



# CS 540 Introduction to Artificial Intelligence

## **Classification - KNN and Naive Bayes**

University of Wisconsin—Madison  
Fall 2025, Section 3  
October 1, 2025



# Announcements

- HW3 due Friday 10/3 at 11:59 PM

- Class roadmap:

ML: Unsupervised Learning

ML Linear Regression

**Machine Learning: K - Nearest Neighbors  
& Naive Bayes**

Machine Learning: Neural Networks I  
(Perceptron)

Machine Learning: Neural Networks II

Supervised Learning

# Outline

- K-Nearest Neighbors
- Maximum likelihood estimation
- Naive Bayes





# Part I: K-nearest neighbors





WIKIPEDIA  
The Free Encyclopedia

[Main page](#)

Article

[Talk](#)

# *k*-nearest neighbors algorithm

From Wikipedia, the free encyclopedia

*Not to be confused with [k-means clustering](#).*

(source: wiki)

# Example 1: Predict if a user likes a song or not



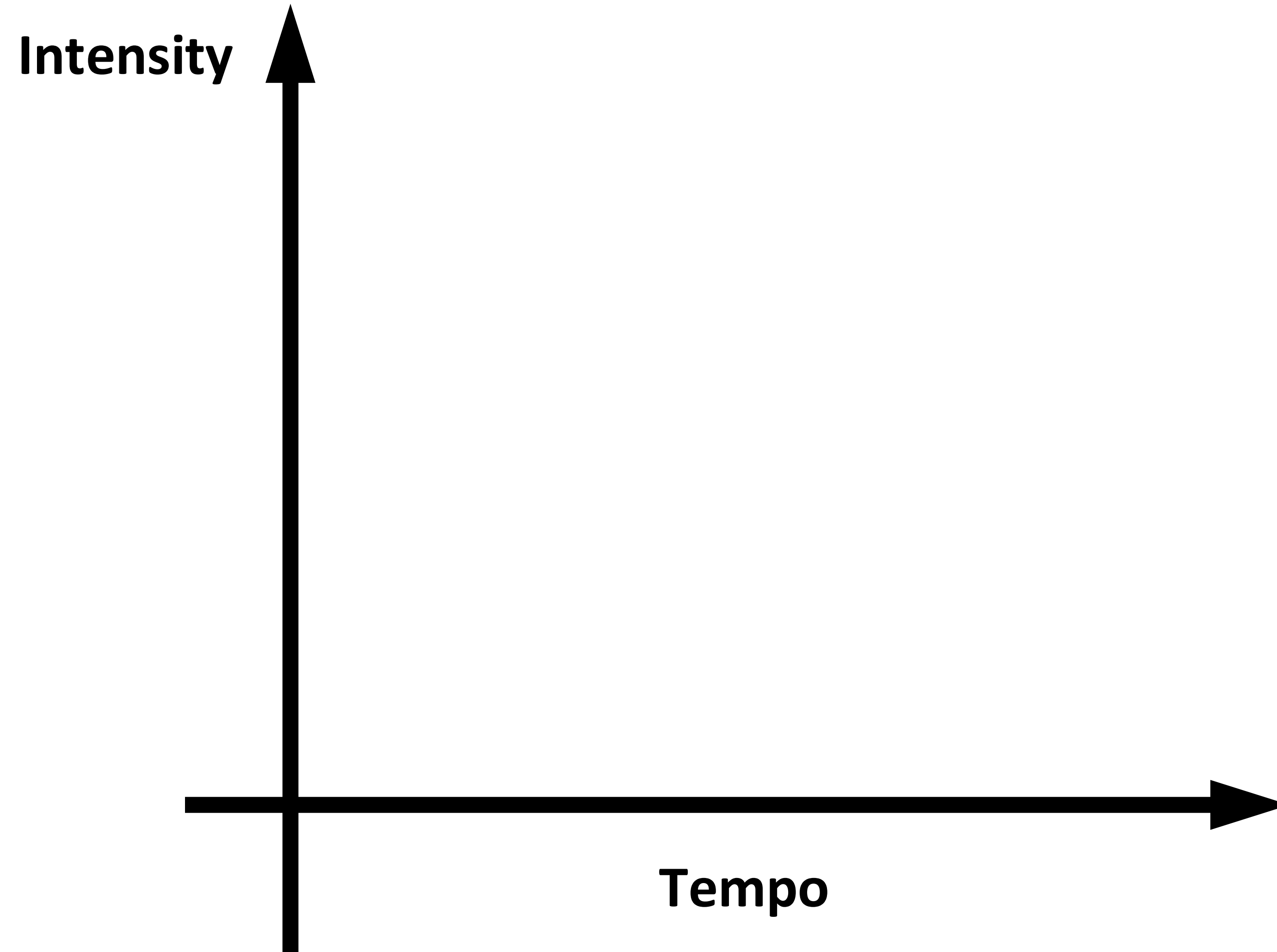
model



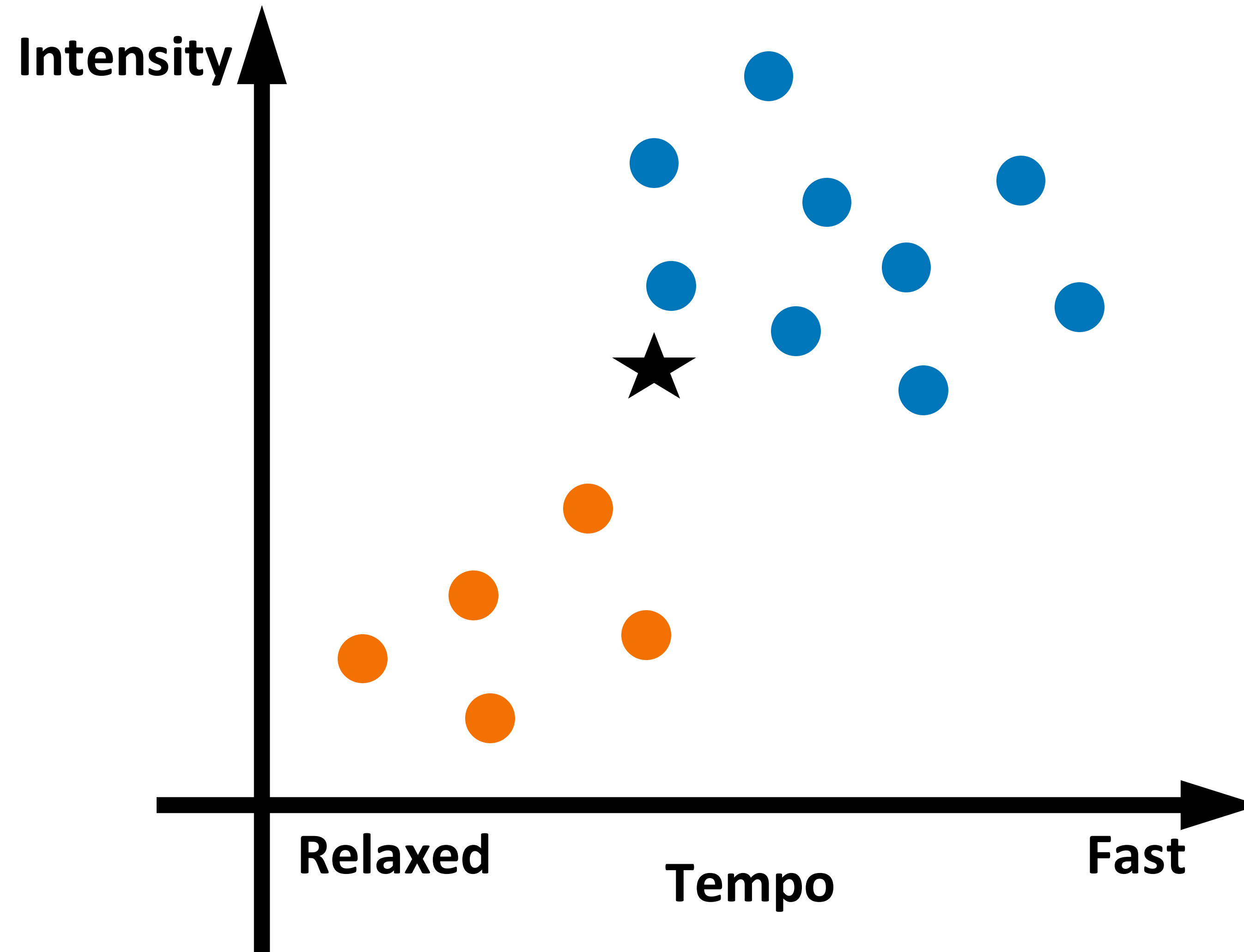
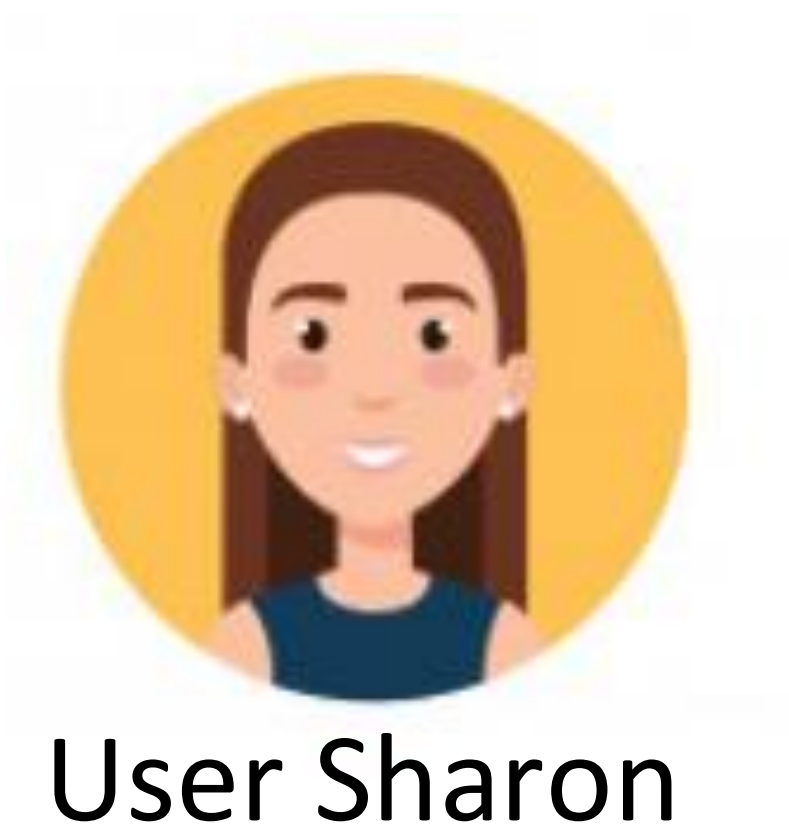
# Example 1: Predict if a user likes a song or not



