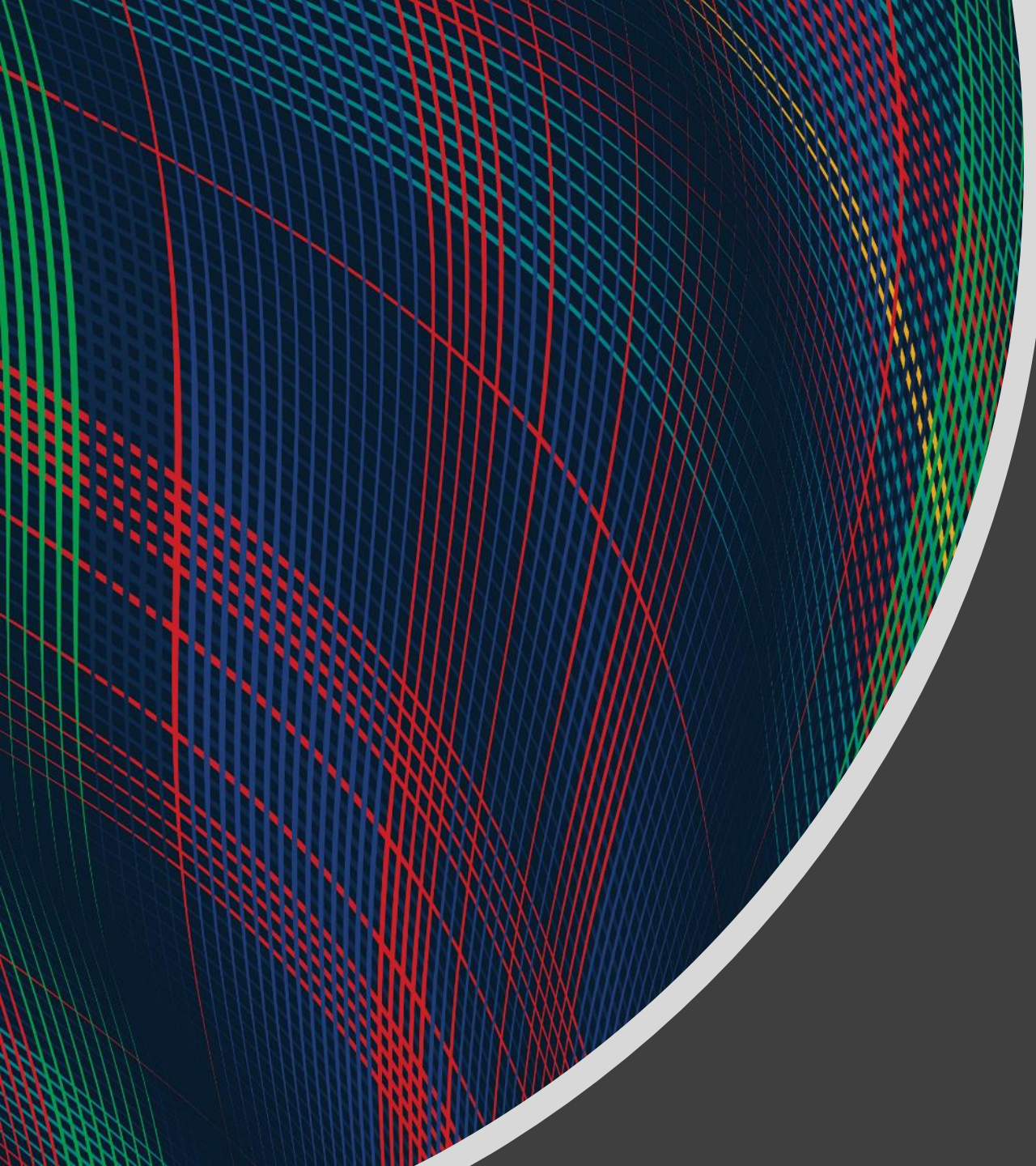# Announcements

## HW3

- Out late tonight
- Due Sat 10/9

## Quizzes

- Mon 10/4,   last 15 min. of class (calculus, optimization, Lagrange)
- Mon 10/11, last 15 min. of class (probability, statistics)

