

# Data Mining Techniques

CS 6220 - Section 3 - Fall 2016

## Lecture 4: Classification

Jan-Willem van de Meent

(*credit*: Zhao, CS 229, Bishop, Hastie)



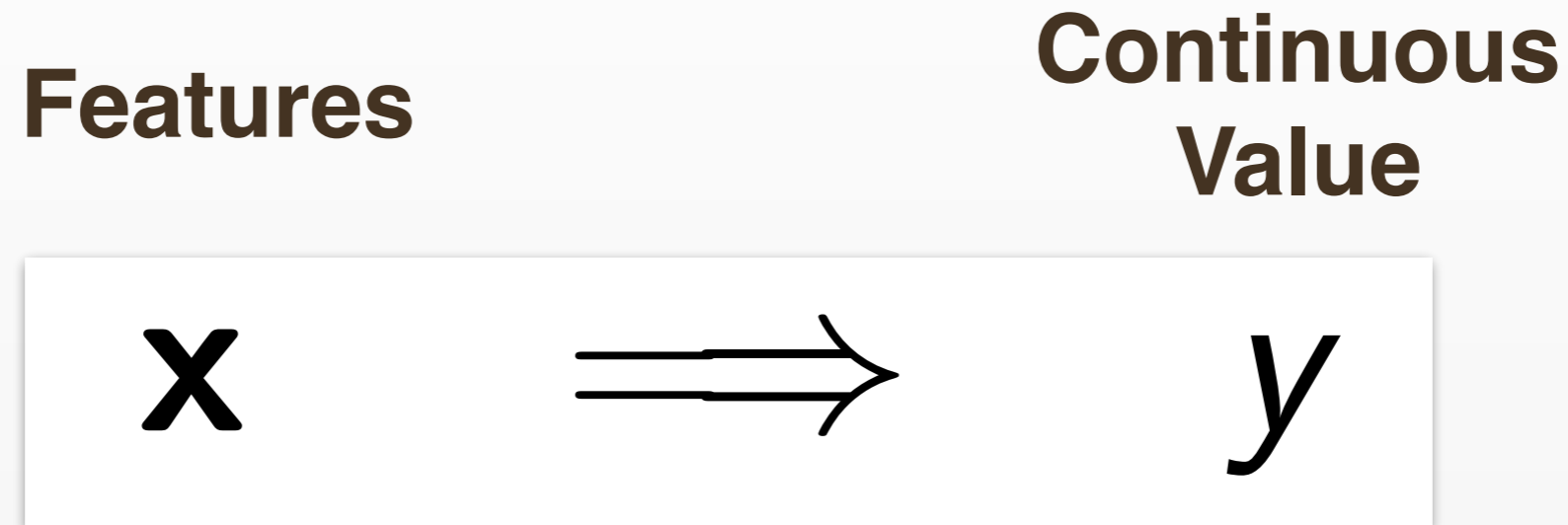
# Homework Problems

## Homework 1 is out on blackboard

- Due: **Friday 30 Sep**
- Can use any language (within reason)
- **Discussion is encouraged, but submissions must be completed individually**  
(absolutely **no** sharing of code)
- Submit as **zip** file ~~via e-mail~~ **on blackboard** by **11.59pm** on day of deadline (no late submissions)
- Please follow *submission guidelines* on website  
(TA's have authority to deduct points)

# Classification

# Regression Examples



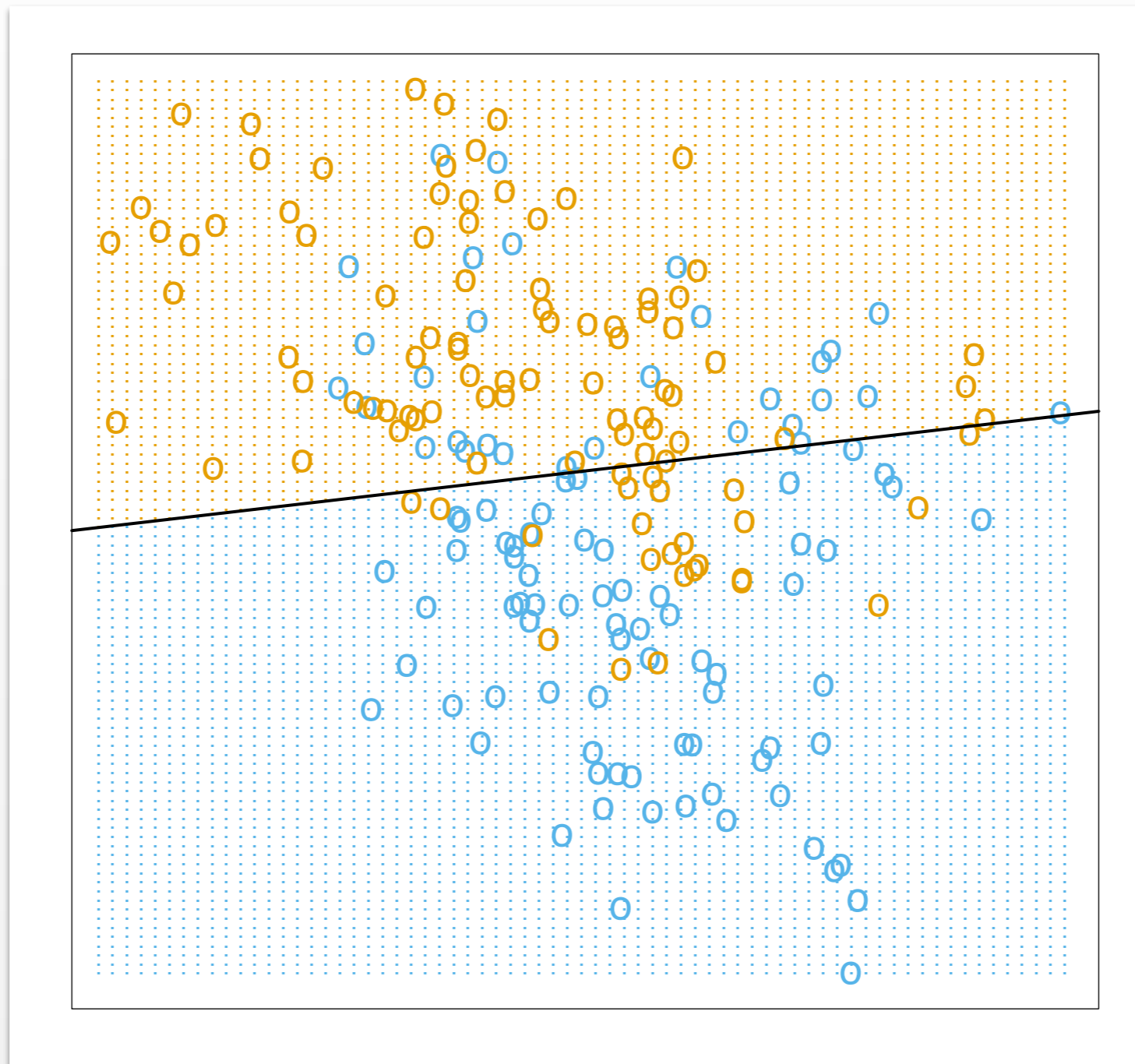
- {age, major, gender, race} ⇒ GPA
- {income, credit score, profession} ⇒ Loan Amount
- {college, major, GPA} ⇒ Future Income