User Sharon



# Example 1: Predict if a user likes a song or not 1-NN





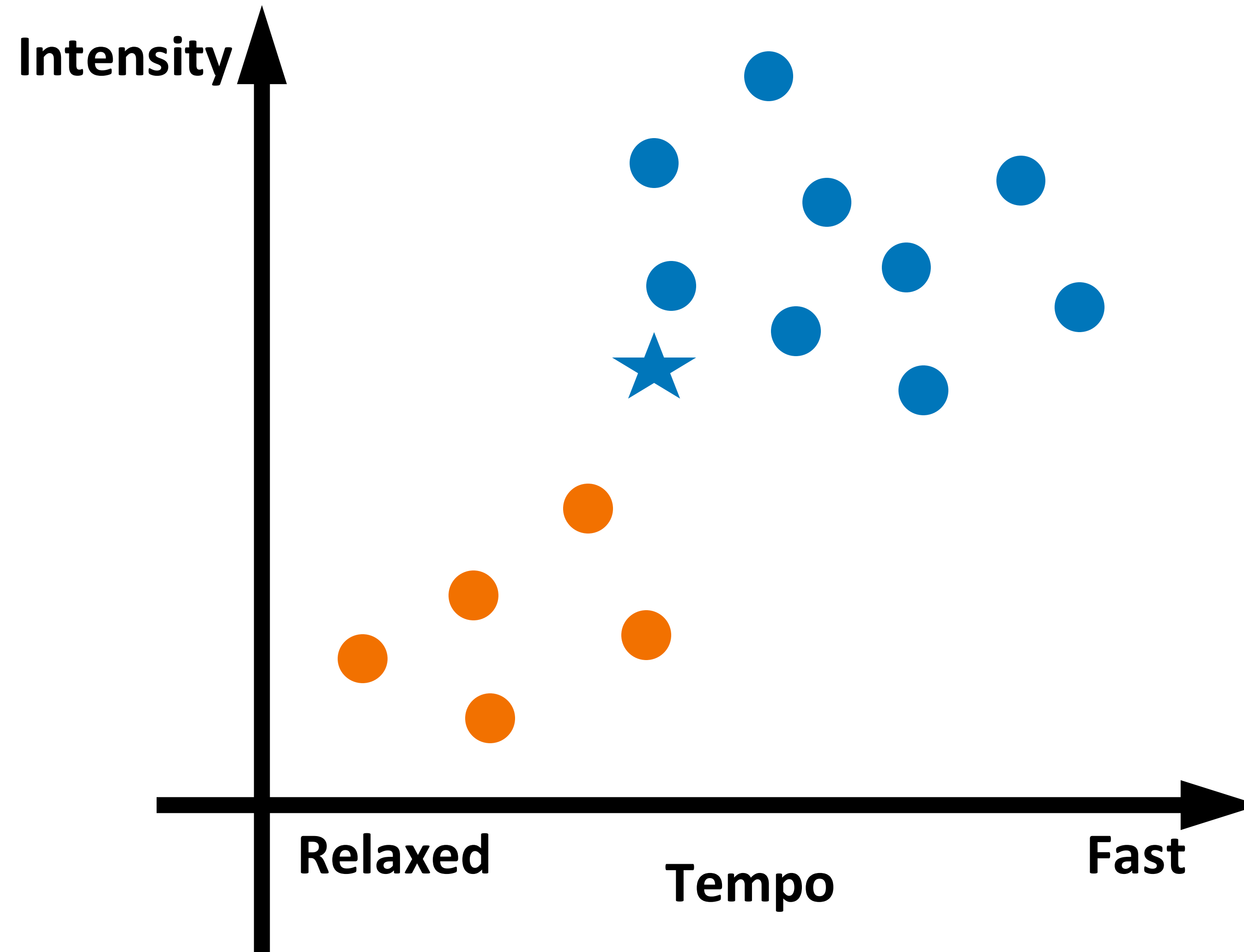
# Example 1: Predict if a user likes a song or not 1-NN



