

# Announcements

- pset1 out today, due in two weeks.
  - Submit on Gradescope
- No screens (laptops, tablets, phones) during the class.

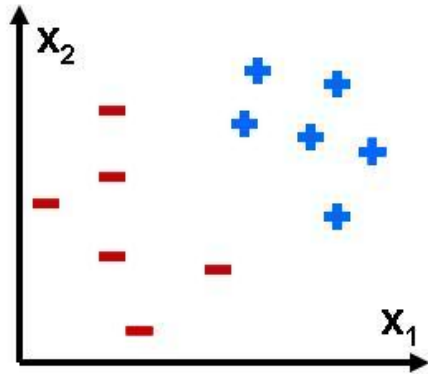
# Last time: Supervised Learning

Output

**Categorical**



**Classification**



$y \in \{0,1\}$   
0: Negative Class  
1: Positive Class

**Continuous**






**Regression**

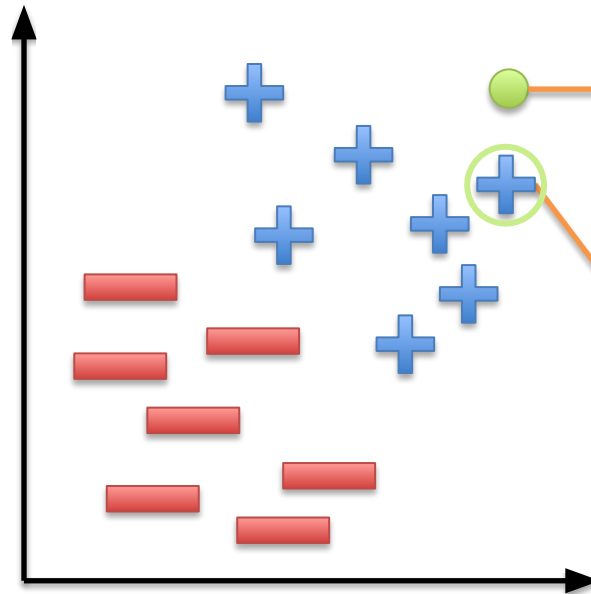


# Nearest Neighbor Classifier



Training data

-  = Panda
-  = Not panda
-  = Test sample



## Variants of this algorithm:

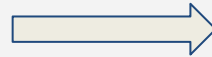
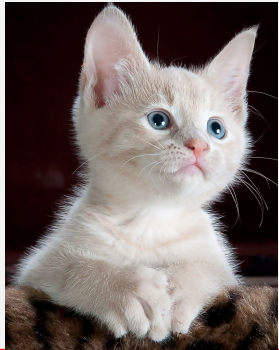
- K-nearest neighbor classifier (k-NN)
- Approx nearest neighbor

# K-Nearest Neighbor Classifier

## Training Data



Feature vector

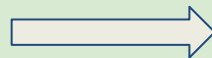


Feature vector



$$\sum_{i=1}^d (x_i - y_i)^2$$

L2  
distance



Feature vector



# Diving deeper into model analysis

	Predicted "1"	Predicted "0"
GT Label "1"	<b>True Positive (TP)</b>	<b>False Negative (FN)</b>
GT Label "0"	<b>False Positive (FP)</b>	<b>True Negative (TN)</b>

slido



## What is precision?

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.



## What is precision?

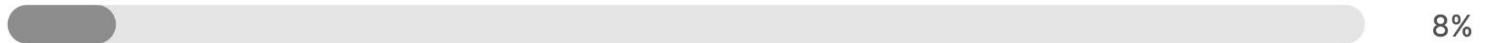
Quiz question   77 answers   77 participants

TP / (TP + FP) - 64 answers



83%

TP / (TP + FN) - 6 answers



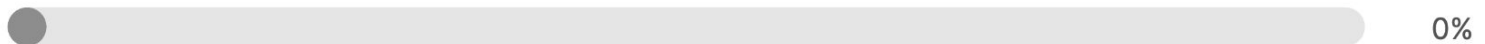
8%

TP / (FP + TN) - 7 answers



9%

TP / (FN + TN) - 0 answers



0%

# Precision and Recall

- **Precision:**

- a. The percentage of predictions that are correct
- b. The ability to identify **only** relevant data points

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **Recall:**

- a. The percentage of relevant data points that are correctly identified
- b. The ability to identify **all** relevant data points

$$\text{Recall} = \frac{TP}{TP + FN}$$

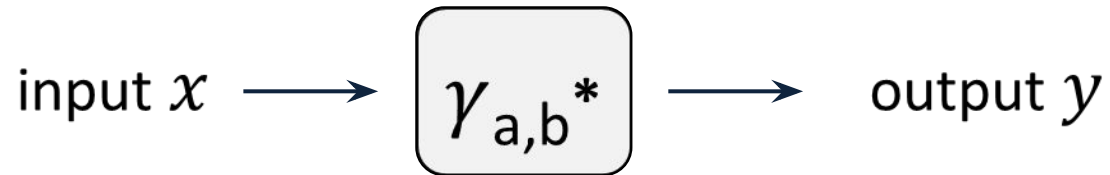


# Hypothesis $\gamma$

$\gamma$  : function parametrized by  $\theta$ , e.g.,

$$\gamma(x) = \text{sign}(\underbrace{a}_{\theta_{0,1}}x + \underbrace{b}_{\theta_2})$$

**Goal:**



# Linear Regression

Linear hypothesis:

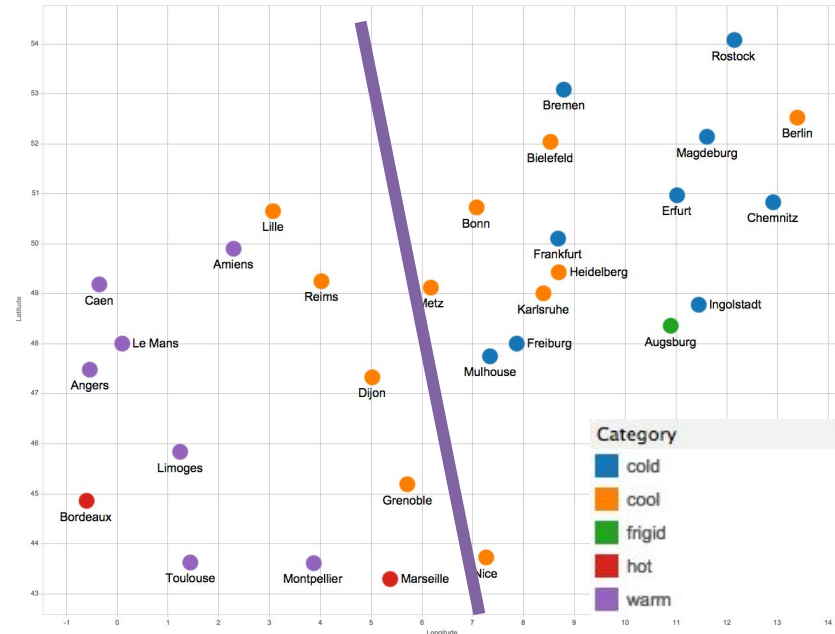
$$\gamma_{a,b}(x) = \text{sign}(ax + b)$$

Error Function:

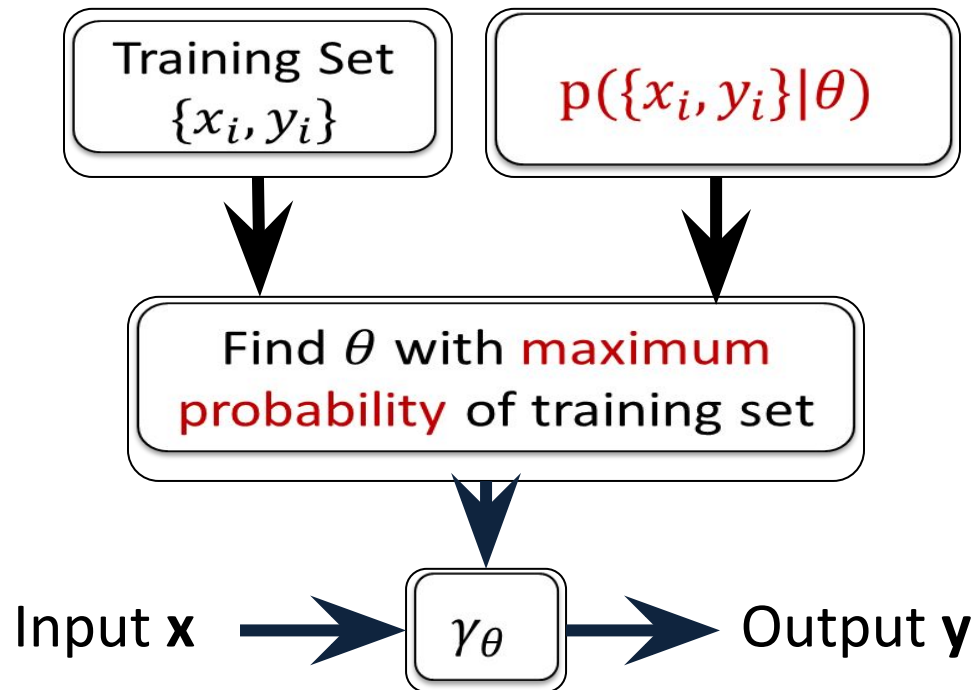
Portion of incorrect predictions

$$\text{Error}(\gamma_{a,b}, D\{x, y\}) = \frac{1}{N} \sum_{i=1}^N \gamma_{a,b}(x_i) \neq y_i$$

**Goal:** minimize  $\text{Error}(\gamma_{a,b}, D\{x, y\})$



# Alternative View: “Maximum Likelihood”



# Maximum likelihood way of estimating model parameters $\theta$

In general, assume data is generated by some distribution

$$U \sim p(U|\theta)$$

Observations (i.i.d.)

$$D = \{u^{(1)}, u^{(2)}, \dots, u^{(m)}\}$$

Maximum likelihood estimate

$$\mathcal{L}(D) = \prod_{i=1}^N p(u^{(i)}|\theta)$$

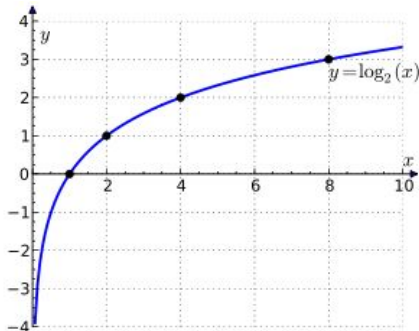
Likelihood

$$\theta_{ML} = \operatorname{argmax}_{\theta} \mathcal{L}(D)$$

Log likelihood

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log p(u^{(i)}|\theta)$$

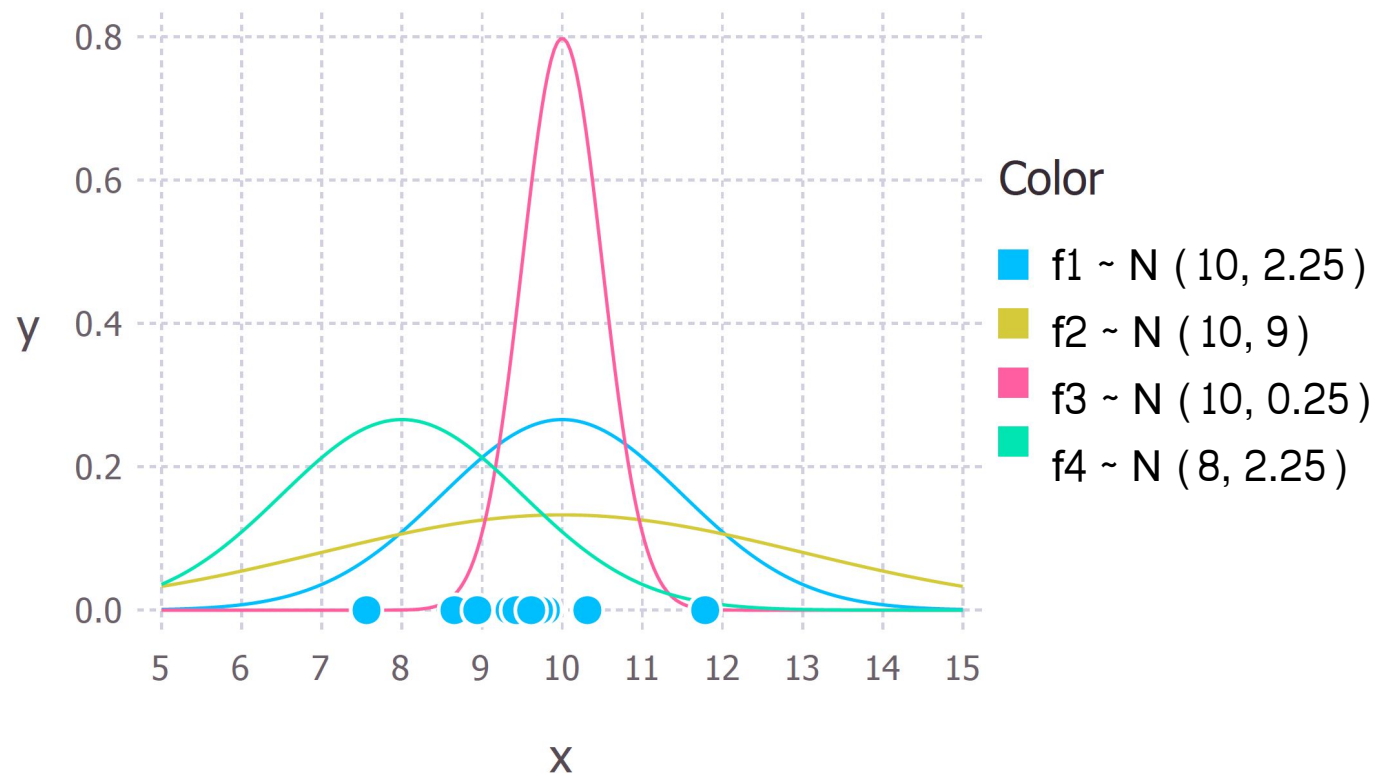
Note:  $p$  replaces  $\gamma$ ,  
and max replaces min