# Classification Example



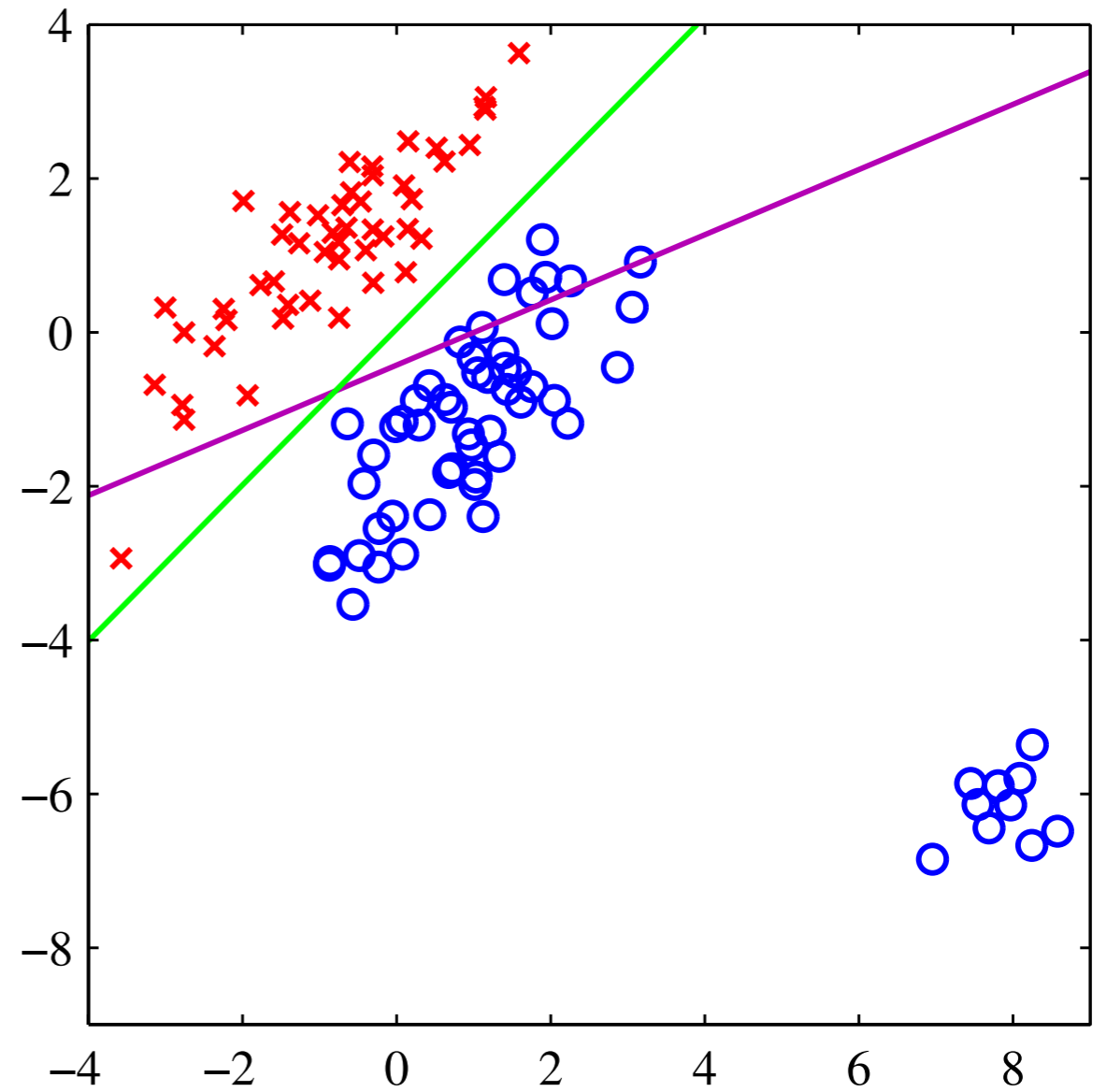
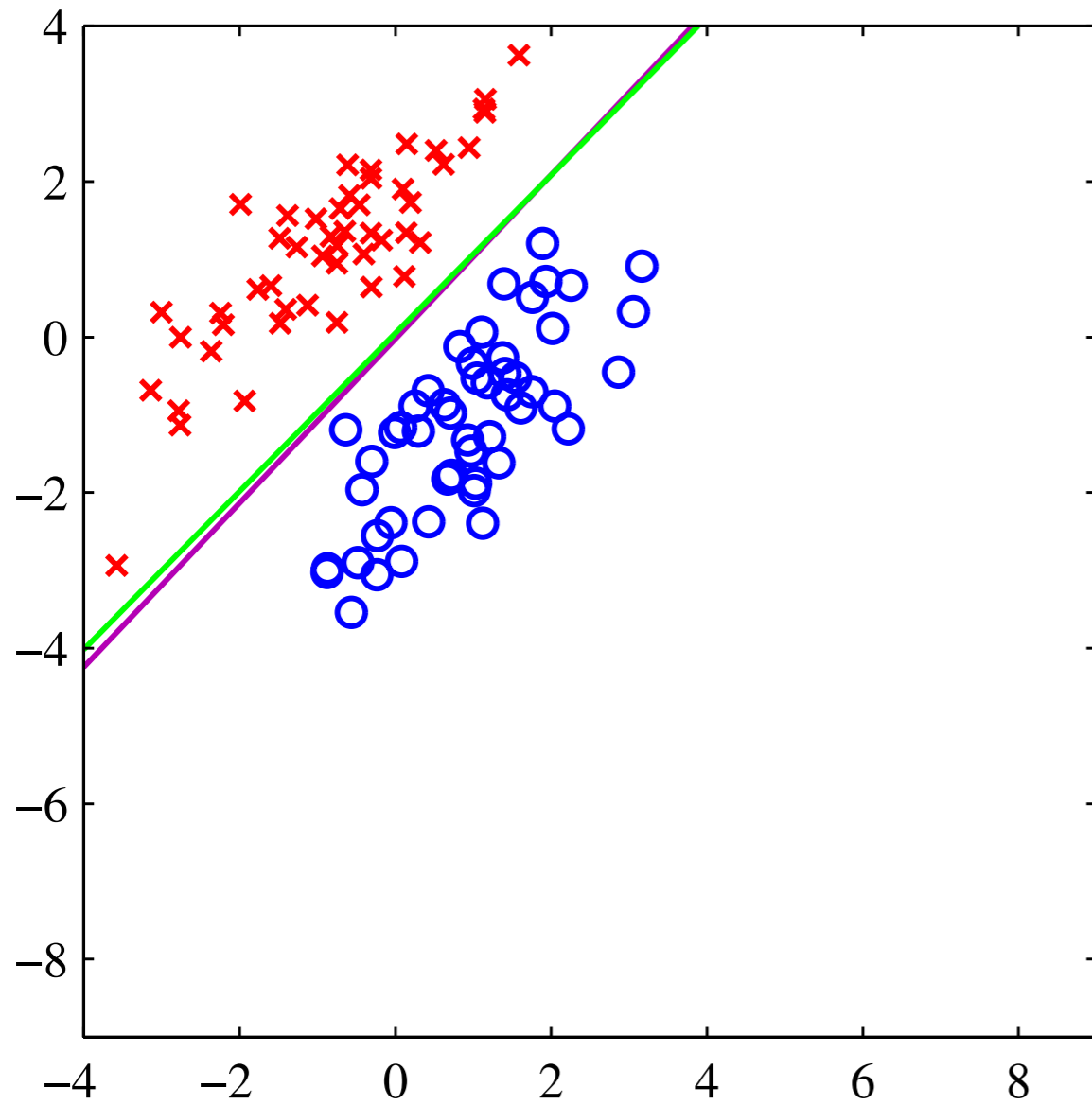
- {words in e-mail}  $\Rightarrow$  SPAM or not?
- {income, credit score, profession}  $\Rightarrow$  Award Loan?
- {pixels in image}  $\Rightarrow$  Cat, dog, or pony?

# Perceptron (Linear Regression)

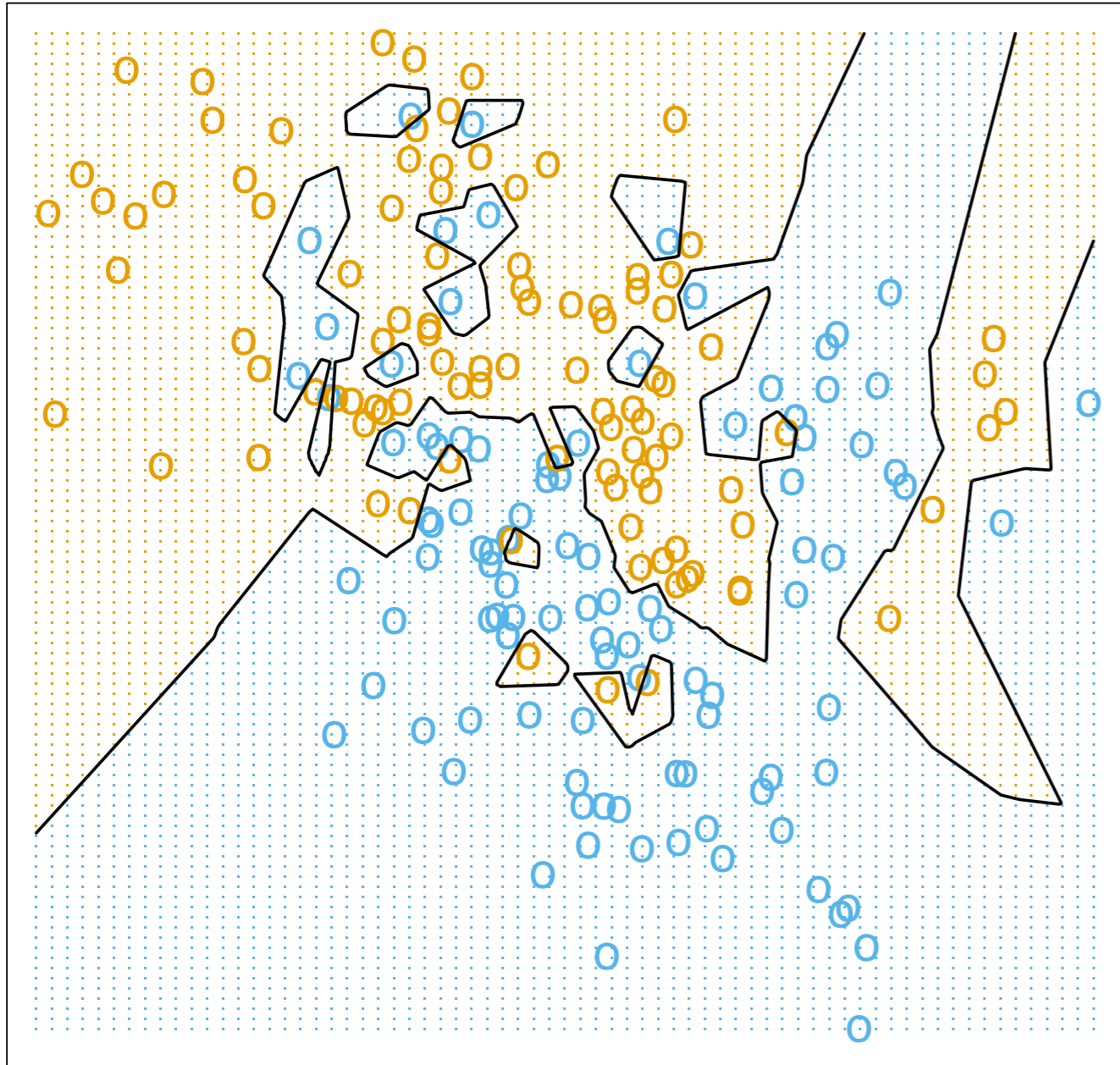


For what types of datasets would this not work well?