User Sharon

● Dislike

● Like



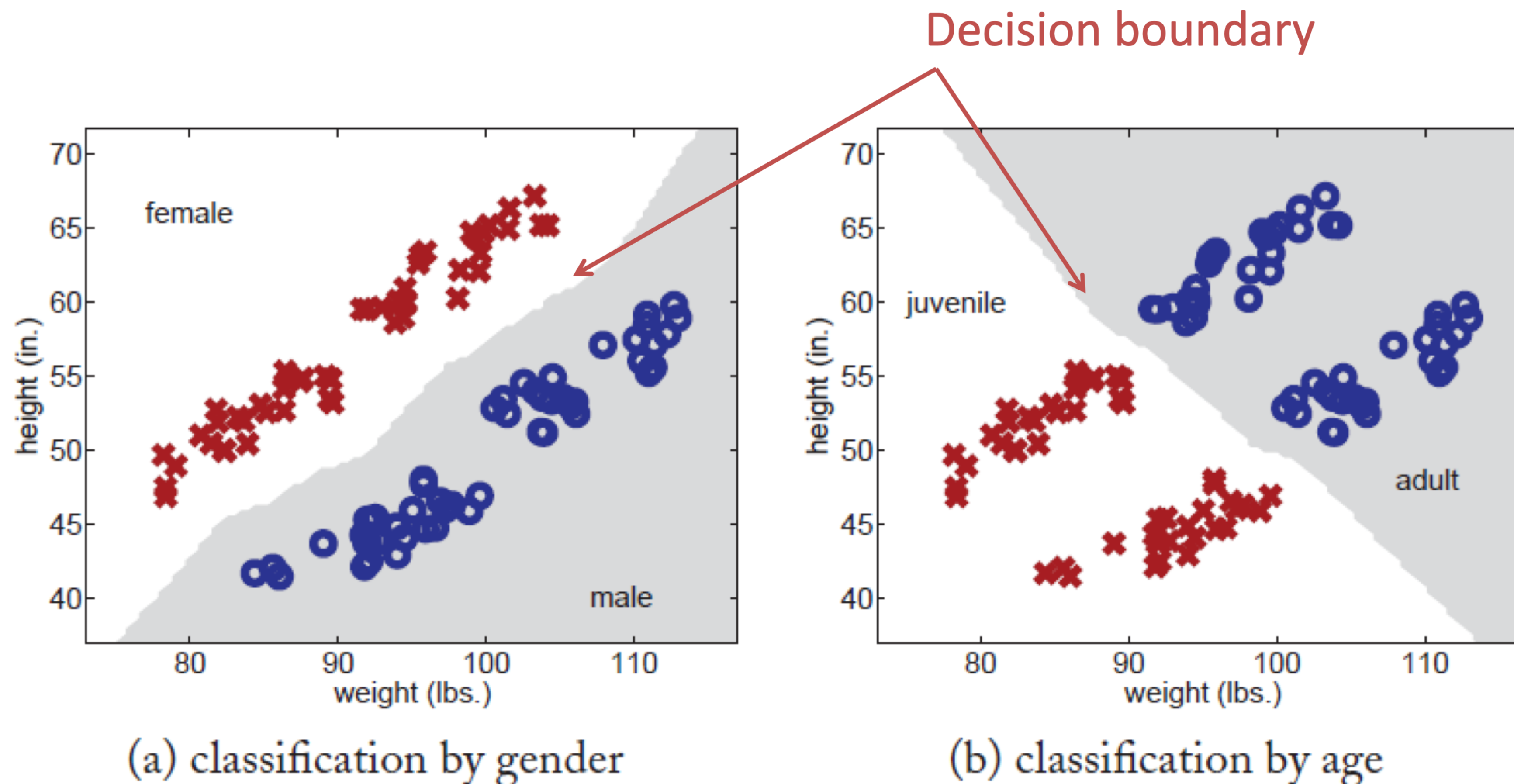
# K-nearest neighbors for classification

- Input: **Training data**  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$   
**Distance function**  $d(\mathbf{x}_i, \mathbf{x}_j)$ ; **number of neighbors**  $k$ ; **test data**  $\mathbf{x}^*$ 
  1. Find the  $k$  training instances  $\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}$  closest to  $\mathbf{x}^*$  under  $d(\mathbf{x}_i, \mathbf{x}_j)$
  2. Output  $y^*$ , the majority class of  $y_{i_1}, \dots, y_{i_k}$ . Break ties randomly.



# Example 2: 1-NN for little green man

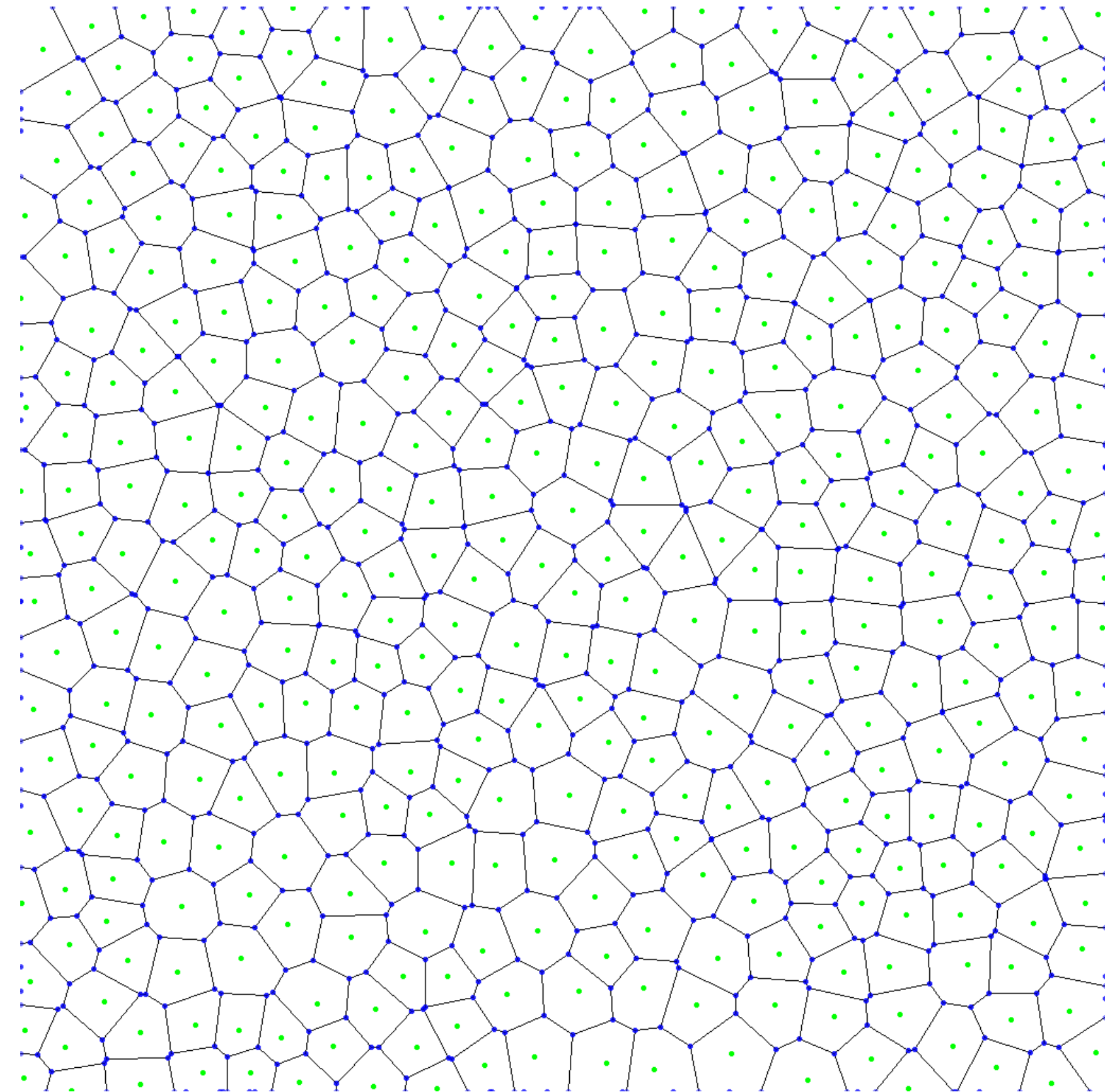
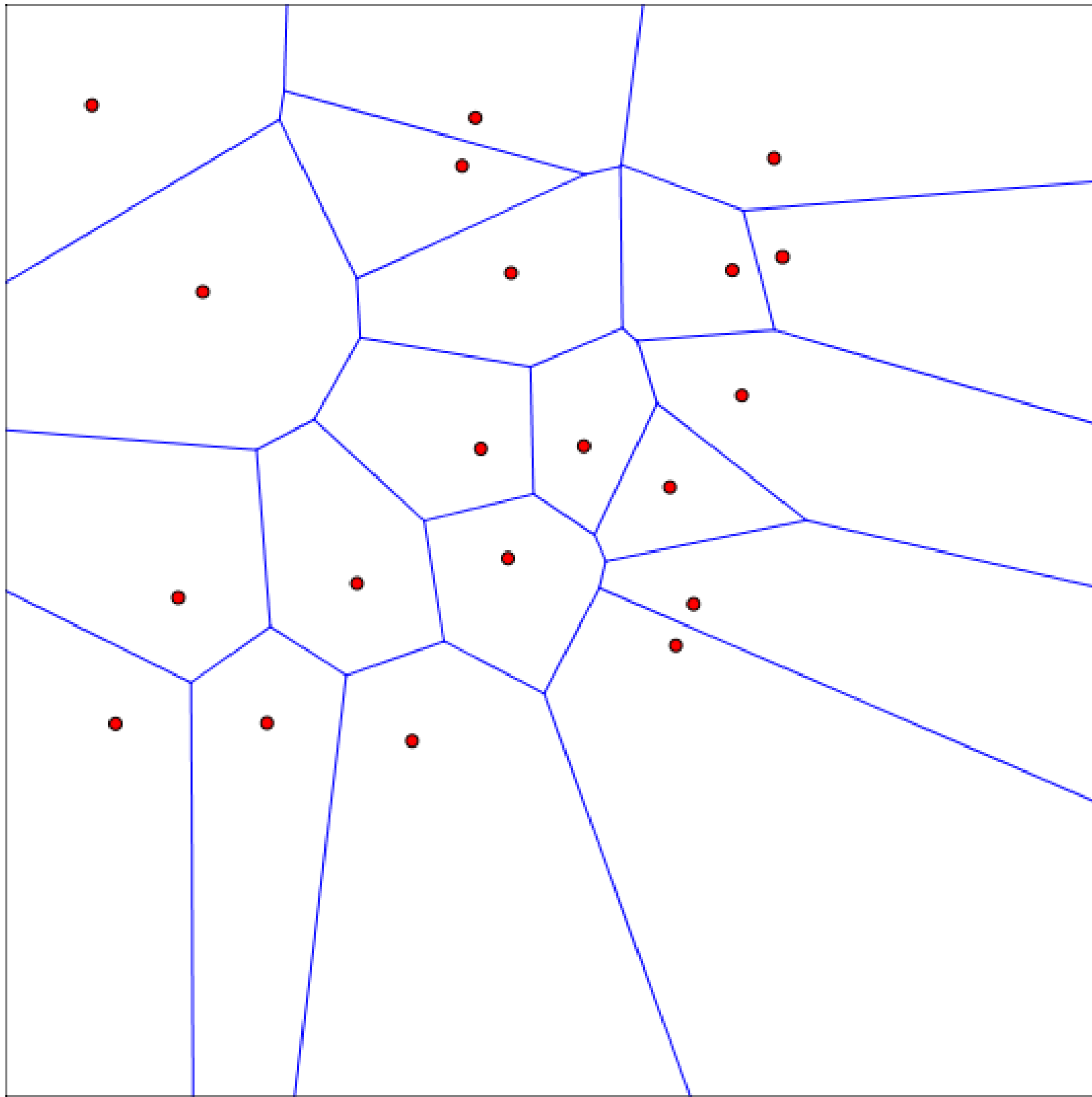
- Predict gender (M,F) from weight, height
- Predict age (adult, juvenile) from weight, height