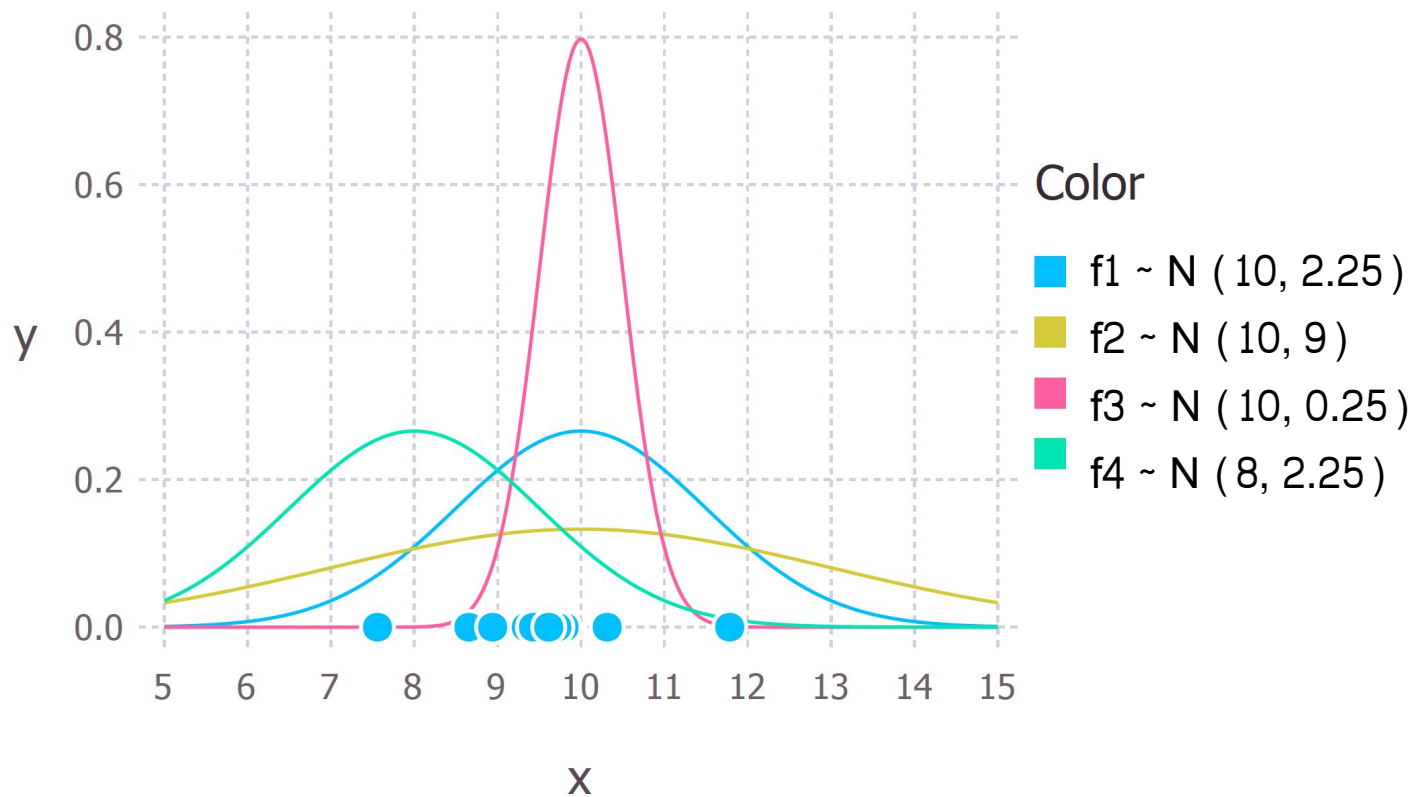


$\log(f(x))$  is monotonic/increasing, same  $\operatorname{argmax}$  as  $f(x)$

# ML: Another example

- Observe a dataset of points  $D = \{x^i\}_{i=1:10}$
- Assume  $x$  is generated by Normal distribution,  $x \sim N(x|\mu, \sigma)$
- Find parameters  $\theta_{ML} = [\mu, \sigma]$  that maximize  $\prod_{i=1}^{10} N(x^i|\mu, \sigma)$





# What is the right distribution that fits the given datapoints?

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.



## What is the right distribution that fits?

Multiple Choice Poll   79 votes   79 participants

Blue  $N(10, 2.25)$  - 56 votes



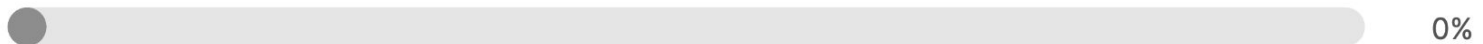
Yellow-green:  $N(10, 9)$  - 12 votes



Magenta:  $N(10, 0.25)$  - 11 votes



Green:  $N(8, 2.25)$  - 0 votes



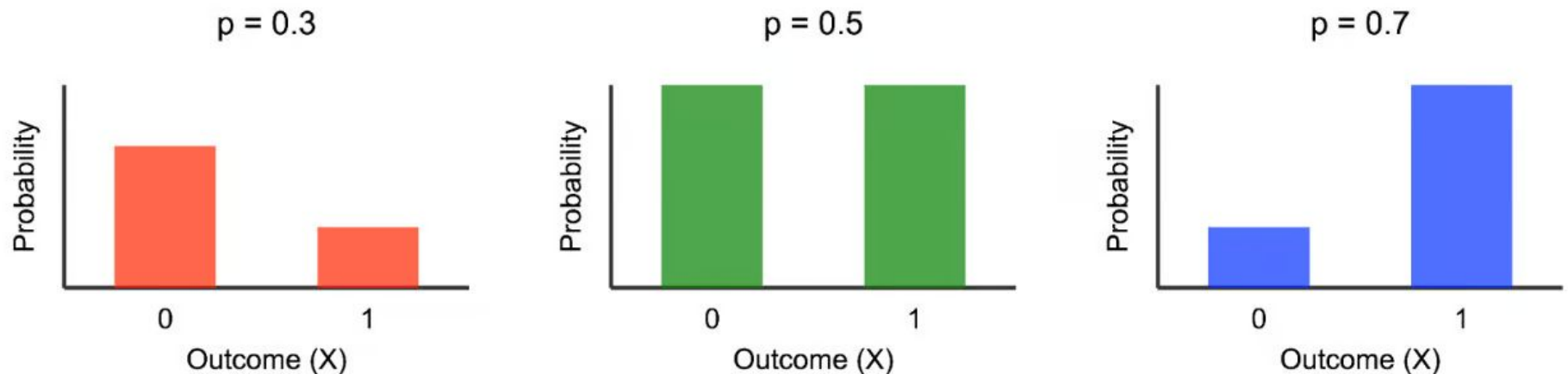


# Types of data distribution: Bernoulli

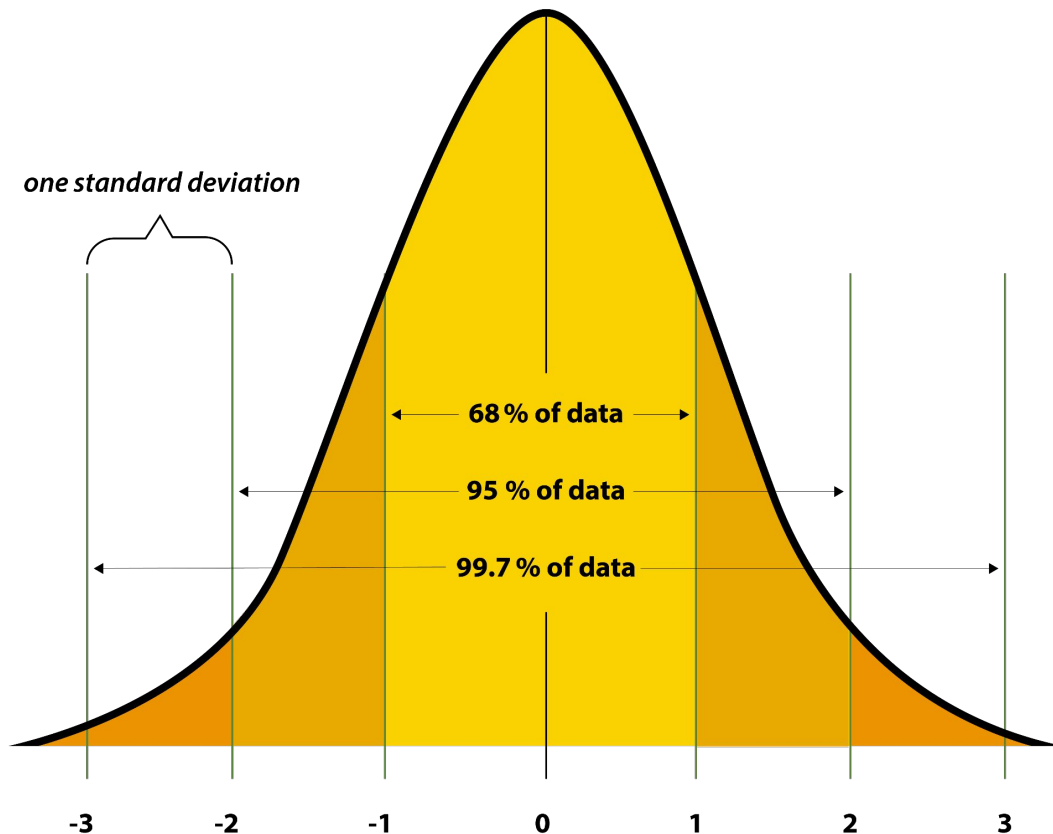
$$P(X = 1) = p, \quad P(X = 0) = 1 - p$$