# Perceptron (Linear Regression)



# k-Nearest Neighbors



1 Nearest Neighbor

## Algorithm

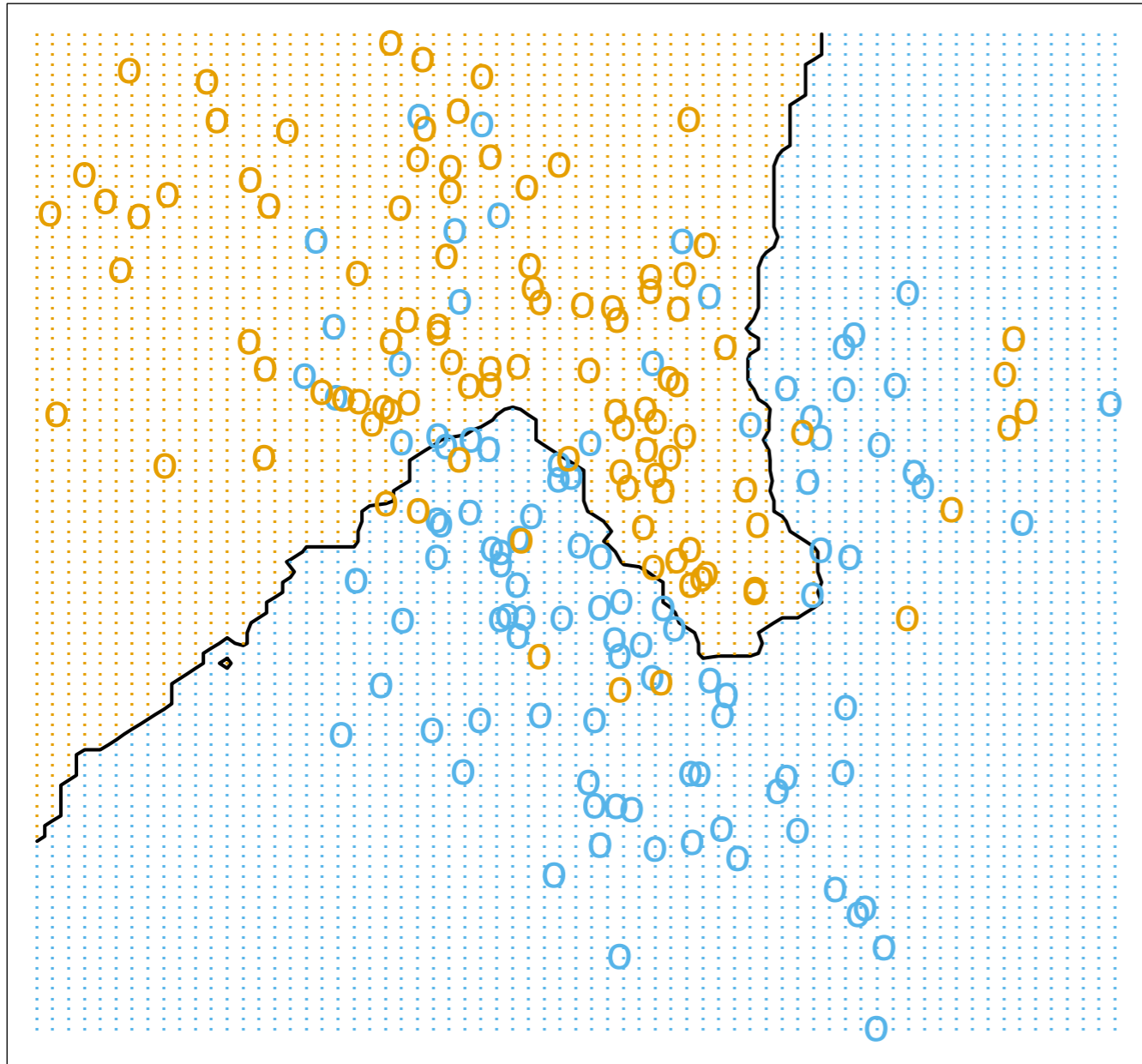
- Find  $k$  nearest points to  $\mathbf{x}^*$
- Predict  $y^*$  according to majority vote

## Properties

- *Lazy*  
No learned parameters.
- *Instance Learner*  
Needs all data at test time.



# k-Nearest Neighbors



15 Nearest Neighbors

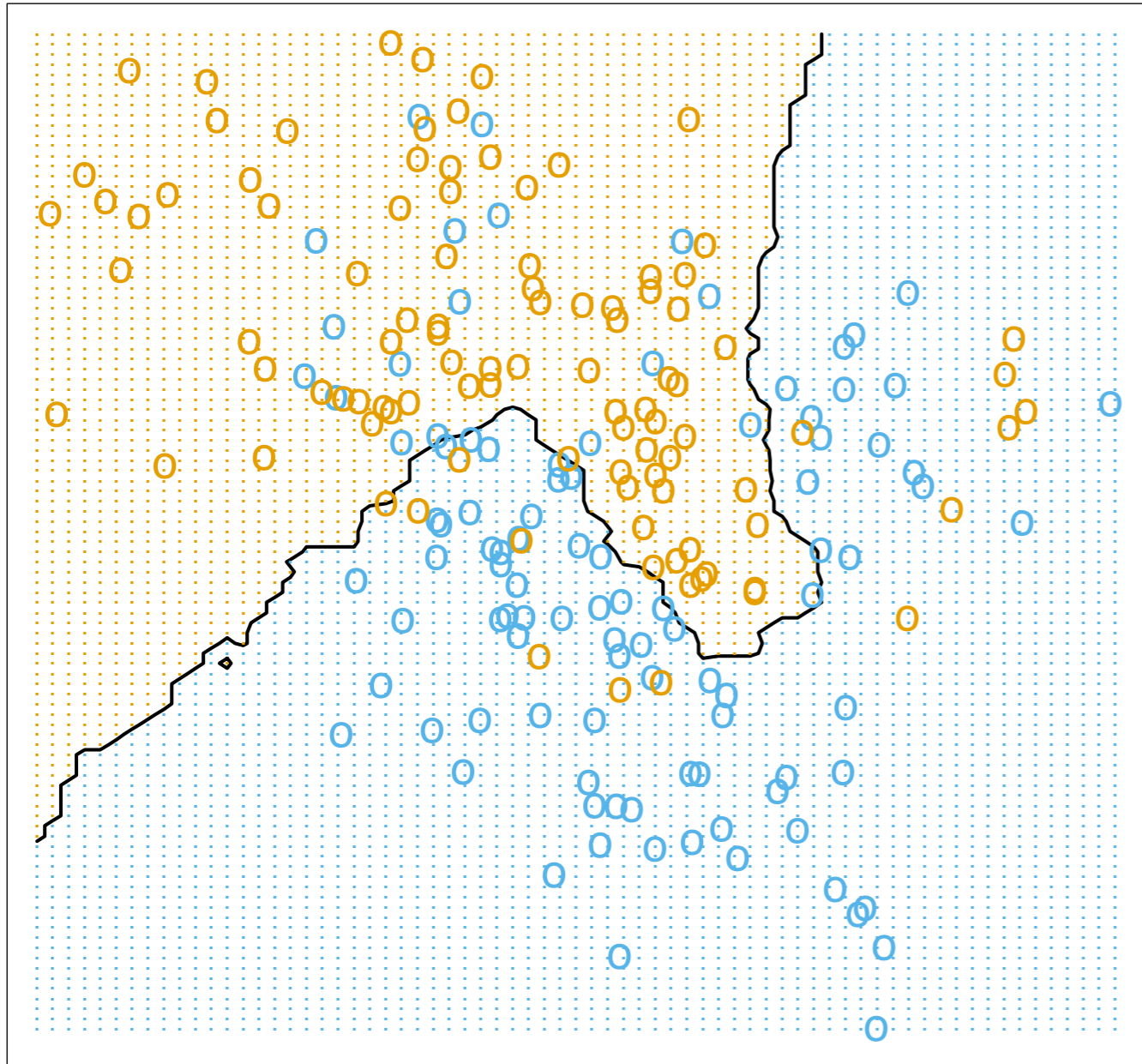
## Algorithm

- Find  $k$  nearest points to  $\mathbf{x}^*$
- Predict  $y^*$  according to majority vote

## Properties

- *Lazy*  
No learned parameters.
- *Instance Learner*  
Needs all data at test time.