# 1NN: Decision Regions

Defined by “**Voronoi Diagram**”

- Each cell contains points closer to a particular training point





# k-Nearest Neighbors: Distances

**Discrete features:** Hamming distance

$$d_H(x^{(i)}, x^{(j)}) = \sum_{a=1}^d 1\{x_a^{(i)} \neq x_a^{(j)}\}$$

**Continuous features:**

- Euclidean distance:

$$d(x^{(i)}, x^{(j)}) = \left( \sum_{a=1}^d (x_a^{(i)} - x_a^{(j)})^2 \right)^{\frac{1}{2}}$$

- L1 (Manhattan) dist.:

$$d(x^{(i)}, x^{(j)}) = \sum_{a=1}^d |x_a^{(i)} - x_a^{(j)}|$$

# k-Nearest Neighbors: Regression

**Training/learning:** given

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

**Prediction:** for  $x$ , find  $k$  most similar training points

Return

$$\hat{y} = \frac{1}{k} \sum_{i=1}^k y^{(i)}$$

- I.e., among the  $k$  points, output mean label.



# More on distance functions...

- Be careful with **scale**
- Same feature but different units may change relative distance (fixing other features)
- Sometimes OK to normalize each feature dimension

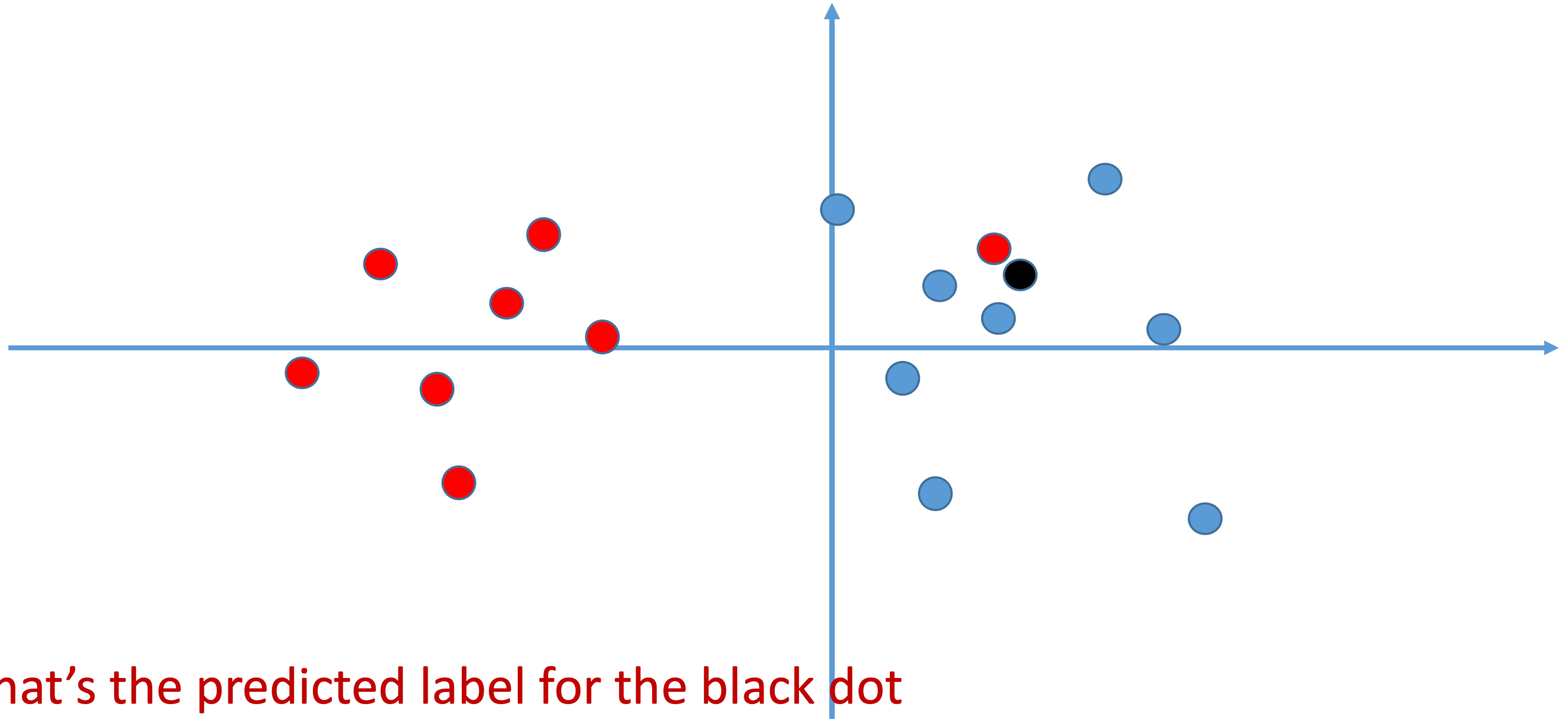
$$x'_{id} = \frac{x_{id} - \mu_d}{\sigma_d}, \forall i = 1 \dots n, \forall d$$

Training set mean for dimension d

Training set standard deviation for dimension d

- Other times not OK: e.g. dimension contains small random noise

# Effect of $k$



What's the predicted label for the black dot using 1 neighbor? 3 neighbors?



