - This step often ignored in expository derivations
 - So, we'll ignore it here too (and won't require it in this class)

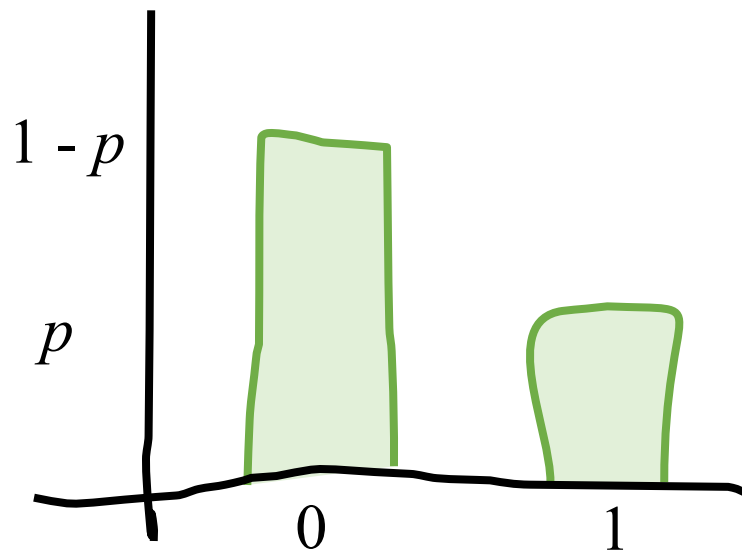
Maximizing Likelihood with Bernoulli

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_j \sim \text{Ber}(p)$

Maximizing Likelihood with Bernoulli

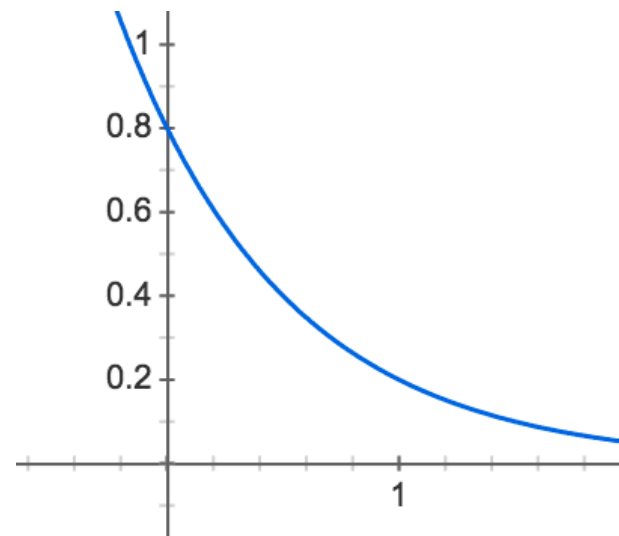
- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Ber}(p)$
 - Probability mass function, $f(X_i | p)$:

PMF of Bernoulli



$$f(X_i | p) = p^{x_i} (1 - p)^{1-x_i}$$

PMF of Bernoulli ($p = 0.2$)



$$f(x) = 0.2^x (1 - 0.2)^{1-x}$$

Maximizing Likelihood with Bernoulli

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_j \sim \text{Ber}(p)$
 - Probability mass function, $f(X_j | p)$, can be written as:

$$f(X_i | p) = p^{x_i} (1 - p)^{1-x_i} \quad \text{where } x_i = 0 \text{ or } 1$$

- Likelihood: $L(\theta) = \prod_{i=1}^n p^{X_i} (1 - p)^{1-X_i}$

- Log-likelihood:

$$\begin{aligned} LL(\theta) &= \sum_{i=1}^n \log(p^{X_i} (1 - p)^{1-X_i}) = \sum_{i=1}^n [X_i (\log p) + (1 - X_i) \log(1 - p)] \\ &= Y (\log p) + (n - Y) \log(1 - p) \quad \text{where } Y = \sum_{i=1}^n X_i \end{aligned}$$

- Differentiate w.r.t. p , and set to 0:

$$\frac{\partial LL(p)}{\partial p} = Y \frac{1}{p} + (n - Y) \frac{-1}{1 - p} = 0 \quad \Rightarrow \quad p_{MLE} = \frac{Y}{n} = \frac{1}{n} \sum_{i=1}^n X_i$$

Maximizing Likelihood with Poisson

- Consider I.I.D. random variables X_1, X_2, \dots, X_n

- $X_i \sim \text{Poi}(\lambda)$

- PMF: $f(X_i | \lambda) = \frac{e^{-\lambda} \lambda^{X_i}}{X_i!}$ Likelihood: $L(\theta) = \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{X_i}}{X_i!}$

- Log-likelihood:

$$\begin{aligned} LL(\theta) &= \sum_{i=1}^n \log\left(\frac{e^{-\lambda} \lambda^{X_i}}{X_i!}\right) = \sum_{i=1}^n \left[-\lambda \log(e) + X_i \log(\lambda) - \log(X_i!)\right] \\ &= -n\lambda + \log(\lambda) \sum_{i=1}^n X_i - \sum_{i=1}^n \log(X_i!) \end{aligned}$$