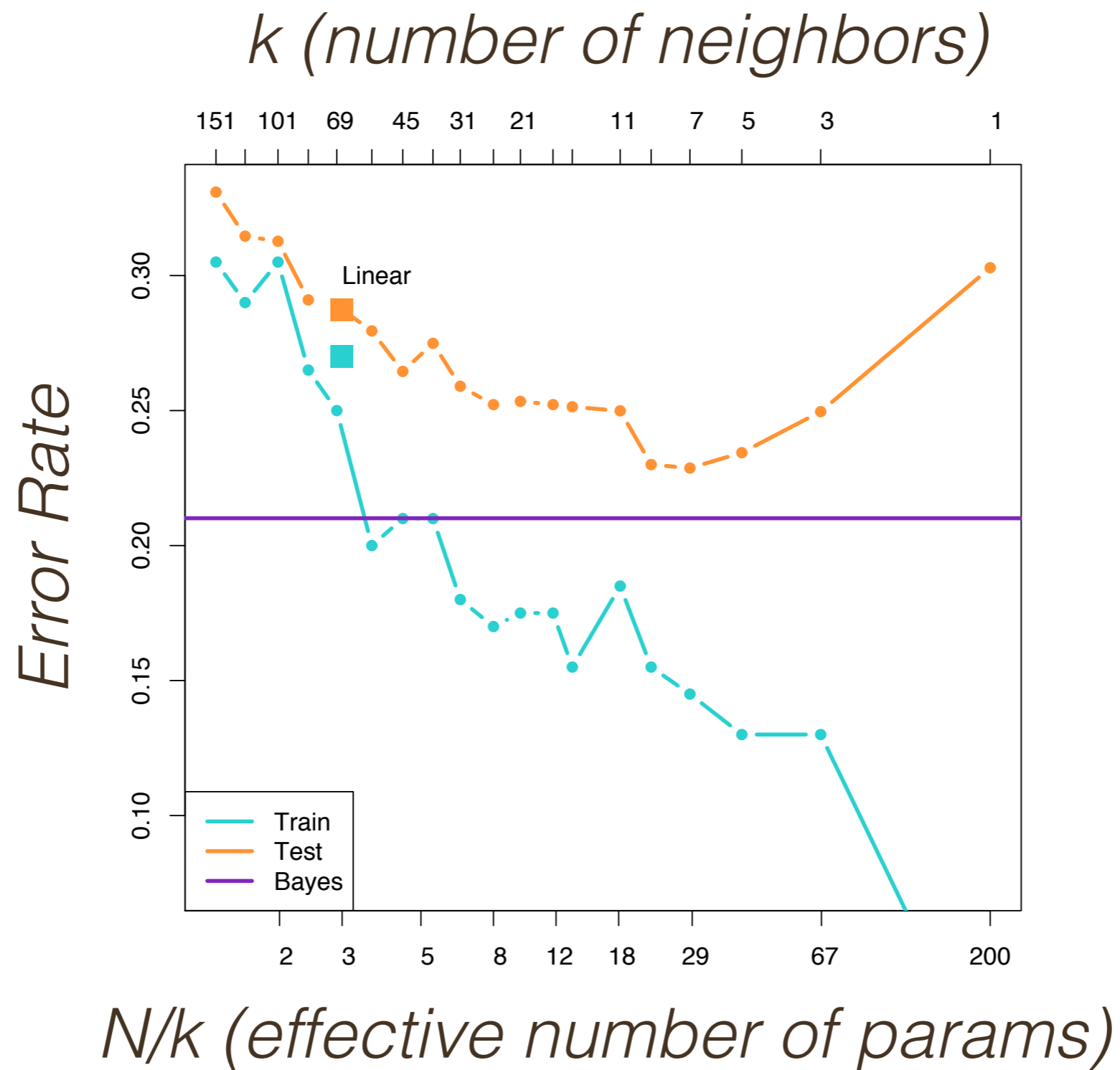
# k-Nearest Neighbors



15 Nearest Neighbors

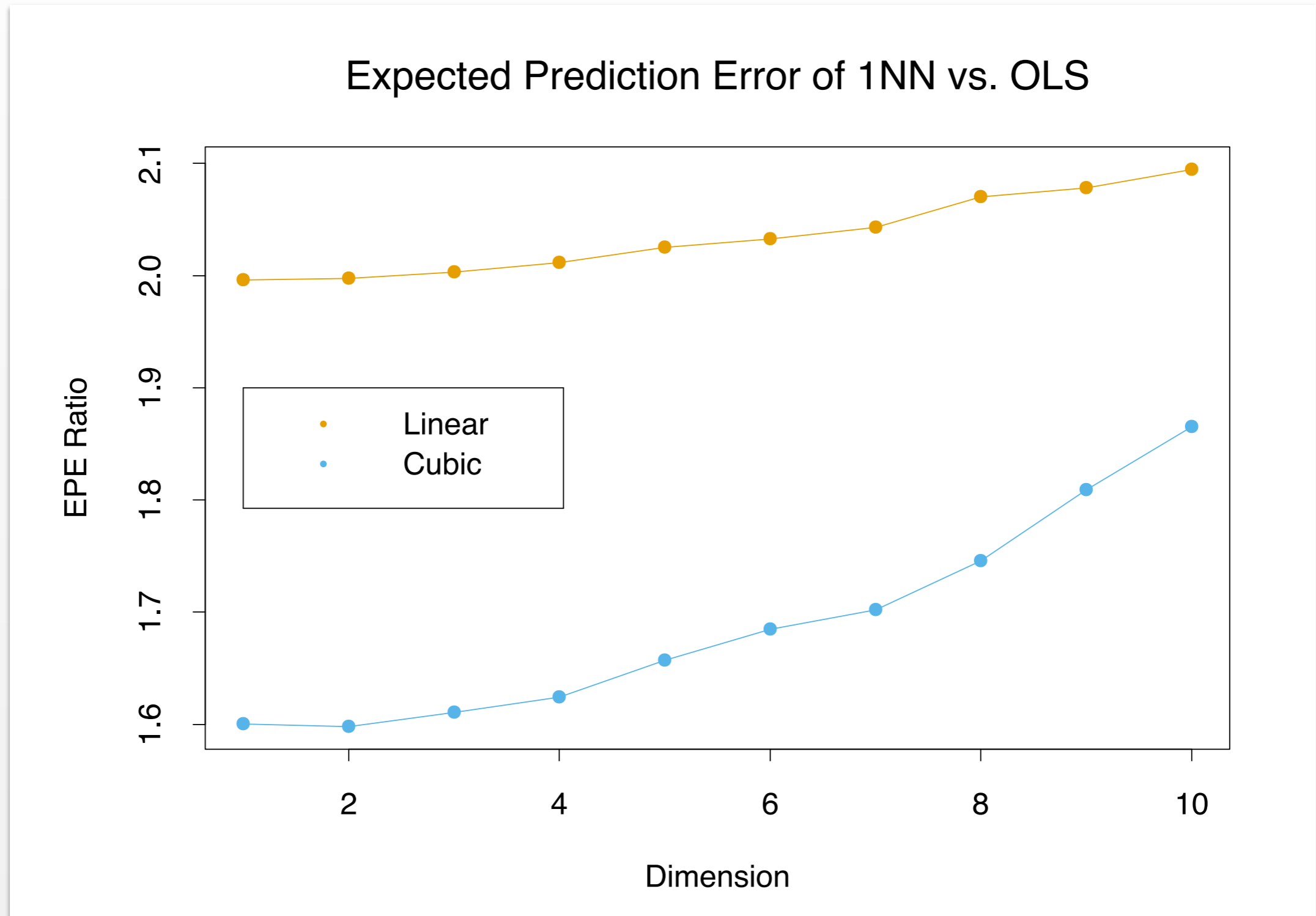
**How should we choose  $k$ ?**

# Choosing $k$

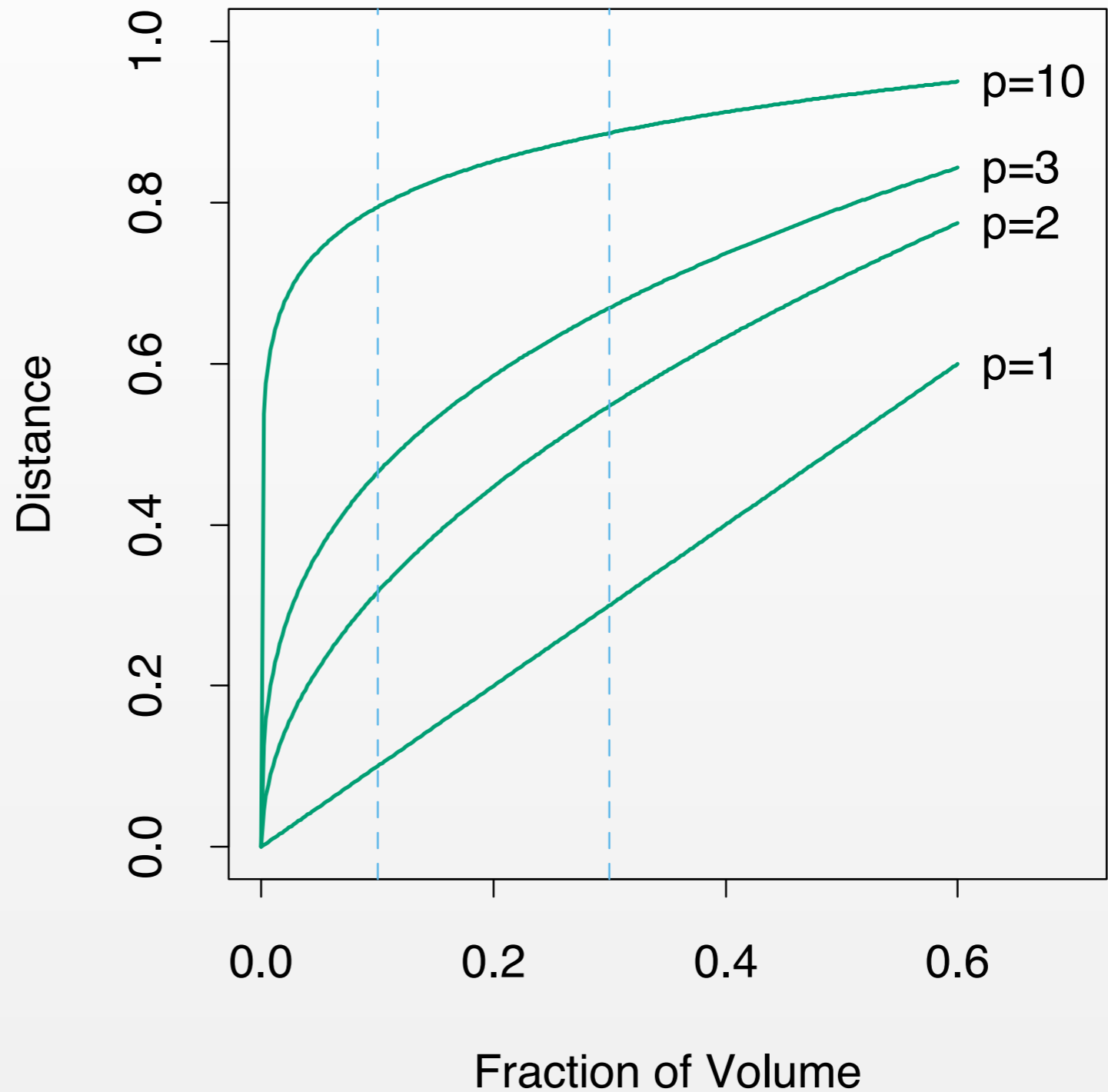
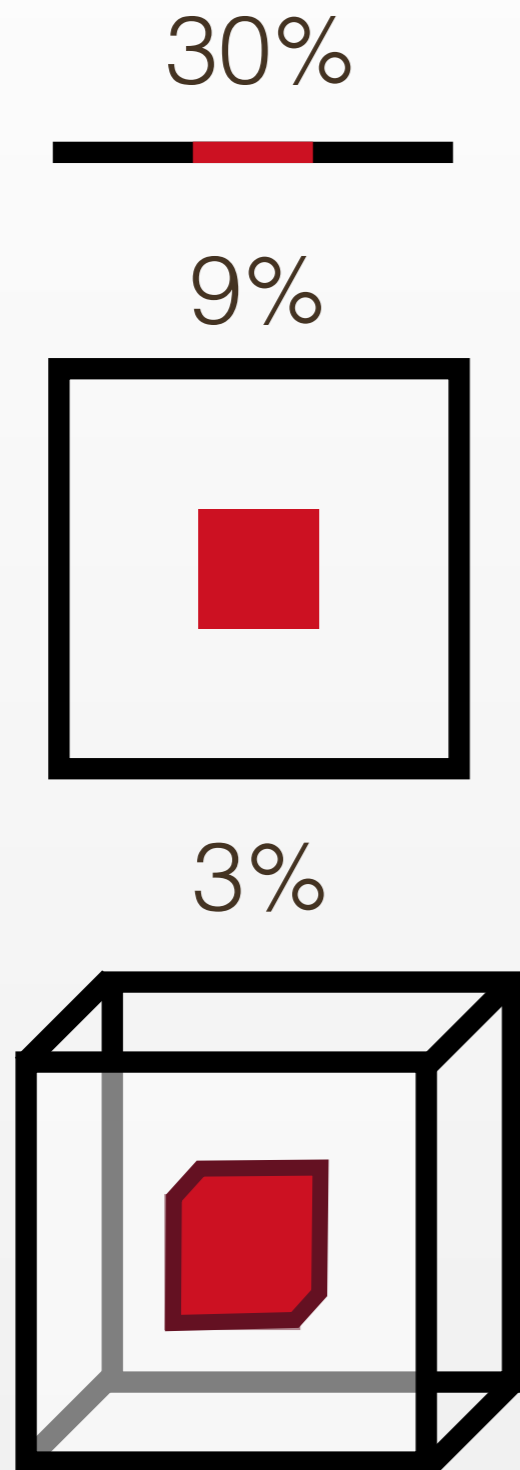


When might  
k-Nearest Neighbors  
not perform well?