# Mathematical Foundations for Machine Learning

## Statistics

Instructor: Pat Virtue

# Plan

# Likelihood

**Likelihood**: The probability (or density) of random variable $Y$ taking on value $y$ given the distribution parameters, $\boldsymbol{\theta}$.

$$P(Y = y \mid \theta)$$

Additional notation

$$P(y \mid \theta)$$

$$P(y ; \theta)$$

$$P_\theta(y)$$

$$P(D \mid \theta)$$

# Likelihood and i.i.d

$$D = \{y^{(1)}, y^{(2)}, y^{(3)}\}$$

**Likelihood**: The probability (or density) of random variable $Y$ taking on value $y$ given the distribution parameters, $\boldsymbol{\theta}$.

**i.i.d.**: Independent and identically distributed

$$P\left(Y^{(1)} = y^{(1)},\ Y^{(2)} = y^{(2)},\ Y^{(3)} = y^{(3)} \mid \theta^{(1)}, \theta^{(2)}, \theta^{(3)}\right)$$

identical $\downarrow$

$$P\left(Y = y^{(1)},\ Y = y^{(2)},\ Y = y^{(3)} \mid \theta\right)$$

independent $\downarrow$

$$= P\left(Y = y^{(1)} \mid \theta\right) P\left(Y = y^{(2)} \mid \theta\right) P\left(Y = y^{(3)} \mid \theta\right)$$

# Bernoulli Likelihood

$P(\mathcal{D} \mid \theta)$

Bernoulli distribution:

$$Y \sim Bern(\phi) \qquad p(y \mid \phi) = \begin{cases} \phi, & y = 1 \\ 1 - \phi, & y = 0 \end{cases}$$

What is the likelihood for three i.i.d. samples, given parameter $\phi$:

$$\mathcal{D} = \{y^{(1)} = 1, y^{(2)} = 1, y^{(3)} = 0\}$$

$$P(\mathcal{D} \mid \phi) = \prod_{i=1}^{N} P(Y = y^{(i)} \mid \phi)$$

$$= \phi \cdot \phi \cdot (1 - \phi)$$

# MLE

# Estimating Parameters with Likelihood

We model the outcome of a single mysterious weighted-coin flip as a Bernoulli random variable:

$$Y \sim Bern(\phi)$$

$$p(y \mid \phi) = \begin{cases} \phi, & y = 1 \ (Heads) \\ 1 - \phi, & y = 0 \ (Tails) \end{cases}$$

Given the ordered sequence of coin flip outcomes:

$$[1, 0, 1, 1]$$

What is the estimate of parameter $\hat{\phi}$?
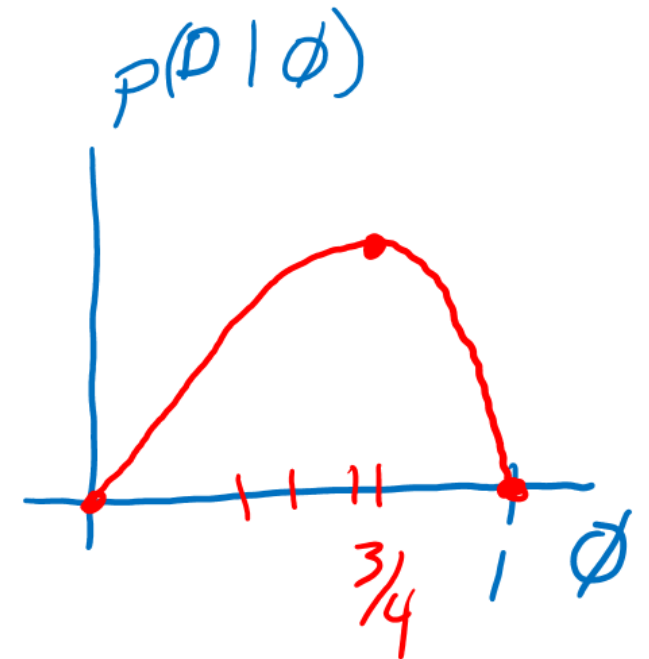
# Estimating Parameters with Likelihood