## Comparison of Bernoulli Distributions

Probability Mass Functions for different  $p$  values



# Types of data distribution: Normal



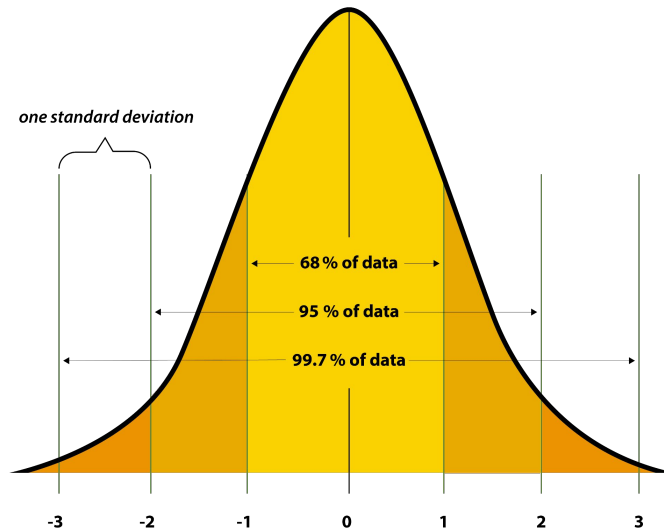
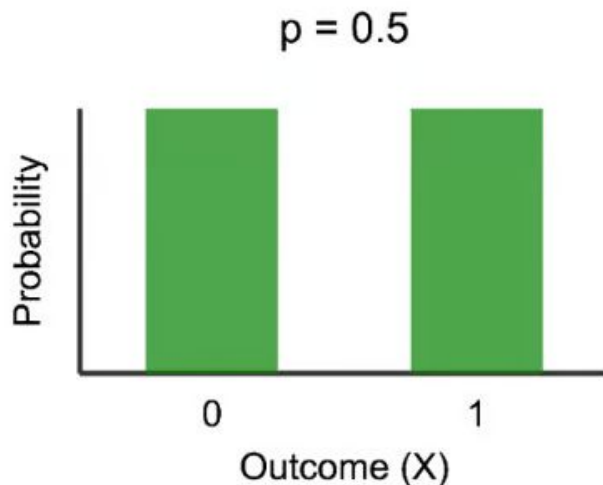
$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} .$$

## Parameters:

- Mean ( $\mu$ )
- Deviation ( $\sigma$ )

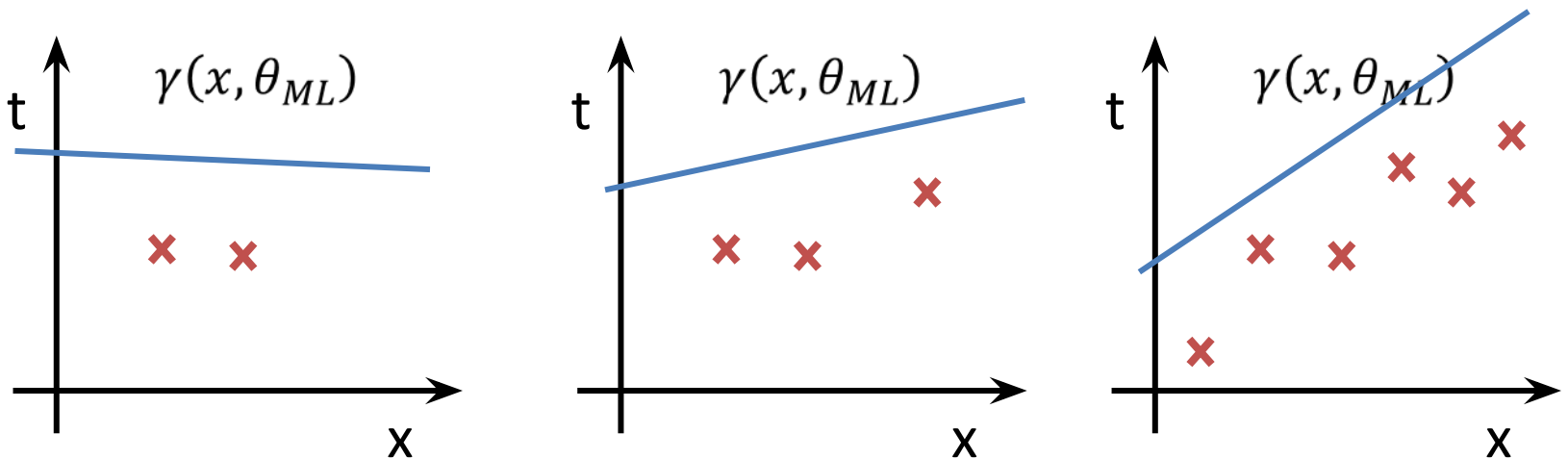
# How to choose a particular distribution?

Distribution	
Bernoulli	<ul style="list-style-type: none"><li>● Discrete data samples, with binary outcomes.<ul style="list-style-type: none"><li>○ Eg: Coin tosses</li></ul></li></ul>
Normal	<ul style="list-style-type: none"><li>● Data is primarily centered around a value, eg: mean</li><li>● Most other data points fall very close to the mean.</li><li>● Eg: Weights of a mouse</li></ul>



# Problem with Maximum Likelihood: **Bias**

$$\text{MLE} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log p(u^{(i)} | \theta)$$

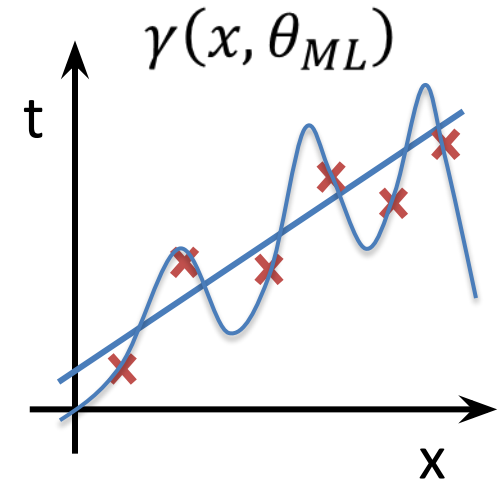


- Depends on the number of datapoints (N)
- When N is small, the estimate may not be an accurate reflection of the true model

# One issue: **Overfitting**

$$\text{MLE} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log p(u^{(i)} | \theta)$$

- Let  $K$  denote the complexity of the model estimator:  $\theta = \{\theta_1, \theta_2, \dots, \theta_K\}$ 
  - $K = 1$  denotes a line.
  - $K = 15$  denotes a 15-degree polynomial function.

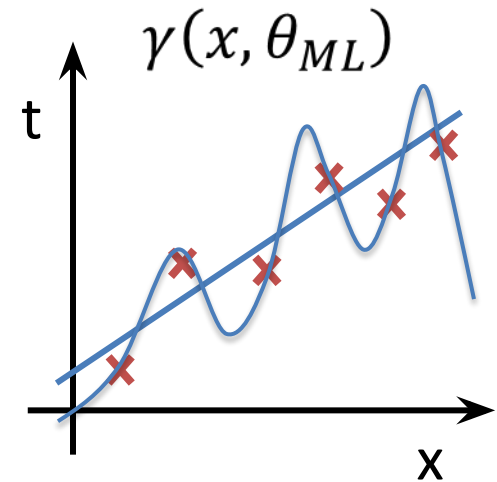


- How to control the value of  $K$ ?
- MLE does not offer a way to directly control model complexity.
  - It will always choose the solution with  $K = 15$ .

# One issue: **Overfitting**

$$\text{MLE} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log p(u^{(i)} | \theta)$$

- Let  $K$  denote the complexity of the model estimator:  $\theta = \{\theta_1, \theta_2, \dots, \theta_K\}$ 
  - $K = 1$  denotes a line.
  - $K = 15$  denotes a 15-degree polynomial function.