# How to pick k, the number of neighbors

- Split data into training and **tuning/validation sets**
- Classify tuning set with different k
- Pick k that produces least tuning-set error

(Shuffle whole dataset first)







## Part II: Maximum Likelihood Estimation



# Supervised Machine Learning

Non-parametric  
(e.g., KNN)

vs.

Parametric

# Supervised Machine Learning

Statistical modeling approach

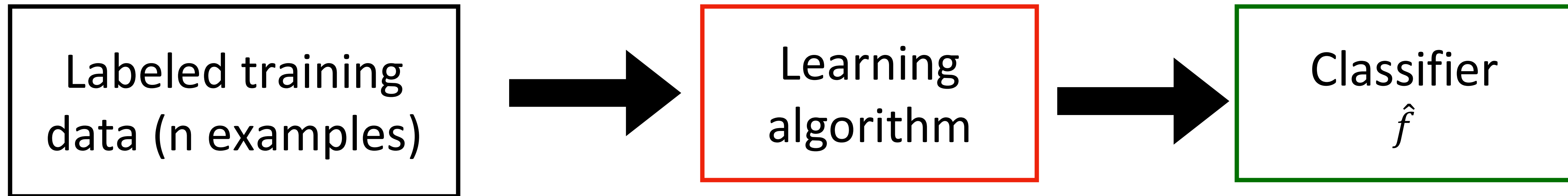
Labeled training  
data (n examples)

$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$

drawn **independently** from  
a fixed distribution  
(also called the i.i.d. assumption)

# Supervised Machine Learning

Statistical modeling approach



$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$

drawn **independently** from  
a fixed underlying distribution  
(also called the i.i.d. assumption)

select  $\hat{f}(\theta)$  from a pool of models  $\mathcal{F}$   
that **best describe the data observed**

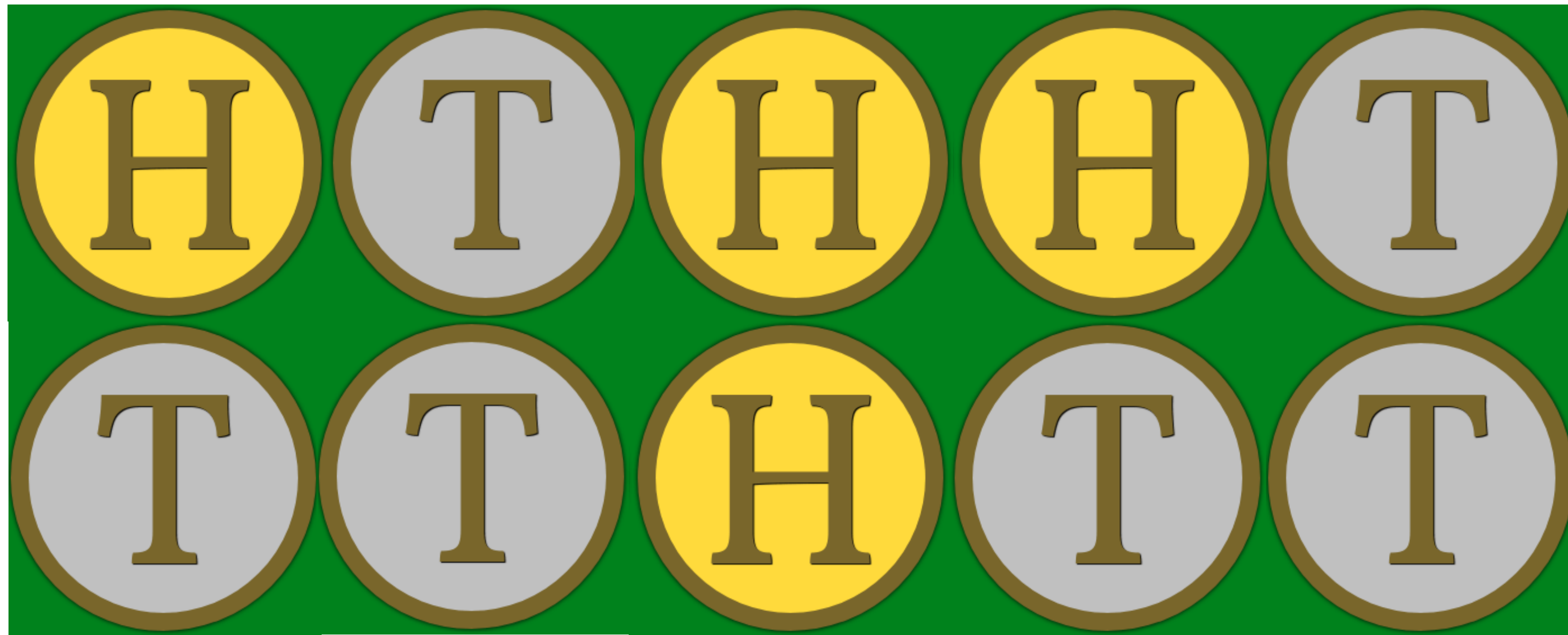


How to select  $\hat{f} \in \mathcal{F}$ ?

- **Maximum likelihood (best fits the data)**
- Maximum a posteriori  
(best fits the data but incorporates prior assumptions)
- Optimization of 'loss' criterion (best discriminates the labels)

# Maximum Likelihood Estimation: An Example

Flip a coin 10 times, how can you estimate  $\theta = p(\text{Head})$ ?



Intuitively,  $\theta = 4/10 = 0.4$


# How good is $\theta$ ?

It depends on how likely it is to generate the observed data

$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$

(Let's forget about label for a second)

**Likelihood function**

$$L(\theta) = \prod_i p(\mathbf{x}_i | \theta)$$


Under i.i.d assumption

Interpretation: How **probable** (or how likely) is the data given the probabilistic model  $p_\theta$ ?



# How good is $\theta$ ?

It depends on how likely it is to generate the observed data

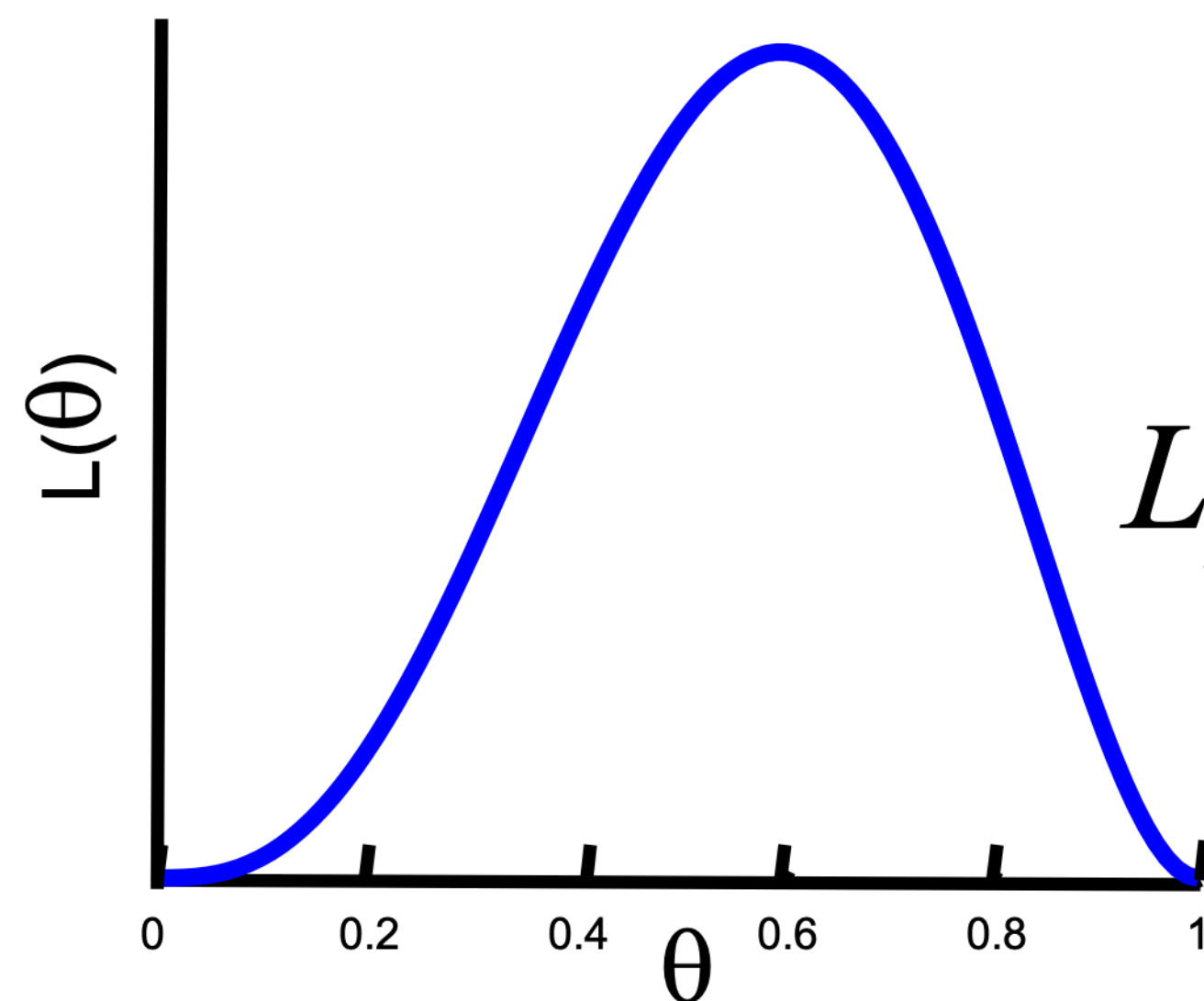
$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$