We model the outcome of a single mysterious weighted-coin flip as a
Bernoulli random variable:

$$Y \sim Bern(\phi)$$

$$p(y \mid \phi) = \begin{cases} \phi, & y = 1 \ (Heads) \\ 1 - \phi, & y = 0 \ (Tails) \end{cases}$$

Given the ordered sequence of coin flip outcomes:

$$[1, 0, 1, 1]$$

What is the estimate of parameter $\hat{\phi}$?

$$p(\,D \mid \phi\,) = \phi \cdot \phi \cdot (1 - \phi) \cdot \phi$$
$$= \phi^3 (1 - \phi)^1$$

# Likelihood and Maximum Likelihood Estimation

**Likelihood**: The probability (or density) of random variable $Y$ taking on value $y$ given the distribution parameters, $\boldsymbol{\theta}$.

**Likelihood function**: The value of likelihood as we change theta

(same as likelihood, but conceptually we are considering many different values of the parameters)

**Maximum Likelihood Estimation (MLE)**: Find the parameter value that maximizes the likelihood.

# MLE as Data Increases

Given the ordered sequence of coin flip outcomes:
$$[1, 0, 1, 1]$$

$$p(\mathcal{D} \mid \phi) = \prod_{i}^{N} p(y^{(i)} \mid \phi) = \phi^{N_{y=1}}(1-\phi)^{N_{y=0}}$$

What happens as we flip more coins?

# MLE for Gaussian

Gaussian distribution:

$$Y \sim \mathcal{N}(\mu, \sigma^2)$$

$$p(y \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

What is the log likelihood for three i.i.d. samples, given parameters $\mu, \sigma^2$?

$$\mathcal{D} = \{y^{(1)} = 65, y^{(2)} = 95, y^{(3)} = 85\}$$

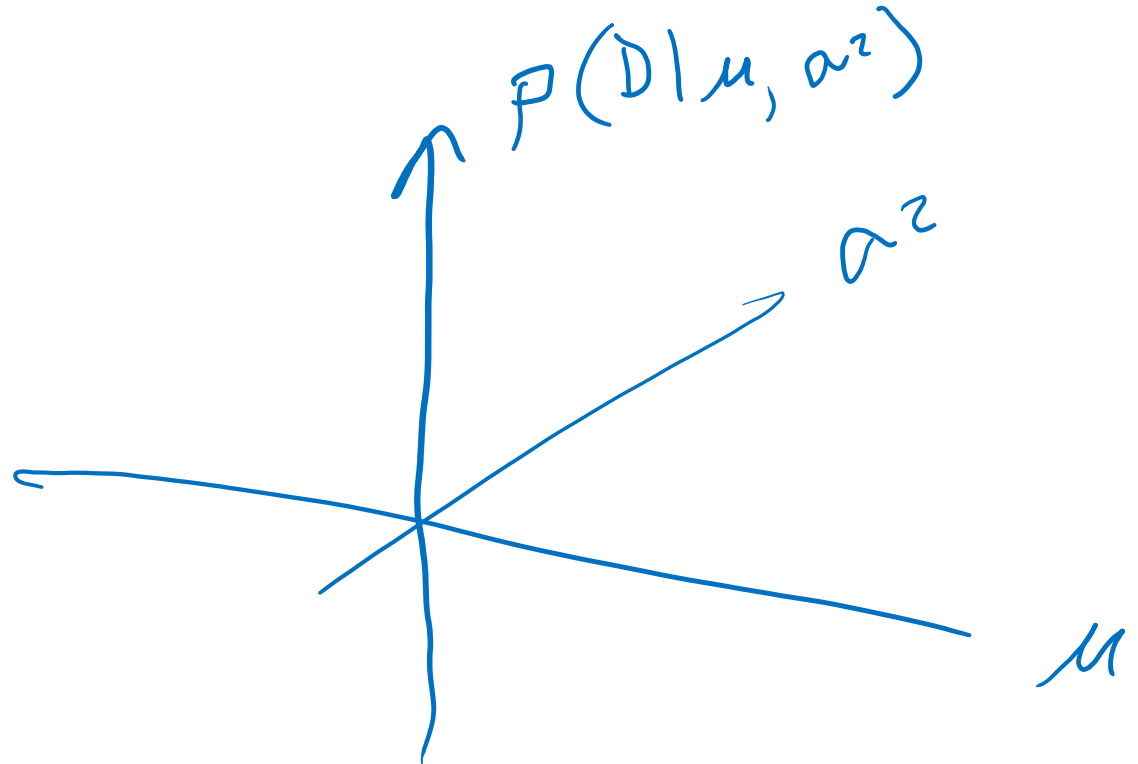$$L(\mu, \sigma^2) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\left(y^{(i)}-\mu\right)^2}{2\sigma^2}}$$