- How to control the value of  $K$ ?
- **Solution:** use a **Bayesian method**--define a prior distribution over the parameters (results in regularization)

# Bayesian Modeling

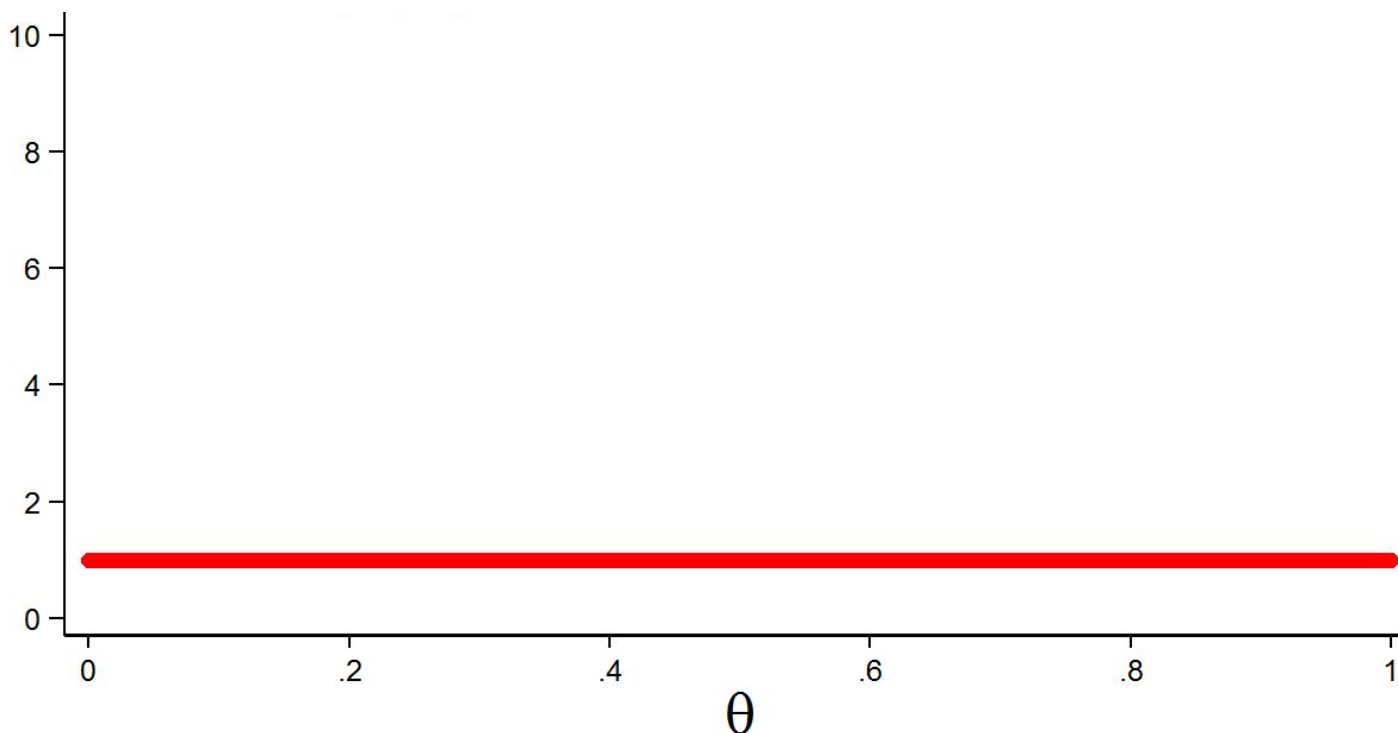


# Prior distribution

- Prior distributions  $p(\theta)$  are probability distributions of model parameters based on some **a priori knowledge** about the problem at hand.
- Prior distributions are assumed **before** any data is observed.
- **Examples:**
  - A coin toss is random (50% heads, 50% tails).
  - Assumed prior : **Binomial distribution**
    - bernoulli is a special case of binomial.

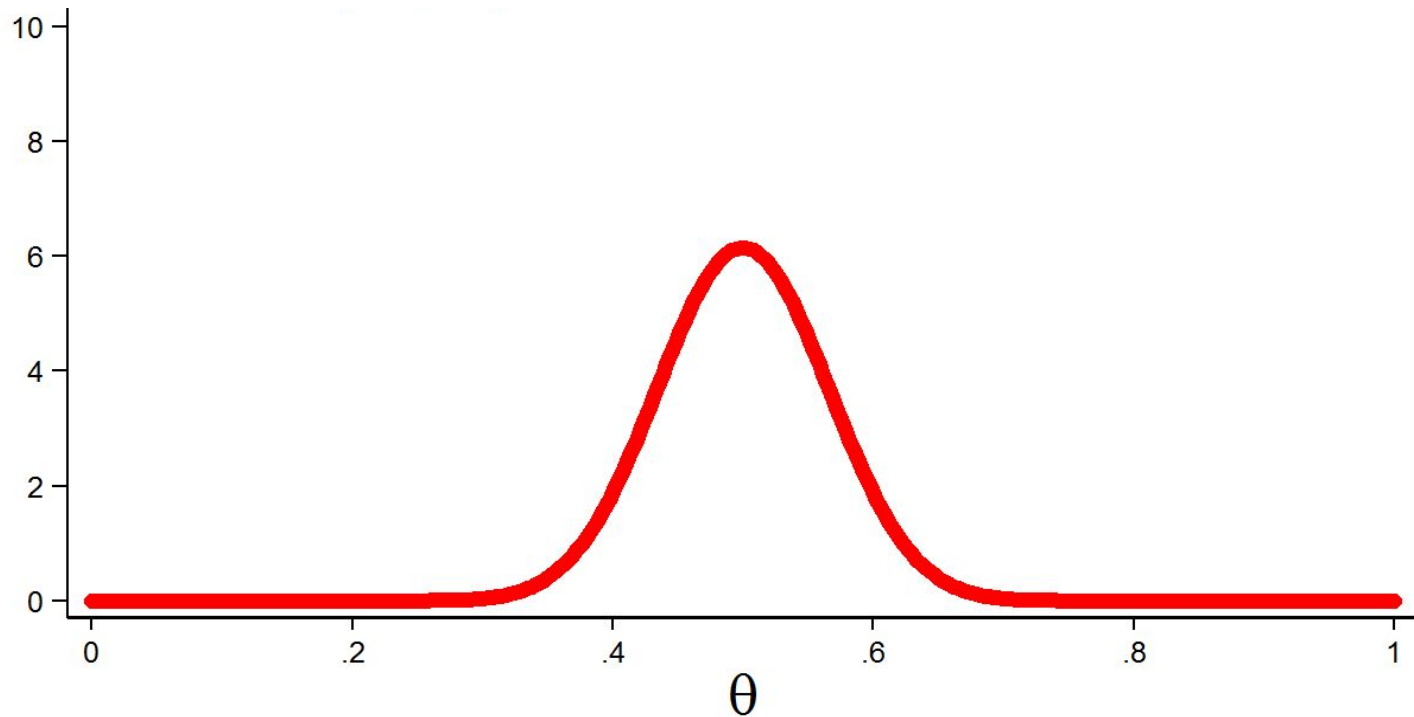


# What information is this prior ( $\theta$ ) is giving?



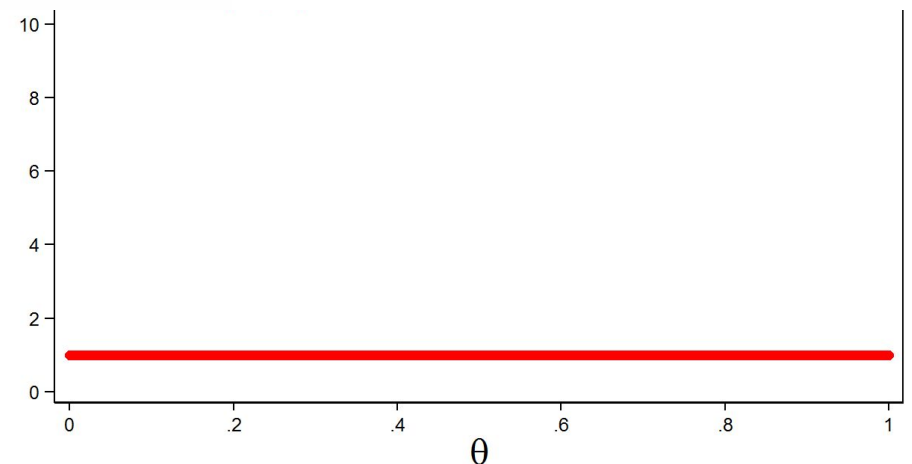
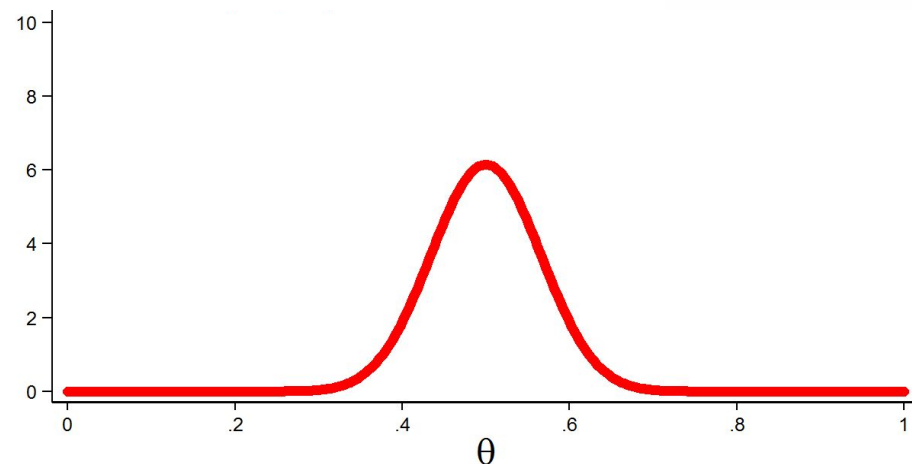
The prior takes a uniform probability distribution.