# Curse of Dimensionality

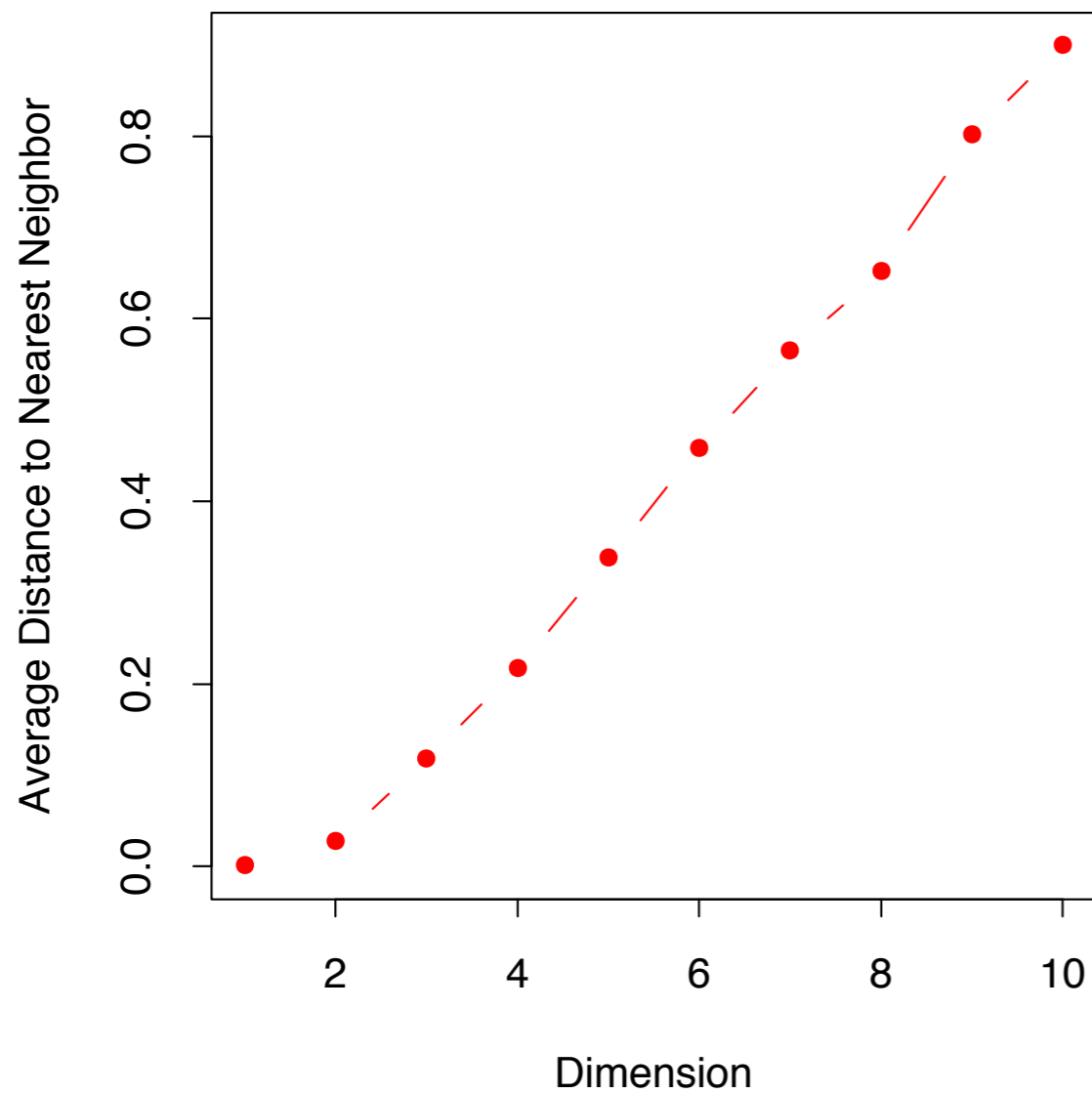


# Curse of Dimensionality

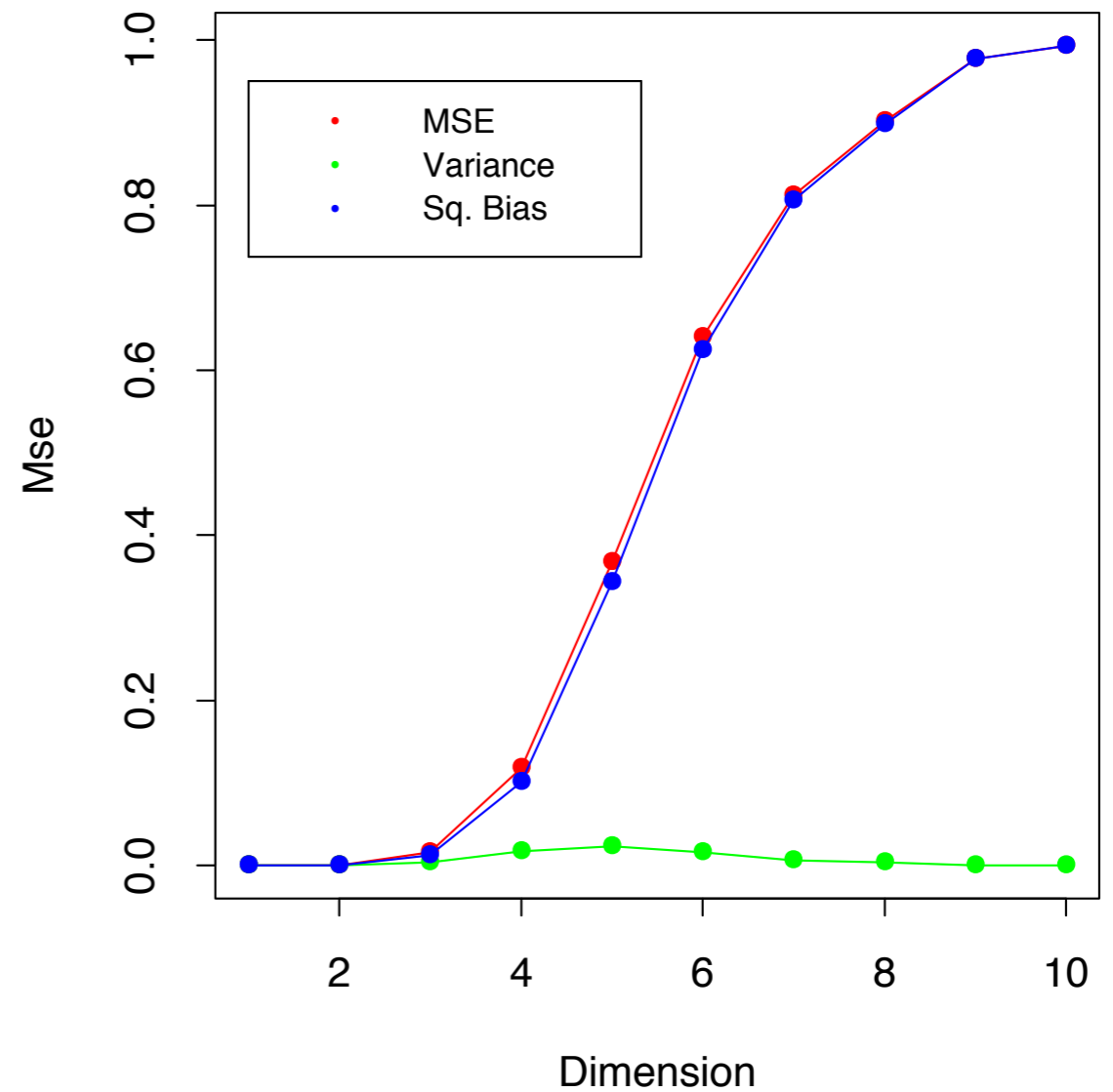


# Curse of Dimensionality

Distance to 1-NN vs. Dimension



MSE vs. Dimension



# Distance Norms

- Euclidean Distance

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

- Mahattan Distance

$$\sum_{i=1}^k |x_i - y_i|$$

- Minkowski Distance

$$\left( \sum_{i=1}^k (|x_i - y_i|)^q \right)^{\frac{1}{q}}$$



# Distance Norms

For nominal variables: Hamming Distance

$$D_H(x, y) = \sum_{i=1}^k I(x_i, y_i)$$

where

$$I(a, b) = \begin{cases} 1 & \text{if } a \neq b \\ 0 & \text{if } a = b \end{cases}$$

Examples:

- $D_H(\text{"Karolin"}, \text{"Kathrin"}) = 3$
- $D_H(\text{"Karolin"}, \text{"Kerstin"}) = 3$
- $D_H(\text{"1011101"}, \text{"1001001"}) = 2$

# Sensitivity to Scaling

Age	Loan	Default	Distance	
25	\$40,000	N	102000	
35	\$60,000	N	82000	
45	\$80,000	N	62000	
20	\$20,000	N	122000	
35	\$120,000	N	22000	2
52	\$18,000	N	124000	
23	\$95,000	Y	47000	
40	\$62,000	Y	80000	
60	\$100,000	Y	42000	3
48	\$220,000	Y	78000	
33	\$150,000	Y	8000	1
<b>48</b>	<b>\$142,000</b>	<b>?</b>		

Euclidean Distance

$$D = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

# Sensitivity to Scaling

Age	Loan	Default	Distance
0.125	0.11	N	0.7652
0.375	0.21	N	0.5200