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_{i}^{N} p\left(y^{(i)} \mid \boldsymbol{\theta}\right)$$

# MLE for Gaussian

Assume that exam scores are drawn independently from the same Gaussian (Normal) distribution.

Given three exam scores 75, 80, 90, which pair of parameters is the best fit (the highest likelihood)?

$$p(\mathcal{D}|\mu,\sigma^2) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)}-\mu)^2}{2\sigma^2}}$$
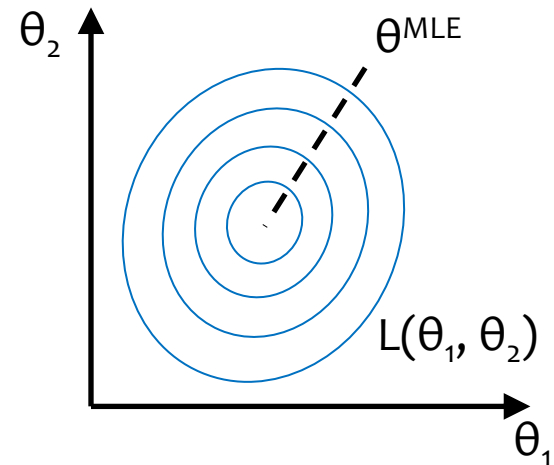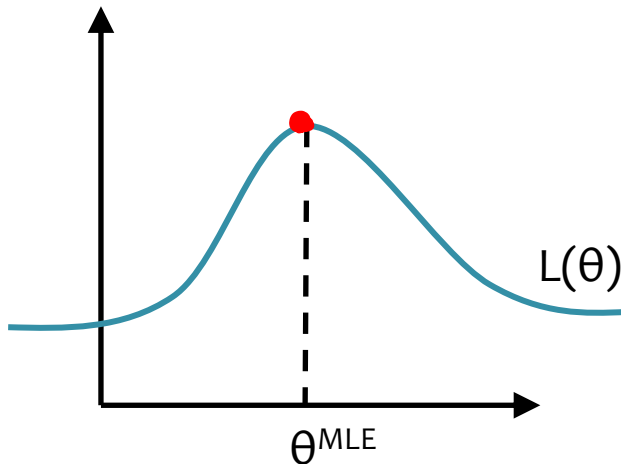
# MLE

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likelihood Estimation:**
Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\text{MLE}} = \underset{\boldsymbol{\theta}}{\text{argmax}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)



L($\theta$)

$\theta^{\text{MLE}}$



$\theta_2$

$\theta^{\text{MLE}}$

L($\theta_1$, $\theta_2$)

$\theta_1$

# MLE Recipe

# Likelihood and Log Likelihood

$$\log xz = \log x + \log z$$

**Likelihood**: The probability (or density) of random variable $Y$ taking on value $y$ given the distribution parameters, $\boldsymbol{\theta}$.