# What information is this prior ( $\theta$ ) is giving?



- The coin is likely to have a outcome of heads with 50% probability

# Informative v/s uninformative priors



- Important to choose a prior that is more likely to be aligned with the expected outcomes.

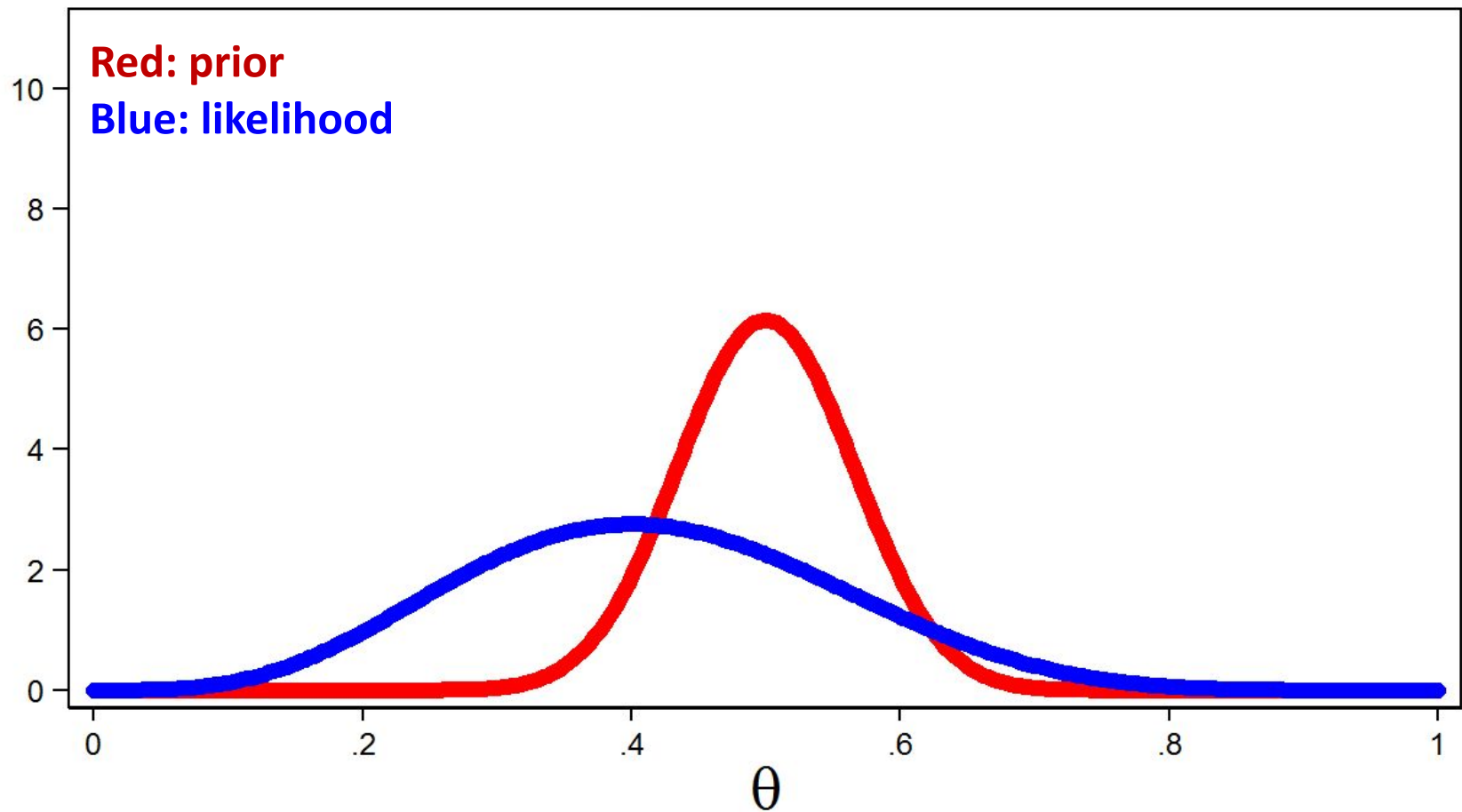
# Reminder: Coin Toss Experiment

- $n=10$  coin tosses
- $y$  (target) = 4 number of heads
- $P(\text{head}) = 4/10 = 0.4$
- $P(\text{tail}) = 1 - 0.4 = 0.6$