0.625	0.31	N	0.3160
0	0.01	N	0.9245
0.375	0.50	N	0.3428
0.8	0.00	N	0.6220
0.075	0.38	Y	0.6669
0.5	0.22	Y	0.4437
1	0.41	Y	0.3650
0.7	1.00	Y	0.3861
0.325	0.65	Y	0.3771
<b>0.7</b>	<b>0.61</b>	<b>?</b>	

Standardized Variable

$$X_s = \frac{X - Min}{Max - Min}$$

# Normalization Strategies

- Min-Max:  $\frac{X - X_{min}}{X_{max} - X_{min}}$

$$a + \frac{(X - X_{min})(b - a)}{X_{max} - X_{min}}$$

- Z-score:  $X \sim N(\mu, \sigma): \frac{X - \mu}{\sigma}$

- Scaling:  $\frac{X}{X_{max}}$

# Linear Methods

# Linear Classifiers

$$y_n = f(\mathbf{w}^\top \mathbf{x}_n)$$

## Examples

- Logistic Regression
- Linear Discriminant Analysis
- Naive Bayes

# Logistic Regression

$$\log \frac{p(y = k | \mathbf{x})}{p(y = K | \mathbf{x})} = \mathbf{w}_k^\top \mathbf{x}$$

$$k = 1, \dots, K - 1$$

# Logistic Regression

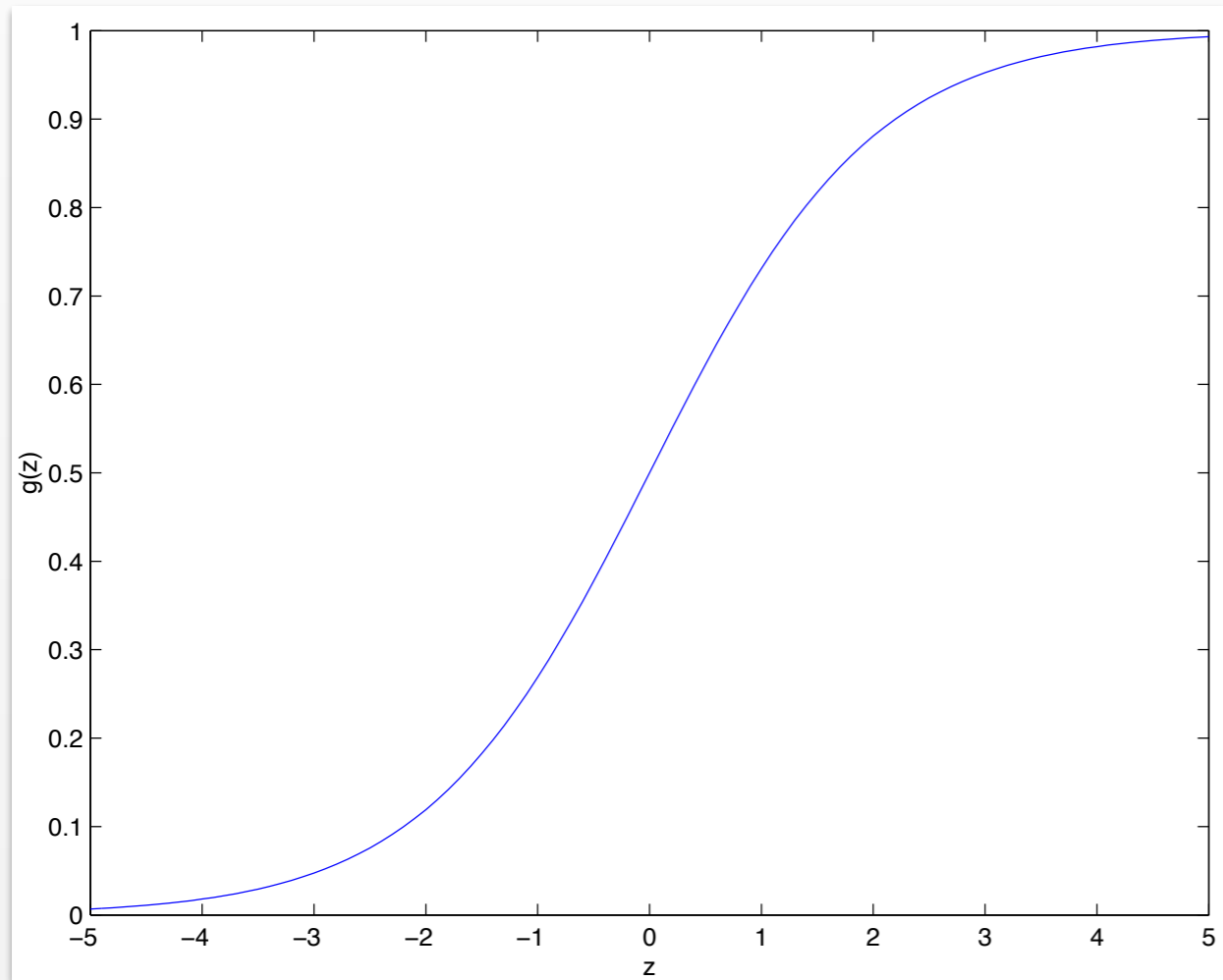
$$\log \frac{p(y = k | \mathbf{x})}{p(y = K | \mathbf{x})} = \mathbf{w}_k^\top \mathbf{x} \quad k = 1, \dots, K - 1$$

$$\frac{p(y = k | \mathbf{x})}{\sum_l p(y = l | \mathbf{x})} = \frac{\exp^{\mathbf{w}_k^\top \mathbf{x}}}{1 + \sum_{l=1}^{K-1} \exp^{\mathbf{w}_l^\top \mathbf{x}}} \quad k = 1, \dots, K - 1$$

$$\frac{p(y = K | \mathbf{x})}{\sum_l p(y = l | \mathbf{x})} = \frac{1}{1 + \sum_{l=1}^{K-1} \exp^{\mathbf{w}_l^\top \mathbf{x}}}$$