(Let's forget about label for a second)

Likelihood function

$$L(\theta) = \prod_i p(\mathbf{x}_i | \theta)$$

H, T, T, H, H



$$L_D(\theta) = \theta \cdot (1 - \theta) \cdot (1 - \theta) \cdot \theta \cdot \theta$$

Bernoulli distribution

## Log-likelihood function

$$\begin{aligned} L_D(\theta) &= \theta \cdot (1 - \theta) \cdot (1 - \theta) \cdot \theta \cdot \theta \\ &= \theta^{N_H} \cdot (1 - \theta)^{N_T} \end{aligned}$$

Log-likelihood function

$$\begin{aligned} \ell(\theta) &= \log L(\theta) \\ &= N_H \log \theta + N_T \log(1 - \theta) \end{aligned}$$

# Maximum Likelihood Estimation (MLE)

Find optimal  $\theta^*$  to maximize the likelihood function (and log-likelihood)

$$\theta^* = \operatorname{argmax} N_H \log \theta + N_T \log(1 - \theta)$$

$$\frac{\partial l(\theta)}{\partial \theta} = \frac{N_H}{\theta} - \frac{N_T}{1 - \theta} = 0 \quad \Rightarrow \quad \theta^* = \frac{N_H}{N_T + N_H}$$

which confirms your intuition!



# Connecting MLE and Loss Minimization

- MLE solves

$$\operatorname{argmax}_{\theta} p(x_1, \dots, x_n \mid \theta) = \operatorname{argmax}_{\theta} \prod_{i=1}^n p(x_i \mid \theta)$$

- Rewrite the problem in an equivalent form

$$\operatorname{argmax}_{\theta} p(x_1, \dots, x_n \mid \theta) = \operatorname{argmin}_{\theta} (-\log p(x_1, \dots, x_n \mid \theta))$$

$$= \operatorname{argmin}_{\theta} \sum_{i=1}^n -\log p(x_i \mid \theta)$$

# Connecting MLE and Loss Minimization

- We call “ $-\log p(x_i \mid \theta)$ ” the **negative log likelihood**
- May define  $\ell(\theta; x_i) := -\log p(x_i \mid \theta)$
- Maximum likelihood estimation is **loss minimization**.  
Different notation, same computation.

$$\operatorname{argmax}_{\theta} p(x_1, \dots, x_n \mid \theta) = \operatorname{argmin}_{\theta} \sum_{i=1}^n \ell(\theta; x_i)$$

Lecture stopped here on 10/1/25



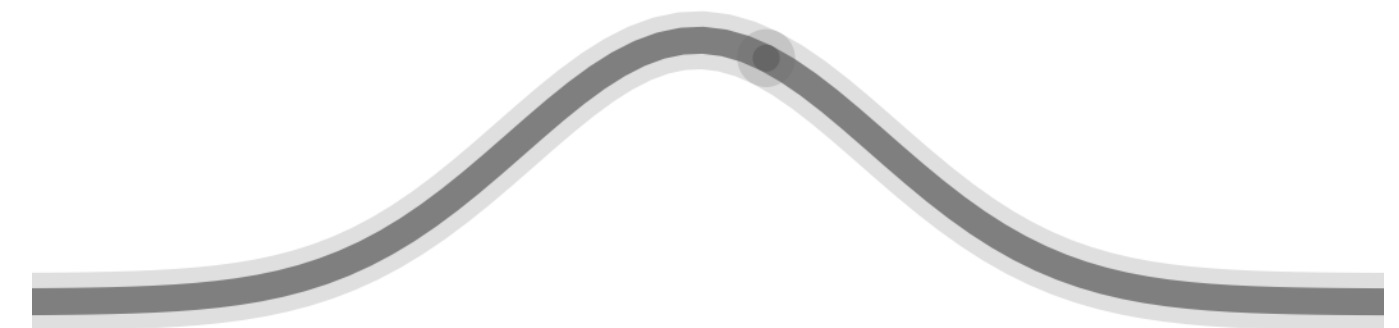
# Maximum Likelihood Estimation: Gaussian Model

Fitting a model to heights of females

**Observed some data** (in inches): 60, 62, 53, 58,...  $\in \mathbb{R}$

$$\{x_1, x_2, \dots, x_n\}$$

**Model class:** Gaussian model



$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

So, what's the MLE for the given data?

# Estimating the parameters in a Gaussian

- **Mean**

$$\mu = \mathbf{E}[x] \quad \text{hence} \quad \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

- **Variance**

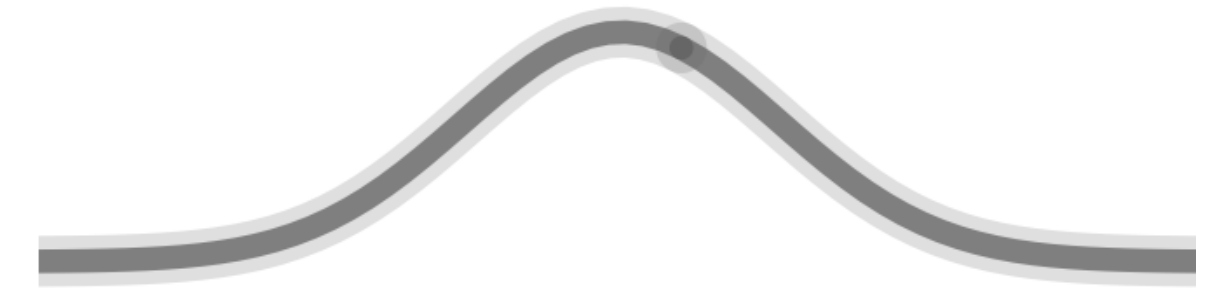
$$\sigma^2 = \mathbf{E}[(x - \mu)^2] \quad \text{hence} \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

Why?