**Likelihood**  
**(estimated purely from data)**

# Prior and Likelihood



# Posterior Distribution

$$\textit{Posterior} = \textit{Prior} \times \textit{Likelihood}$$

Your mental model of  
the target distribution

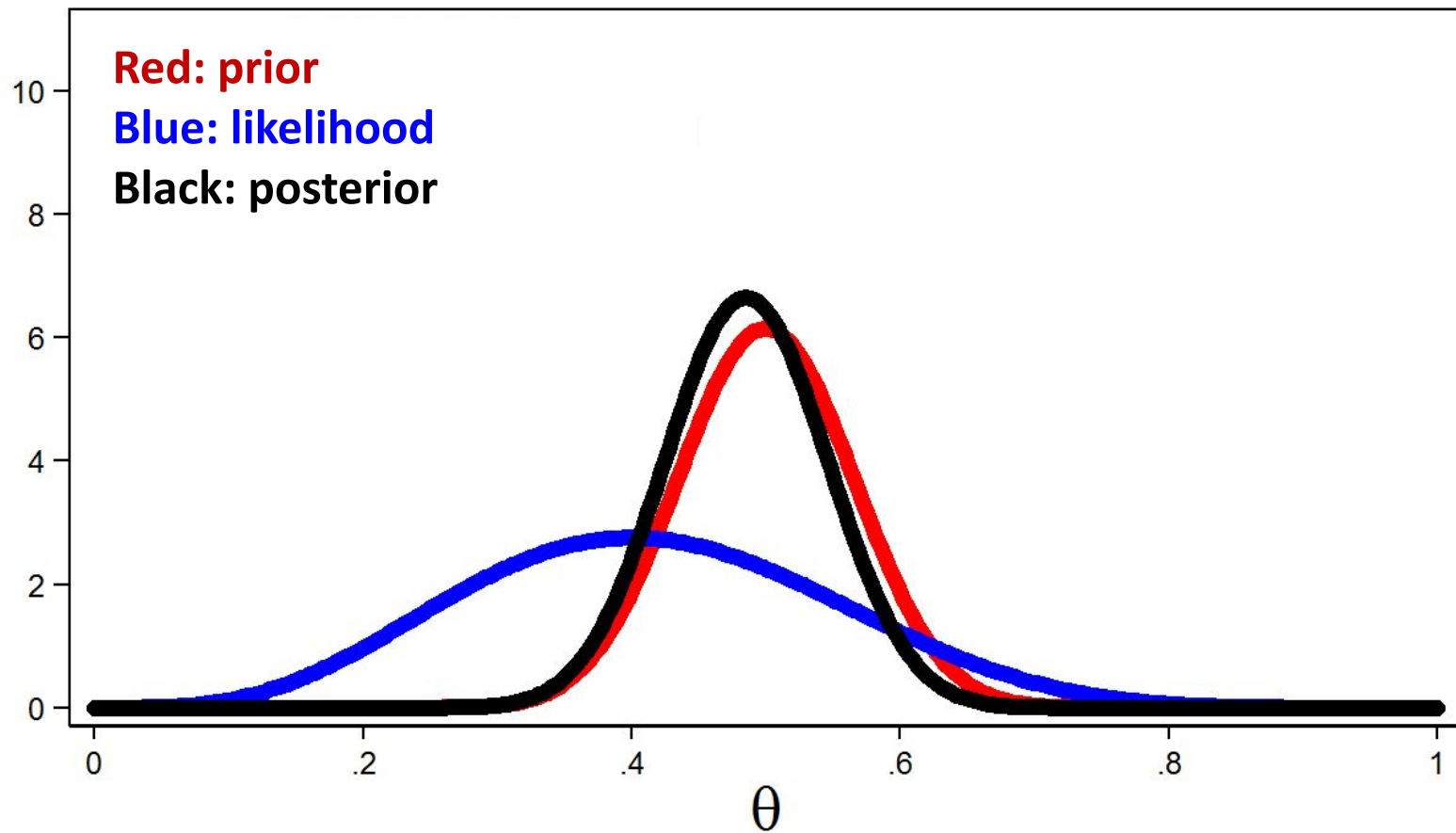


Your estimate from  
the observed data



$$P(\theta|y) = P(\theta)P(y|\theta)$$

# Posterior Distribution



# Evaluating your model

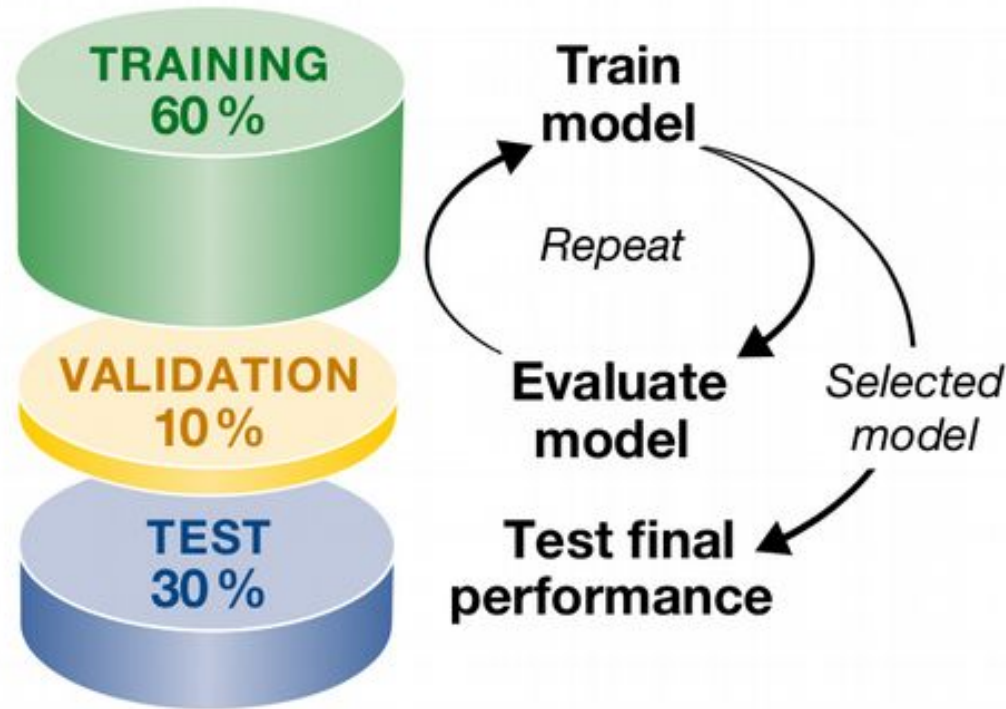


# Setting up an experiment

	Size	Price
train	2104	400
	1600	330
	2400	369
	1416	232
	3000	540
validation	1985	300
	1534	315
	1427	199
test	1380	212
	1494	243

1. For each value of a **hyperparameter**
  - train on the train set
  - evaluate learned parameters on the validation set.
2. Pick the model (hyper parameters) that achieved the **lowest validation error**.
3. Report this model's test set error.

# Train-test-val framework



**Pro:** Computationally efficient

**Con:** This split of data may introduce a bias

slido



## How to avoid bias due to the way training data was split?

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.



## How to avoid bias due to the way training data was split?

Quiz question   81 answers   81 participants

Do not split the data, use all of it! - 2 answers



Split the data randomly many times find hyper parameter that works well most times - 59 answers



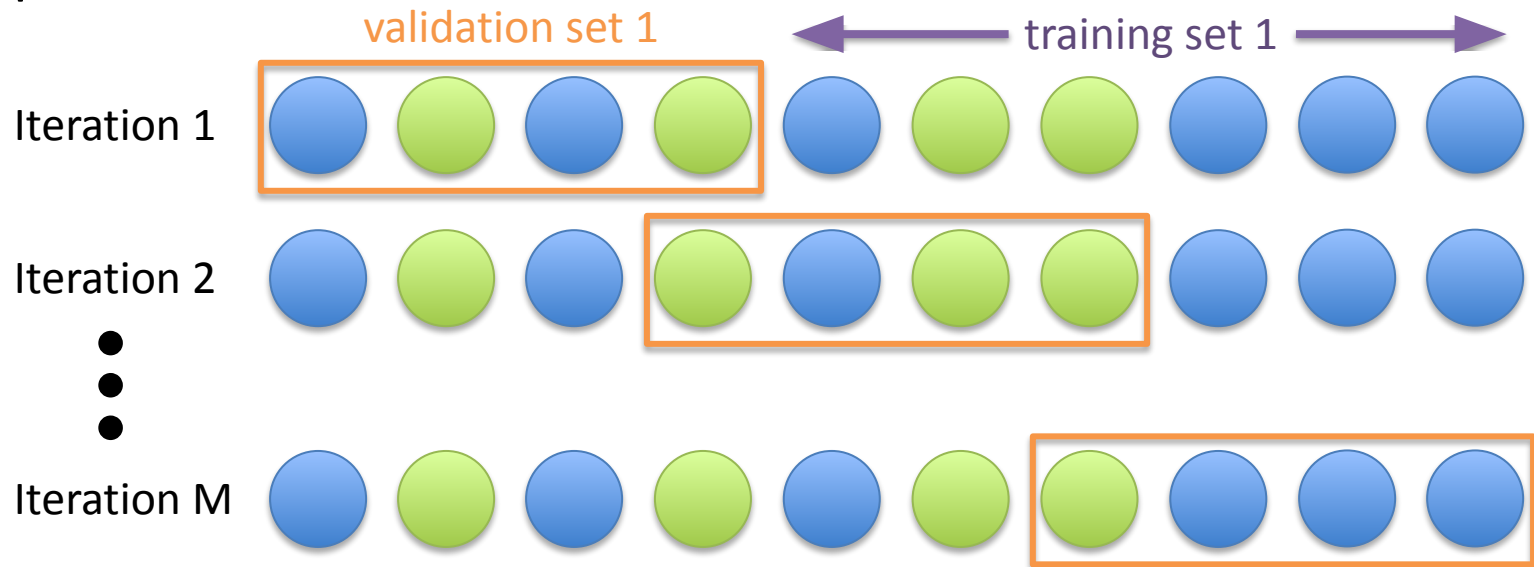
Split the data randomly many times and report the average test metric - 69 answers



slido

# What is the right way to split the data?

- **Solution:** use M-fold cross validation.
- Split training set into train/validation sets M times
- **Option-1:** Report average predictions
- **Option-2:** Pick hyper-parameters that work well across most of M splits



**Pro:** More robust to split bias

**Con:** Requires training M different models

slido



What do you call the scenario when a classifier is performing poorly on training data?

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.

**What do you call the scenario when a classifier is performing poorly on training data?**

Underfitting ✓



Overfitting



Both



None



slido



What do we call a classifier performing well on training data but poorly on test data?