# Logistic Regression



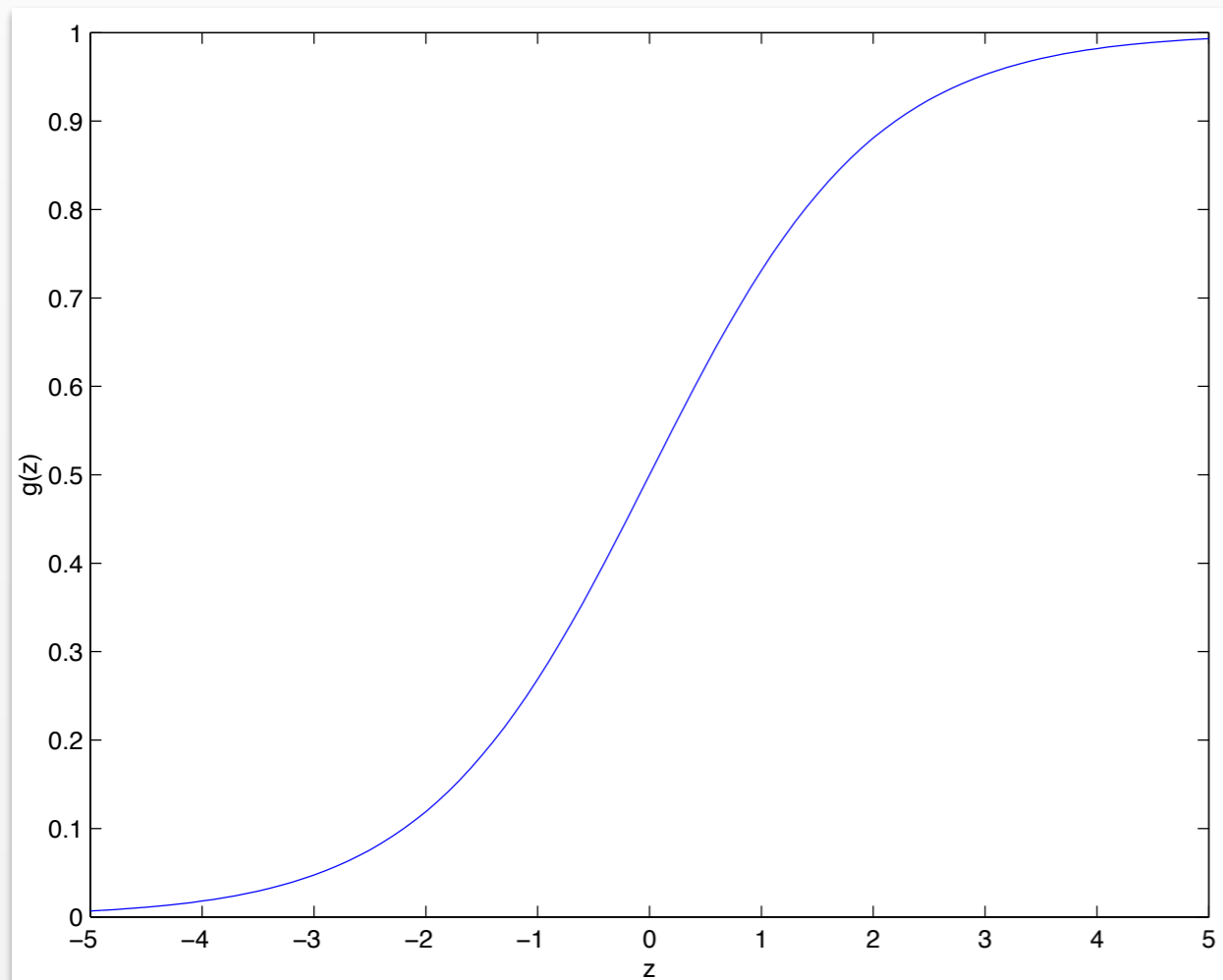
Binary Classification:

$$p(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp^{-\mathbf{w}^\top \mathbf{x}}}$$

Logistic Function:

$$\sigma(z) = \frac{1}{1 + \exp^{-z}}$$

# Logistic Regression



Binary Classification:

$$p(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp^{-\mathbf{w}^\top \mathbf{x}}}$$

Logistic Function:

$$\sigma(z) = \frac{1}{1 + \exp^{-z}}$$

$$y = \sigma(\mathbf{w}^\top \mathbf{x})$$

# Logistic Regression

Maximum Likelihood

$$p(\mathbf{y} | \mathbf{w}) = \prod_{n=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_n)^{y_n} (1 - \sigma(\mathbf{w}^\top \mathbf{x}_n))^{1-y_n}$$

$$E(\mathbf{w}) = -\log p(\mathbf{y} | \mathbf{w})$$

$$= -\sum_{n=1}^N y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n)$$

$$+ (1 - y_n) \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}_n))$$

# Logistic Regression

Gradient

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = - \sum_{n=1}^N \sigma'(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n \left( \frac{y_n}{\sigma(\mathbf{w}^\top \mathbf{x}_n)} - \frac{1 - y_n}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)} \right)$$

$$\sigma(z) = \frac{1}{1 + \exp^{-z}}$$

$$\sigma'(z) = ?$$

# Logistic Regression

## Gradient

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = - \sum_{n=1}^N \sigma'(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n \left( \frac{y_n}{\sigma(\mathbf{w}^\top \mathbf{x}_n)} - \frac{1 - y_n}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)} \right)$$

$$\sigma(z) = \frac{1}{1 + \exp^{-z}}$$

$$\begin{aligned} \sigma'(z) &= \frac{1}{(1 + \exp^{-z})^2} \exp^{-z} \\ &= \frac{1}{(1 + \exp^{-z})} \frac{\exp^{-z}}{(1 + \exp^{-z})} \\ &= \sigma(z)(1 - \sigma(z)) \end{aligned}$$

# Logistic Regression

Gradient

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = - \sum_{n=1}^N \sigma'(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n \left( \frac{y_n}{\sigma(\mathbf{w}^\top \mathbf{x}_n)} - \frac{1 - y_n}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)} \right)$$

$$\sigma(z) = \frac{1}{1 + \exp^{-z}}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

# Logistic Regression

Gradient

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = - \sum_{n=1}^N \sigma'(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n \left( \frac{y_n}{\sigma(\mathbf{w}^\top \mathbf{x}_n)} - \frac{1 - y_n}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)} \right)$$

# Logistic Regression

## Gradient

$$\begin{aligned}\nabla_{\mathbf{w}} E(\mathbf{w}) &= - \sum_{n=1}^N \sigma'(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n \left( \frac{y_n}{\sigma(\mathbf{w}^\top \mathbf{x}_n)} - \frac{1 - y_n}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)} \right) \\ &= - \sum_{n=1}^N \mathbf{x}_n \left( (1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)) y_n - (1 - y_n) \sigma(\mathbf{w}^\top \mathbf{x}_n) \right)\end{aligned}$$