$$P\left(D \mid \theta\right) = \Pi\, p\left(y^{(i)} \mid \theta\right)$$

**Likelihood function**: The value of likelihood as we change theta

(same as likelihood, but conceptually we are considering many different values of the parameters)

$$\text{likelihood } \mathcal{L}\left(\theta ; D\right) = P\left(D \mid \theta\right) = \Pi\, p\left(y^{(i)} \mid \theta\right)$$

$$\log \text{ likelihood } \ell\left(\theta ; D\right) = \log P\left(D \mid \theta\right) = \Sigma\, \log p\left(y^{(i)} \mid \theta\right)$$

# Maximum Likelihood Estimation

MLE of parameter $\theta$ for i.i.d. dataset $\mathcal{D} = \left\{ y^{(i)} \right\}_{i=1}^{N}$

$$\hat{\theta}_{MLE} = \underset{\theta}{\mathrm{argmax}} \, p(\mathcal{D} \mid \theta)$$

# Recipe for Estimation

## MLE

1. Formulate the likelihood, $p(\mathcal{D} \mid \theta)$

2. Set objective $J(\theta)$ equal to negative log likelihood

$$J(\theta) = -\log p(\mathcal{D} \mid \theta)$$

3. Compute derivative of objective, $\partial J / \partial \theta$

4. Find $\hat{\theta}$, either

    a. Set derivate equal to zero and solve for $\theta$

    b. Use (stochastic) gradient descent to step towards better $\theta$

# MLE for Gaussian

Gaussian distribution:

$$Y \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

$$p\left(y \mid \mu, \sigma^2\right) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

What is the log likelihood for three i.i.d. samples, given parameters $\mu, \sigma^2$?

$$\mathcal{D} = \{y^{(1)} = 65, y^{(2)} = 95, y^{(3)} = 85\}$$

$$L\left(\mu, \sigma^2\right) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\left(y^{(i)}-\mu\right)^2}{2\sigma^2}}$$

$$\hat{\theta}_{MLE} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_i p\left(y^{(i)} \mid \boldsymbol{\theta}\right)$$

$$\ell\left(\mu, \sigma^2\right) = \sum_{i=1}^{N} -\log\sqrt{2\pi\sigma^2} - \frac{\left(y^{(i)} - \mu\right)^2}{2\sigma^2}$$

$$\hat{\theta}_{MLE} = \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_i \log p\left(y^{(i)} \mid \boldsymbol{\theta}\right)$$