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.



**What do we call a classifier performing well on training data but poorly on test data?**

---

Overfitting



Underfitting



Both



None



slido



# What do classifiers work best on?

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.

## What do classifiers work best on?

---

Training data



All of the above



Validation data



Test data



# Reasons for underfitting

- Not enough training samples.
- Model too simple.
- Training samples too diverse.

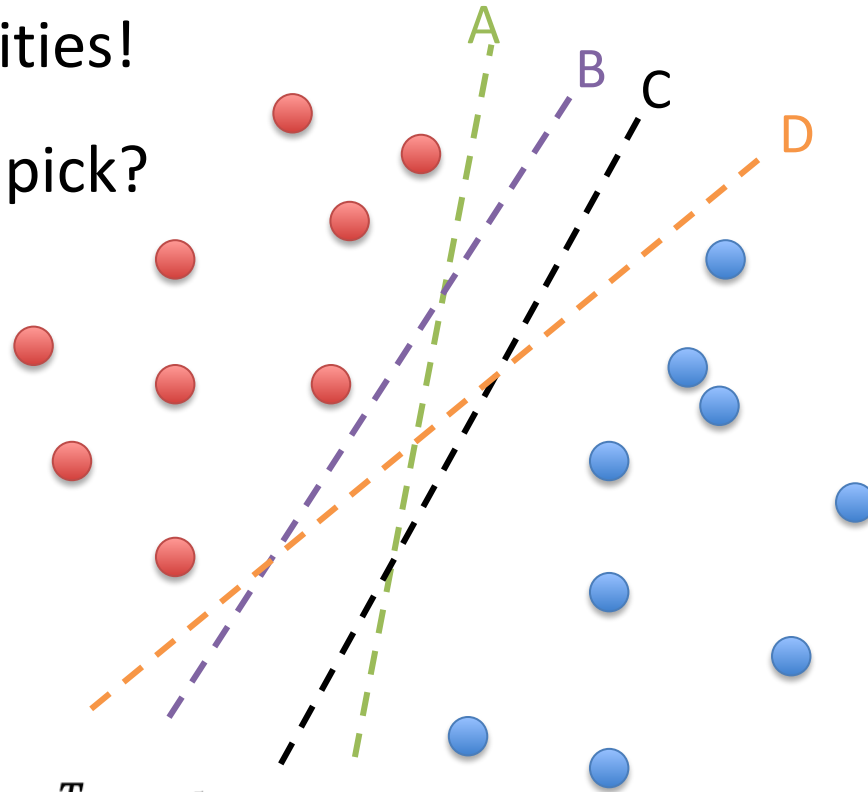
# Reasons for overfitting

- Model overparameterized.
- Not correctly cross-validated.
- Training data from a very different distribution than test data.

# **Learning the right decision boundary**

# Separating two classes

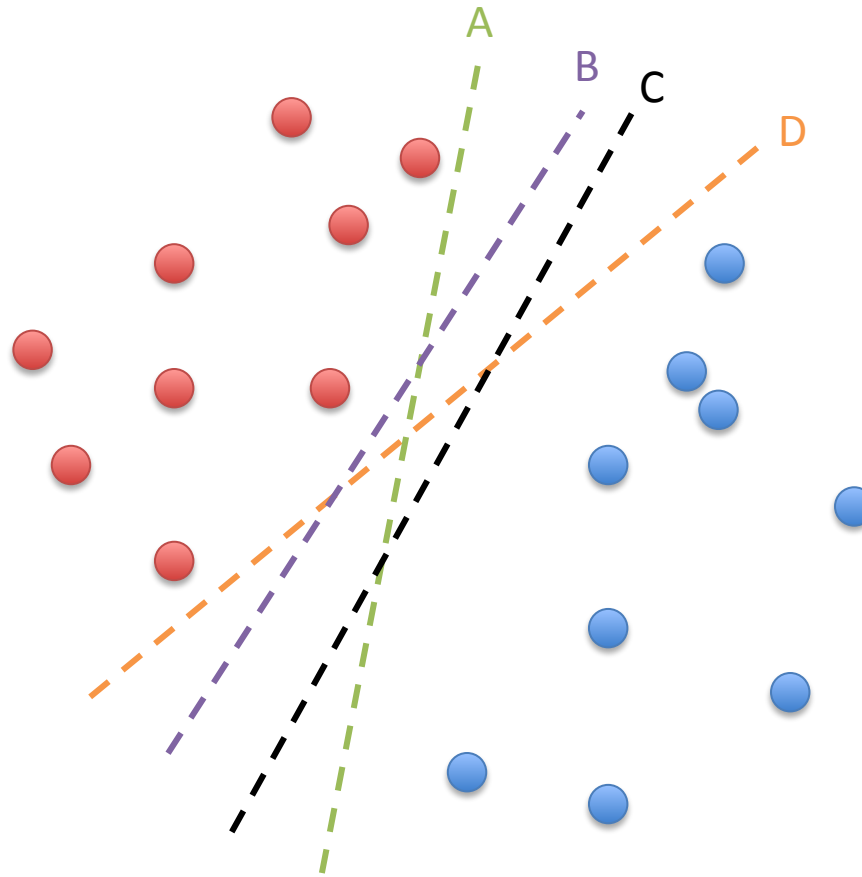
- Many possibilities!
- Which one to pick?



Decision boundary:  $w^T x + b = 0$

$$y = \begin{cases} +1 \text{ [red]} & \text{if } \text{sign}(a^T x + b) \geq 0 \\ -1 \text{ [blue]} & \text{if } \text{sign}(a^T x + b) < 0 \end{cases}$$

# Separating two classes



**Goal:** Avoid misclassifying new test points generated from the same distribution as the training points



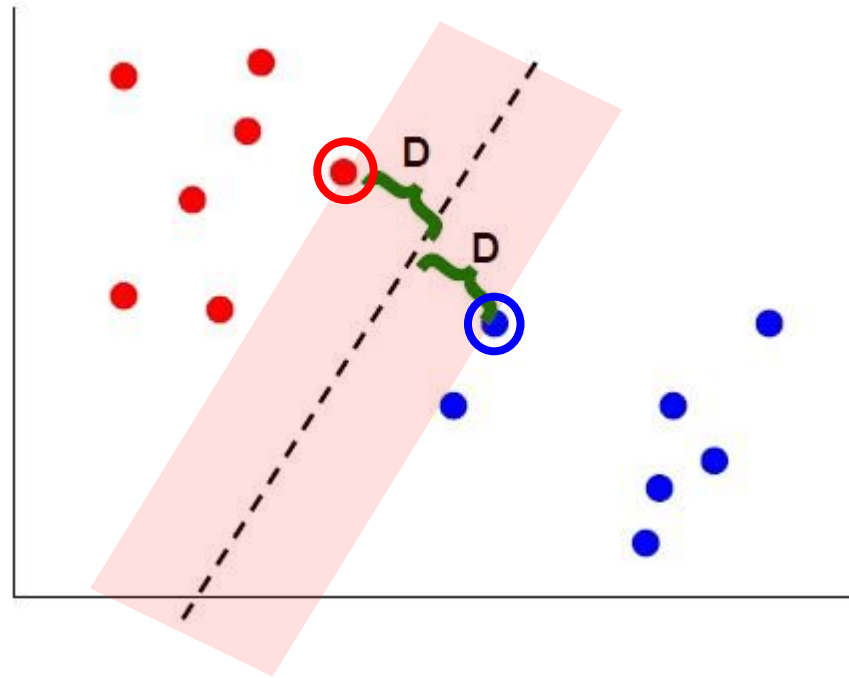
# Maximum margin classification

**Intuition:** Instead of fitting all the points, focus on boundary points

**Aim:** learn a boundary that leads to the largest margin (buffer) from points on both sides

**Why:** make the decision boundary robust to small perturbations

The subset of vectors that support (determine boundary) are called the **support vectors** (circled)



# Next Class

## **Classification III:**

Regularization, stochastic gradient descent,  
multiclass SVM

**Reading:** Forsyth Ch 2.1.3-2.1.7