# Maximum Likelihood Estimation: Gaussian Model

**Observe some data** (in inches):  $x_1, x_2, \dots, x_n \in \mathbb{R}$

Assume that the data is drawn from a Gaussian



$$L(\mu, \sigma^2 | X) = \prod_{i=1}^n p(x_i; \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

**Fitting parameters is maximizing likelihood w.r.t  $\mu, \sigma^2$**   
(maximize likelihood that data was generated by model)

**MLE**

$$\arg \max_{\mu, \sigma^2} \prod_{i=1}^n p(x_i; \mu, \sigma^2)$$



# Maximum Likelihood

- Estimate parameters by finding ones that explain the data

$$\operatorname{argmax}_{\mu, \sigma^2} \prod_{i=1}^n p(x_i; \mu, \sigma^2) = \operatorname{argmin}_{\mu, \sigma^2} -\log \prod_{i=1}^n p(x_i; \mu, \sigma^2)$$

- **Decompose likelihood**

$$\sum_{i=1}^n \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (x_i - \mu)^2 = \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Minimized for  $\mu = \frac{1}{n} \sum_{i=1}^n x_i$

# Maximum Likelihood

- Estimating the variance

$$\frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

# Maximum Likelihood

- Estimating the variance

$$\frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

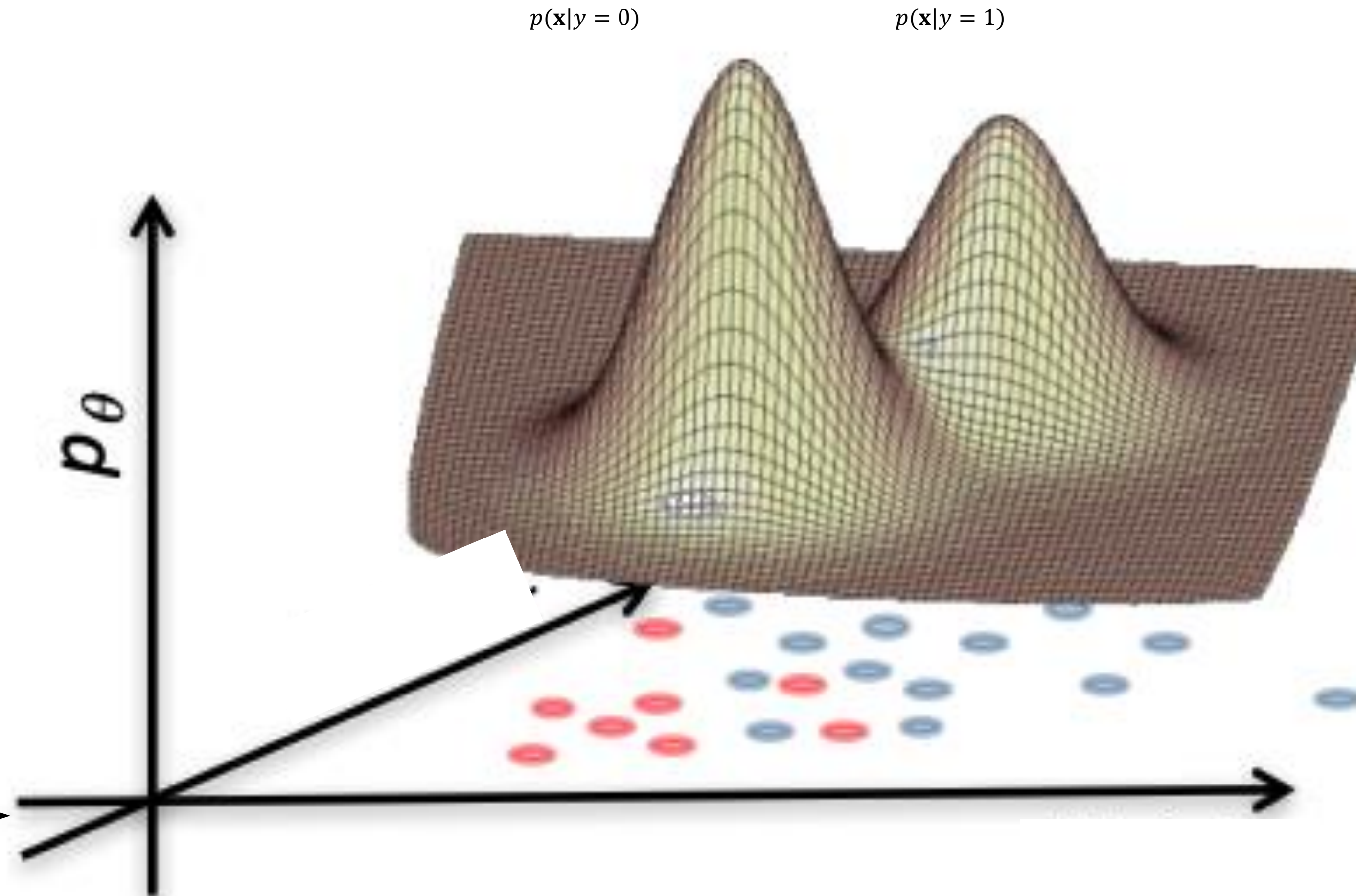
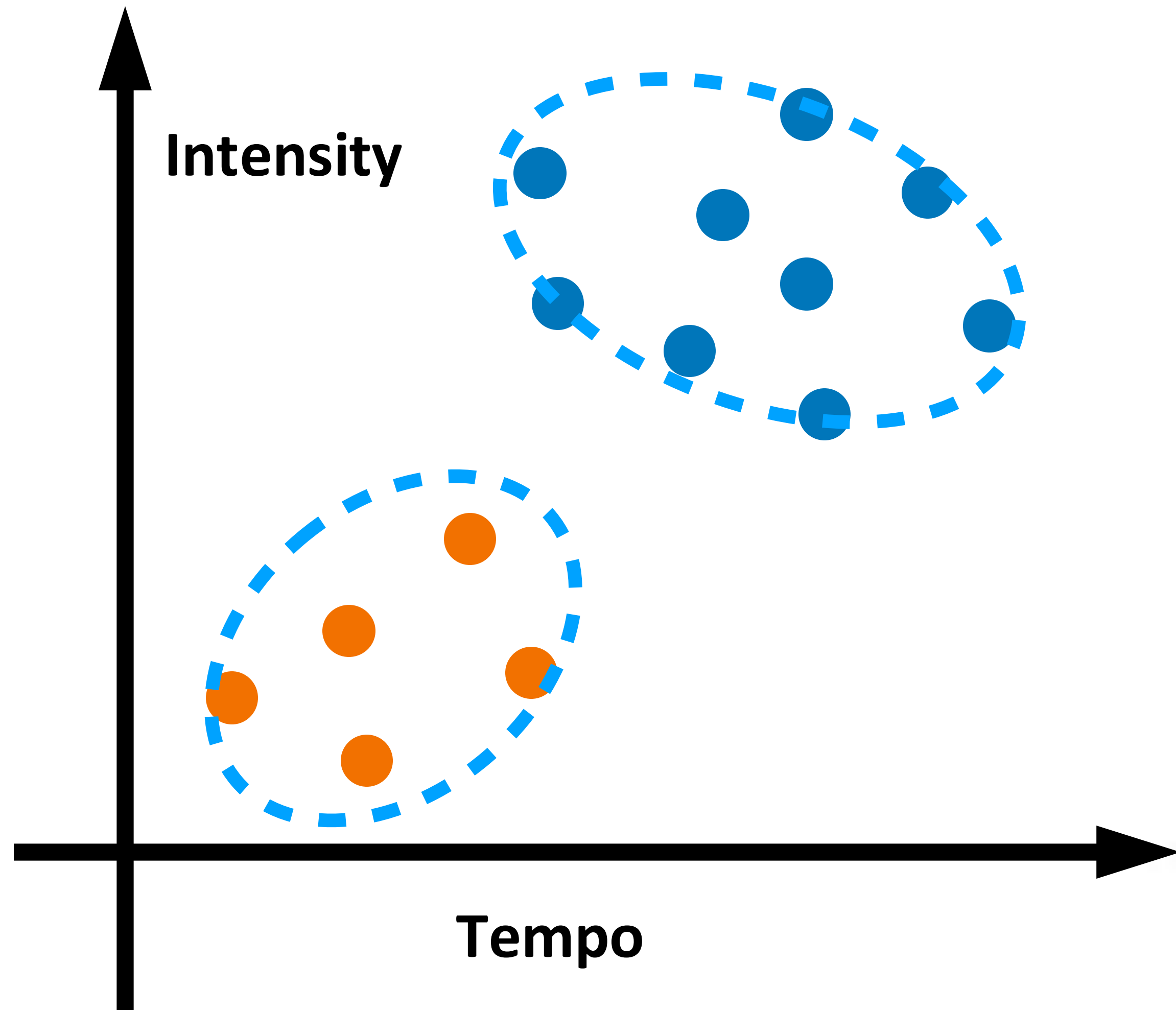
- Take derivatives with respect to it

$$\partial_{\sigma^2} [\cdot] = \frac{n}{2\sigma^2} - \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\Rightarrow \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$