- Differentiate w.r.t. λ , and set to 0:

$$\frac{\partial LL(\lambda)}{\partial \lambda} = -n + \frac{1}{\lambda} \sum_{i=1}^n X_i = 0 \quad \Rightarrow \quad \lambda_{MLE} = \frac{1}{n} \sum_{i=1}^n X_i$$

Maximizing Likelihood with Normal

- Consider I.I.D. random variables X_1, X_2, \dots, X_n

- $X_i \sim N(\mu, \sigma^2)$

- PDF: $f(X_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(X_i - \mu)^2 / (2\sigma^2)}$

- Log-likelihood:

$$LL(\theta) = \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi}\sigma} e^{-(X_i - \mu)^2 / (2\sigma^2)}\right) = \sum_{i=1}^n \left[-\log(\sqrt{2\pi}\sigma) - (X_i - \mu)^2 / (2\sigma^2) \right]$$

- First, differentiate w.r.t. μ , and set to 0:

$$\frac{\partial LL(\mu, \sigma^2)}{\partial \mu} = \sum_{i=1}^n 2(X_i - \mu) / (2\sigma^2) = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = 0$$

- Then, differentiate w.r.t. σ , and set to 0:

$$\frac{\partial LL(\mu, \sigma^2)}{\partial \sigma} = \sum_{i=1}^n -\frac{1}{\sigma} + 2(X_i - \mu)^2 / (2\sigma^3) = -\frac{n}{\sigma} + \sum_{i=1}^n (X_i - \mu)^2 / (\sigma^3) = 0$$

Being Normal, Simultaneously

- Now have two equations, two unknowns:

$$\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = 0 \quad -\frac{n}{\sigma} + \sum_{i=1}^n (X_i - \mu)^2 / (\sigma^3) = 0$$

- First, solve for μ_{MLE} :

$$\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = 0 \Rightarrow \sum_{i=1}^n X_i = n\mu \Rightarrow \mu_{MLE} = \frac{1}{n} \sum_{i=1}^n X_i$$

- Then, solve for σ^2_{MLE} :

$$-\frac{n}{\sigma} + \sum_{i=1}^n (X_i - \mu)^2 / (\sigma^3) = 0 \Rightarrow n\sigma^2 = \sum_{i=1}^n (X_i - \mu)^2$$

$$\sigma^2_{MLE} = \frac{1}{n} \sum_{i=1}^n (X_i - \mu_{MLE})^2$$

- Note: μ_{MLE} unbiased, but σ^2_{MLE} biased (same as MOM)

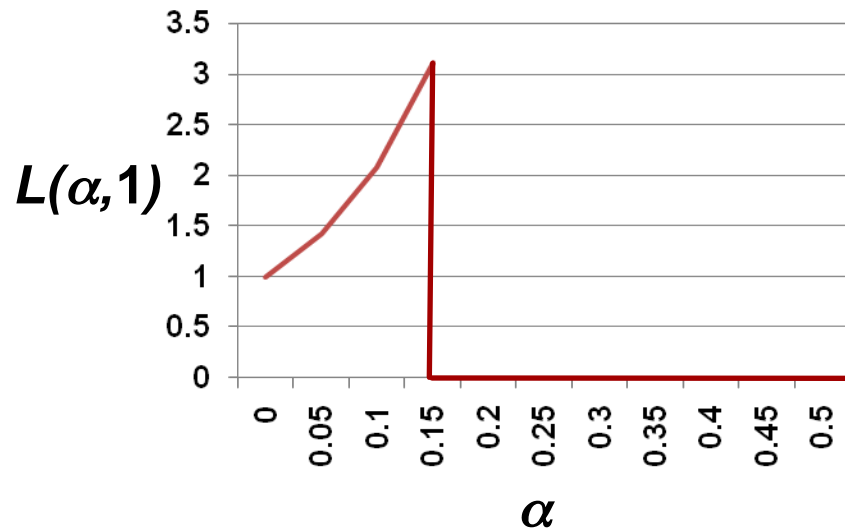
Maximizing Likelihood with Uniform

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Uni}(\alpha, \beta)$
 - PDF: $f(X_i | \alpha, \beta) = \begin{cases} \frac{1}{\beta - \alpha} & \alpha \leq x_i \leq \beta \\ 0 & \text{otherwise} \end{cases}$
 - Likelihood: $L(\theta) = \begin{cases} \left(\frac{1}{\beta - \alpha}\right)^n & \alpha \leq x_1, x_2, \dots, x_n \leq \beta \\ 0 & \text{otherwise} \end{cases}$
 - Constraint $\alpha \leq x_1, x_2, \dots, x_n \leq \beta$ makes differentiation tricky
 - Intuition: want interval size $(\beta - \alpha)$ to be as small as possible to maximize likelihood function for each data point
 - But need to make sure all observed data contained in interval
 - If all observed data not in interval, then $L(\theta) = 0$
 - Solution: $\alpha_{MLE} = \min(x_1, \dots, x_n)$ $\beta_{MLE} = \max(x_1, \dots, x_n)$