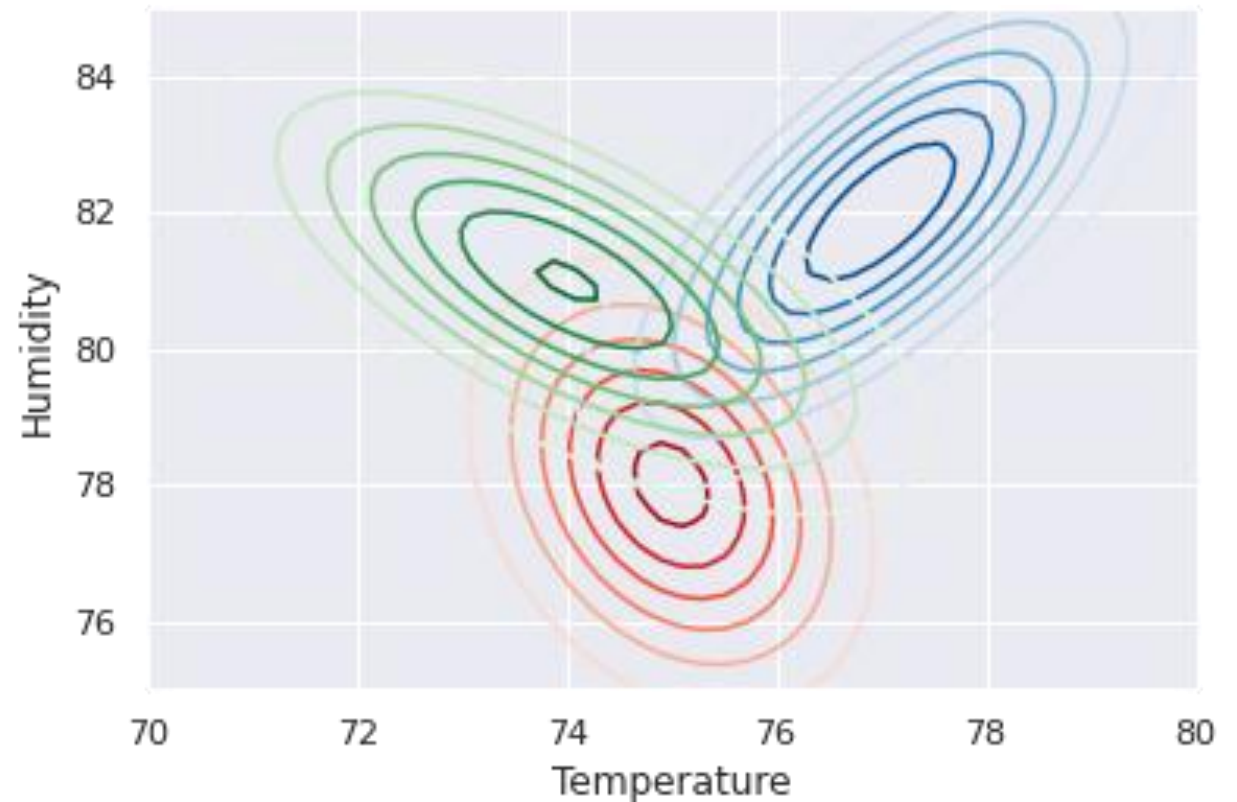# Exercise: MLE for Exponential

Exponential distribution pdf:
$$f(x) = \lambda e^{-\lambda x}$$

# Two Applications of Bayes Rule
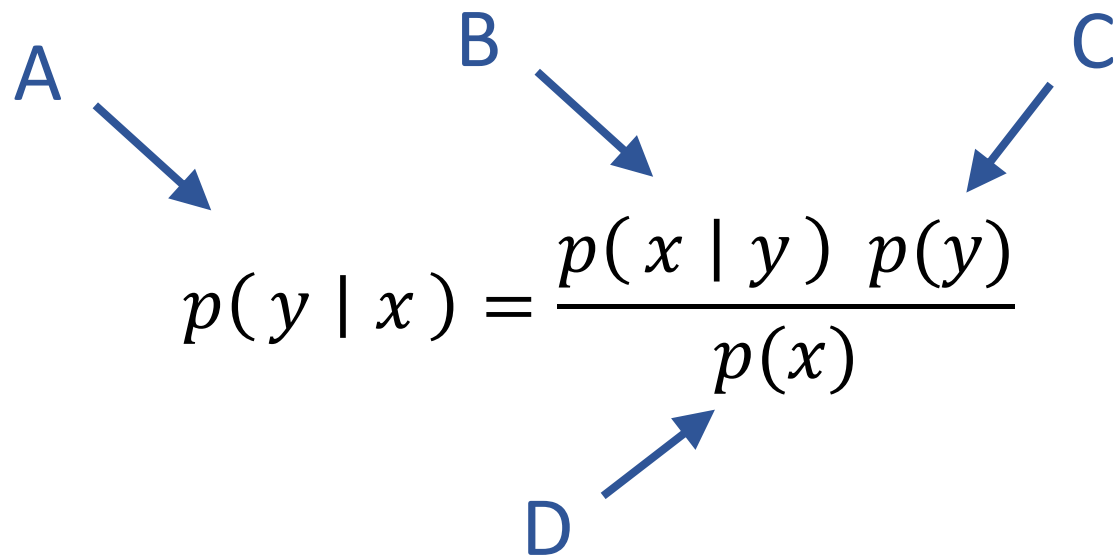
# Prev Recitation

## City classification

There are three nearby cities: $\textcolor{red}{\text{City A}}$, $\textcolor{blue}{\text{City B}}$, and $\textcolor{green}{\text{City C}}$. Let $Y \sim Categorical(a, b, c)$ be a categorical distribution where $Y = 1$ means a randomly sampled sensor is in $\textcolor{red}{\text{City A}}$, $Y = 2$ means it is in $\textcolor{blue}{\text{City B}}$, and $Y = 3$ means it is in $\textcolor{green}{\text{City C}}$.