# Classification via MLE



# Classification via MLE

$$\hat{y} = \hat{f}(\mathbf{x}) = \arg \max p(y | \mathbf{x})$$

(Prediction) (Posterior)

# Classification via MLE

$$\begin{aligned} \hat{y} = \hat{f}(\mathbf{x}) &= \arg \max p(y | \mathbf{x}) && \text{(Posterior)} \\ \text{(Prediction)} &&& \\ &= \arg \max_y \frac{p(\mathbf{x} | y) \cdot p(y)}{p(\mathbf{x})} && \text{(by Bayes' rule)} \\ &= \arg \max_y p(\mathbf{x} | y)p(y) \end{aligned}$$

Using labelled training data, learn **class priors** and **class conditionals**





## Part III: Naïve Bayes



# Example 1: Play outside or not?

- If weather is sunny, will my 2-year-old daughter want to play outside?

**Posterior probability**  $p(\text{Yes} \mid \text{☀️})$  vs.  $p(\text{No} \mid \text{☀️})$

# Example 1: Play outside or not?

- If weather is sunny, will my 2-year-old daughter want to play outside?

**Posterior probability**  $p(\text{Yes} \mid \text{☀️})$  vs.  $p(\text{No} \mid \text{☀️})$

- Weather = {Sunny, Rainy, Overcast}
- Play = {Yes, No}
- Observed data {Weather, play on day  $m$ },  $m=\{1,2,\dots,N\}$

# Example 1: Play outside or not?

- If weather is sunny, would you like to play outside?

**Posterior probability**  $p(\text{Yes} \mid \text{☀})$  vs.  $p(\text{No} \mid \text{☀})$

- Weather = {Sunny, Rainy, Overcast}
- Play = {Yes, No}
- Observed data {Weather, play on day  $m$ },  $m=\{1,2,\dots,N\}$

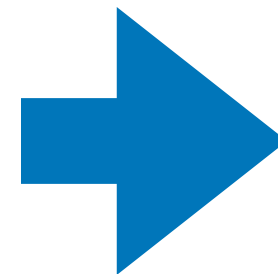
$$p(\text{Play} \mid \text{☀}) = \frac{p(\text{☀} \mid \text{Play}) p(\text{Play})}{p(\text{☀})}$$

**Bayes rule**

# Example 1: Play outside or not?

- **Step 1:** Convert the data to a frequency table of Weather and Play

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No



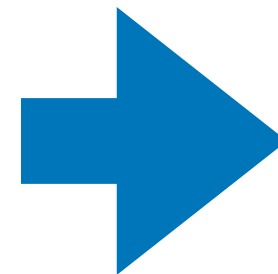
Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9



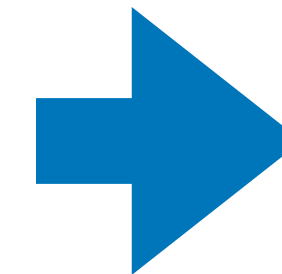
# Example 1: Play outside or not?

- **Step 1:** Convert the data to a frequency table of Weather and Play
- **Step 2:** Based on the frequency table, calculate **likelihoods** and **priors**

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No



Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9



Likelihood table				
Weather	No	Yes		
Overcast		4	=4/14	0.29
Rainy	3	2	=5/14	0.36
Sunny	2	3	=5/14	0.36
All	5	9		
	=5/14	=9/14		
	0.36	0.64		

$$p(\text{Play} = \text{Yes}) = 0.64$$

$$p(\text{☀} | \text{Yes}) = 3/9 = 0.33$$

# Example 1: Play outside or not?

- **Step 3:** Based on the likelihoods and priors, calculate posteriors

$$\begin{aligned} P(\text{Yes} | \text{☀}) \\ = P(\text{☀} | \text{Yes}) * P(\text{Yes}) / P(\text{☀}) \end{aligned} \quad ?$$

$$\begin{aligned} P(\text{No} | \text{☀}) \\ = P(\text{☀} | \text{No}) * P(\text{No}) / P(\text{☀}) \end{aligned} \quad ?$$

# Example 1: Play outside or not?

- **Step 3:** Based on the likelihoods and priors, calculate posteriors

$$\begin{aligned} P(\text{Yes} | \text{☀}) \\ &= P(\text{☀} | \text{Yes}) * P(\text{Yes}) / P(\text{☀}) \\ &= 0.33 * 0.64 / 0.36 \\ &= 0.6 \end{aligned}$$

$$\begin{aligned} P(\text{No} | \text{☀}) \\ &= P(\text{☀} | \text{No}) * P(\text{No}) / P(\text{☀}) \\ &= 0.4 * 0.36 / 0.36 \\ &= 0.4 \end{aligned}$$

$$P(\text{Yes} | \text{☀}) > P(\text{No} | \text{☀}) \quad \text{go outside and play!}$$

# Bayesian classification

$$\hat{y} = \arg \max p(y | \mathbf{x}) \quad (\text{Posterior})$$

(Prediction)

$$= \arg \max \frac{p(\mathbf{x} | y) \cdot p(y)}{p(\mathbf{x})} \quad (\text{by Bayes' rule})$$

$$= \arg \max p(\mathbf{x} | y)p(y)$$



# Bayesian classification

What if  $\mathbf{x}$  has multiple attributes  $\mathbf{x} = \{X_1, \dots, X_k\}$

$$\hat{y} = \arg \max_y p(y | X_1, \dots, X_k) \quad (\text{Posterior})$$

(Prediction)

# Bayesian classification

What if  $\mathbf{x}$  has multiple attributes  $\mathbf{x} = \{X_1, \dots, X_k\}$

$$\begin{aligned} \hat{y} &= \arg \max_y p(y | X_1, \dots, X_k) && \text{(Posterior)} \\ \text{(Prediction)} &&& \\ &= \arg \max_y \frac{p(X_1, \dots, X_k | y) \cdot p(y)}{p(X_1, \dots, X_k)} && \text{(by Bayes' rule)} \\ &&& \uparrow \\ &&& \text{Independent of } y \end{aligned}$$

# Bayesian classification

What if  $\mathbf{x}$  has multiple attributes  $\mathbf{x} = \{X_1, \dots, X_k\}$

$$\hat{y} = \arg \max_y p(y | X_1, \dots, X_k) \quad (\text{Posterior})$$

(Prediction)

$$= \arg \max_y \frac{p(X_1, \dots, X_k | y) \cdot p(y)}{p(X_1, \dots, X_k)} \quad (\text{by Bayes' rule})$$

$$= \arg \max_y p(X_1, \dots, X_k | y) p(y)$$

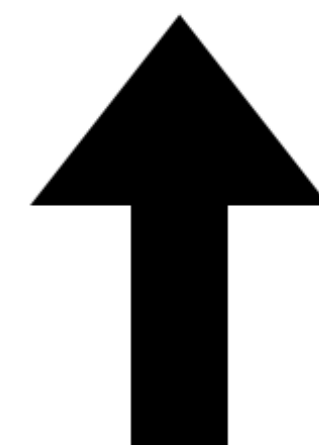
Class conditional  
likelihood

Class prior

# Naïve Bayes Assumption

Conditional independence of feature attributes

$$p(X_1, \dots, X_k | y)p(y) = \prod_{i=1}^k p(X_i | y)p(y)$$



Easier to estimate  
(using MLE!)





**Thanks!**