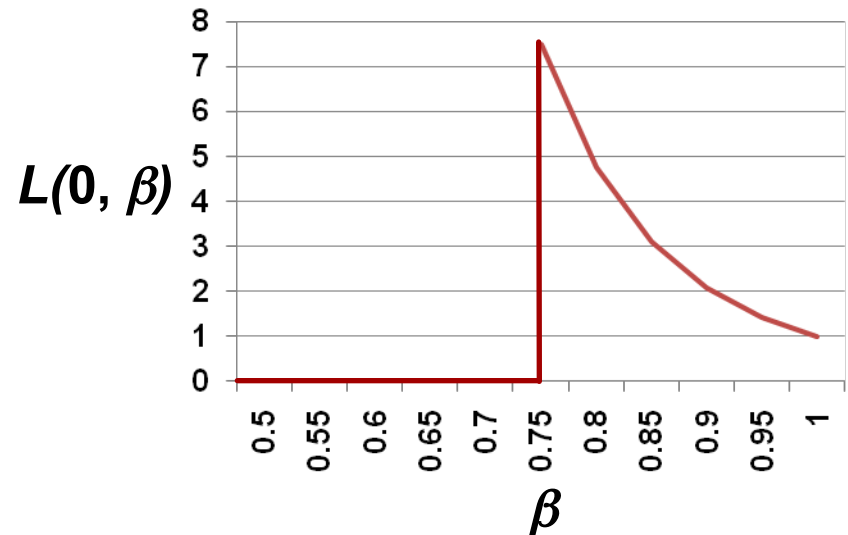
Understanding MLE with Uniform

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Uni}(0, 1)$
 - Observe data:
 - 0.15, 0.20, 0.30, 0.40, 0.65, 0.70, 0.75

Likelihood: $L(\alpha, 1)$



Likelihood: $L(0, \beta)$



Once Again, Small Samples = Problems

- How do small samples affect MLE?

- In many cases, $\mu_{MLE} = \frac{1}{n} \sum_{i=1}^n X_i = \text{sample mean}$
 - Unbiased. Not too shabby...
- As seen with Normal, $\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu_{MLE})^2$
 - Biased. Underestimates for small n (e.g., 0 for $n = 1$)
- As seen with Uniform, $\alpha_{MLE} \geq \alpha$ and $\beta_{MLE} \leq \beta$
 - Biased. Problematic for small n (e.g., $\alpha = \beta$ when $n = 1$)
- Small sample phenomena intuitively make sense:
 - Maximum likelihood \Rightarrow best explain data we've seen
 - Does not attempt to generalize to unseen data

Properties of MLE

- Maximum Likelihood Estimators are generally:
 - Consistent: $\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| < \varepsilon) = 1$ for $\varepsilon > 0$
 - Potentially biased (though asymptotically less so)
 - Asymptotically optimal
 - Has smallest variance of “good” estimators for large samples
 - Often used in practice where sample size is large relative to parameter space
 - But be careful, there are some very large parameter spaces

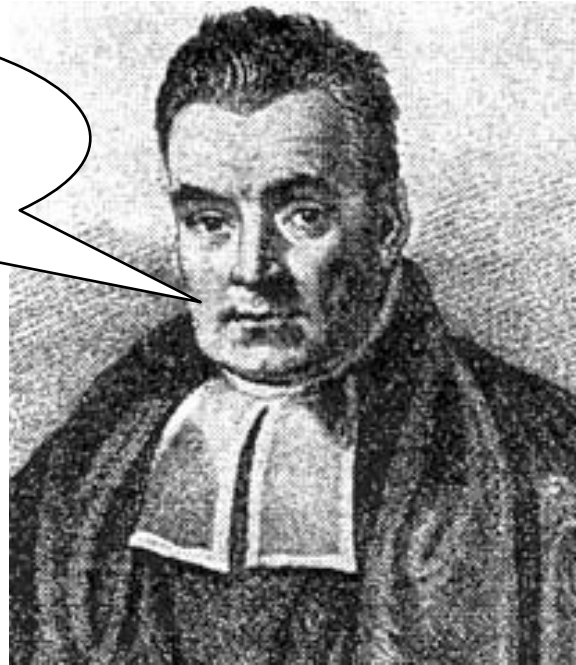
[on board, MLE of line]

From probability theory

To ML algorithm

Need a Volunteer

So good to see
you again!



Two Envelopes

- I have two envelopes, will allow you to have one
 - One contains $\$X$, the other contains $\$2X$
 - Select an envelope
 - Open it!
 - Now, would you like to switch for other envelope?
 - To help you decide, compute $E[\$ \text{ in other envelope}]$
 - Let $Y = \$ \text{ in envelope you selected}$
$$E[\$ \text{ in other envelope}] = \frac{1}{2} \cdot \frac{Y}{2} + \frac{1}{2} \cdot 2Y = \frac{5}{4} Y$$
 - Before opening envelope, think either equally good
 - So, what happened by opening envelope?
 - And does it really make sense to switch?