# Poll 2

Which of these terms is the likelihood?

A

B

C

$$p(y \mid x) = \frac{p(x \mid y) \; p(y)}{p(x)}$$

D

# Bayes Rule

Terminology

$$p(y \mid x) = \frac{p(x \mid y)p(y)}{p(x)}$$

# Bayes Rule

Inserting parameters

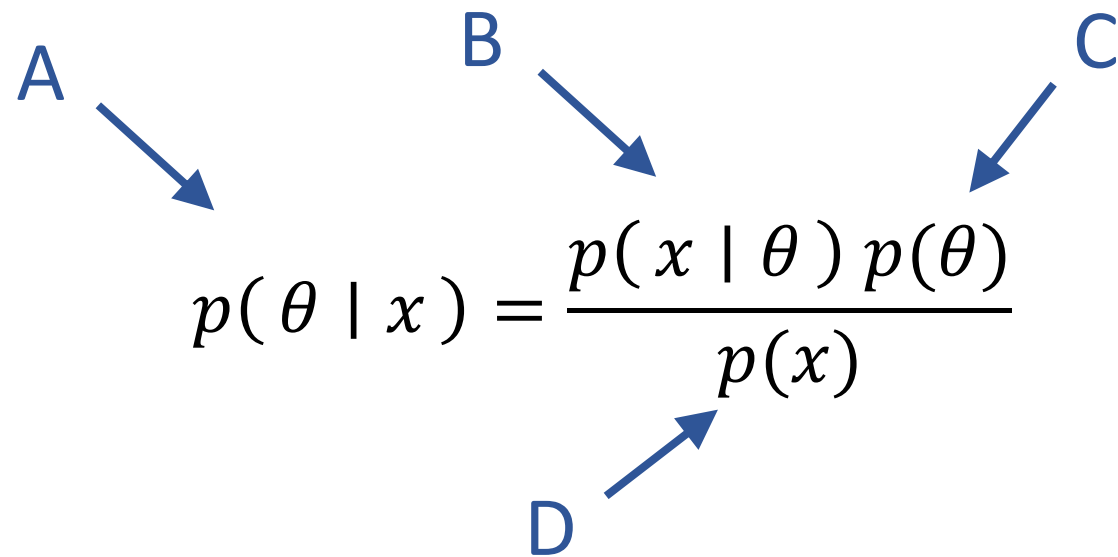$$p(\,y\mid x\,) = \frac{p(\,x\mid y\,)\,p(y)}{p(x)}$$

# Bayes Rule

Another way to use Bayes rule

$$p(\theta \mid x) = \frac{p(x \mid \theta)\, p(\theta)}{p(x)}$$

# Poll 3

Where do we plug in the pdf, $f(x) = \lambda e^{-\lambda x}$

A

B

C

D

$$p(\,\theta \mid x\,) = \frac{p(\,x \mid \theta\,)\;p(\theta)}{p(x)}$$

# MLE vs. MAP

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

$$\boldsymbol{\theta}^{\text{MLE}} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)

$$\boldsymbol{\theta}^{\text{MAP}} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)} | \boldsymbol{\theta}) p(\boldsymbol{\theta})$$

Maximum *a posteriori* (MAP) estimate

Prior

# Expectation and Variance

# Expectation and Variance

The **expected value** of $X$ is $E[X]$. Also called the mean.

- Discrete random variables:

  Suppose $X$ can take any value in the set $\mathcal{X}$.

  $$E[X] = \sum_{x \in \mathcal{X}} x p(x)$$

- Continuous random variables:

  $$E[X] = \int_{-\infty}^{+\infty} x f(x) dx$$

# Expectation and Variance

The **variance** of $X$ is *Var(X).*

$$Var(X) = E[(X - E[X])^2]$$

$$\mu = E[X]$$

- Discrete random variables:

$$Var(X) = \sum_{x \in \mathcal{X}} (x - \mu)^2 p(x)$$

- Continuous random variables:

$$Var(X) = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx$$