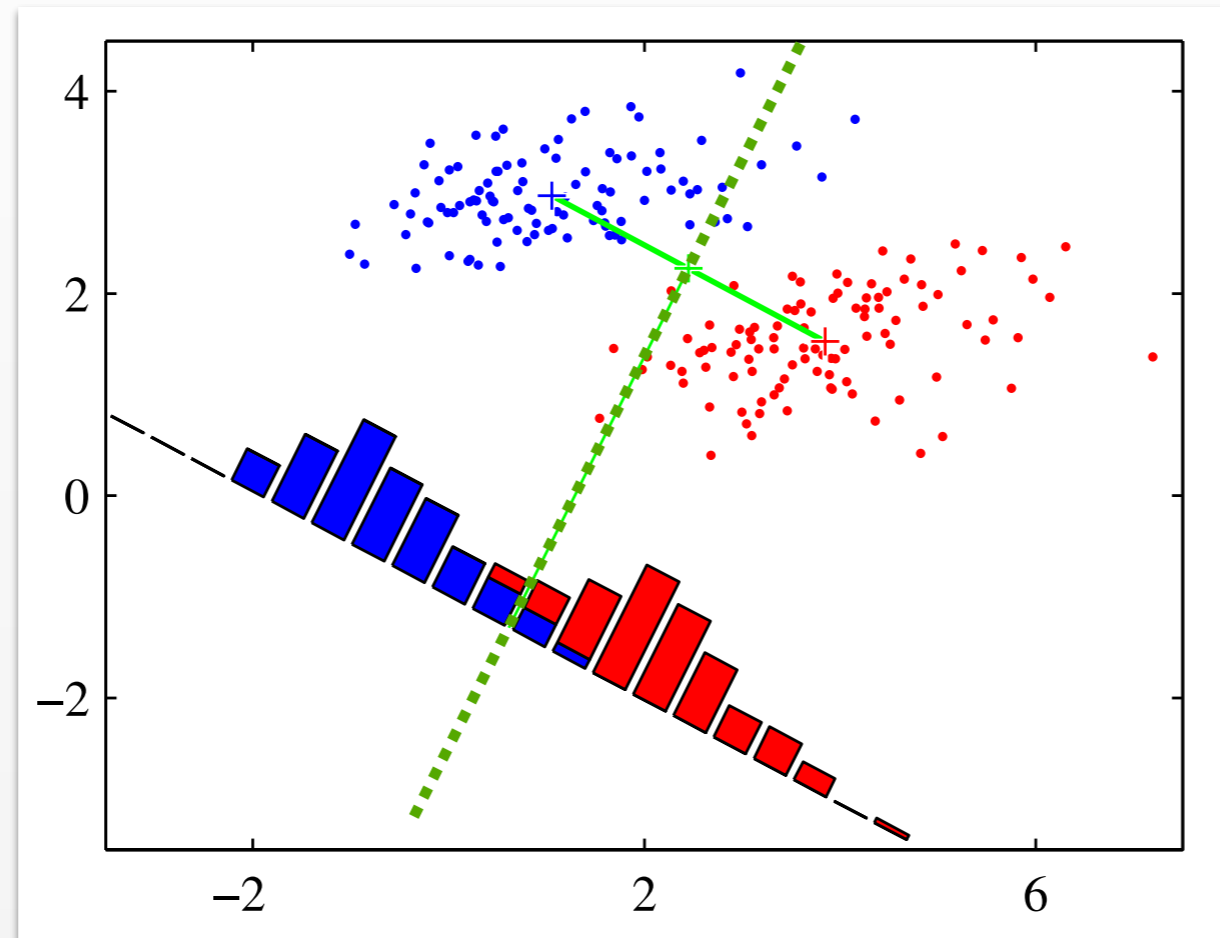


# Logistic Regression

## Gradient

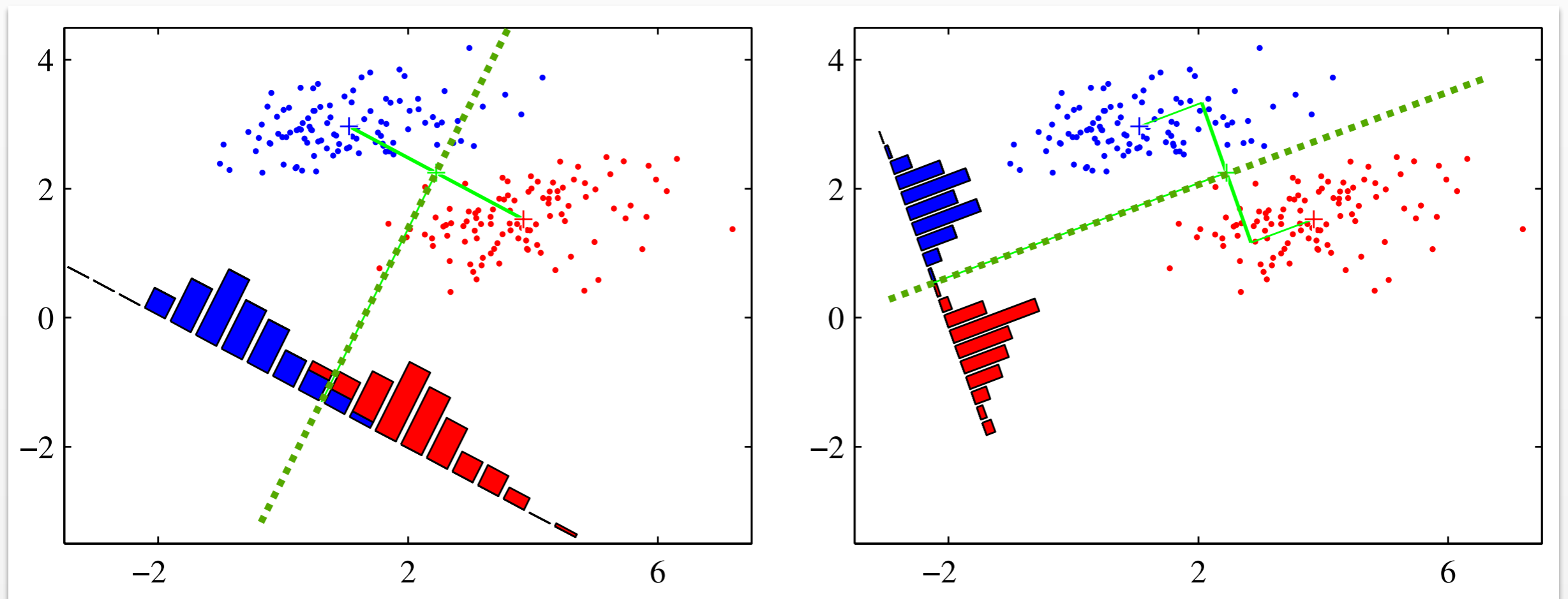
$$\begin{aligned}\nabla_{\mathbf{w}} E(\mathbf{w}) &= - \sum_{n=1}^N \sigma'(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n \left( \frac{y_n}{\sigma(\mathbf{w}^\top \mathbf{x}_n)} - \frac{1 - y_n}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)} \right) \\ &= - \sum_{n=1}^N \mathbf{x}_n \left( (1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)) y_n - (1 - y_n) \sigma(\mathbf{w}^\top \mathbf{x}_n) \right) \\ &= - \sum_{n=1}^N \mathbf{x}_n (y_n - \sigma(\mathbf{w}^\top \mathbf{x}_n))\end{aligned}$$

# Linear Discriminant Analysis



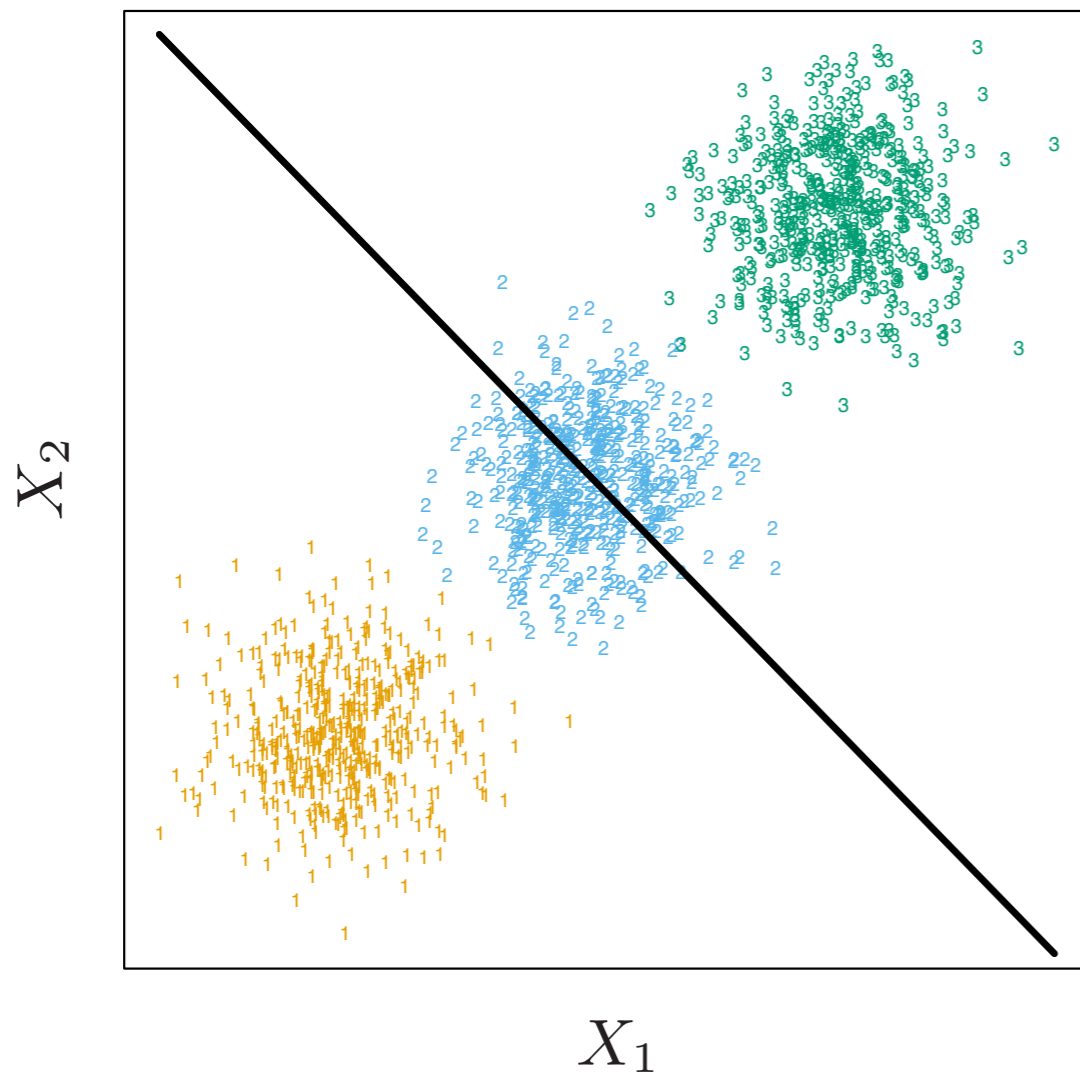
**Idea: Calculate center for each class**

# Linear Discriminant Analysis

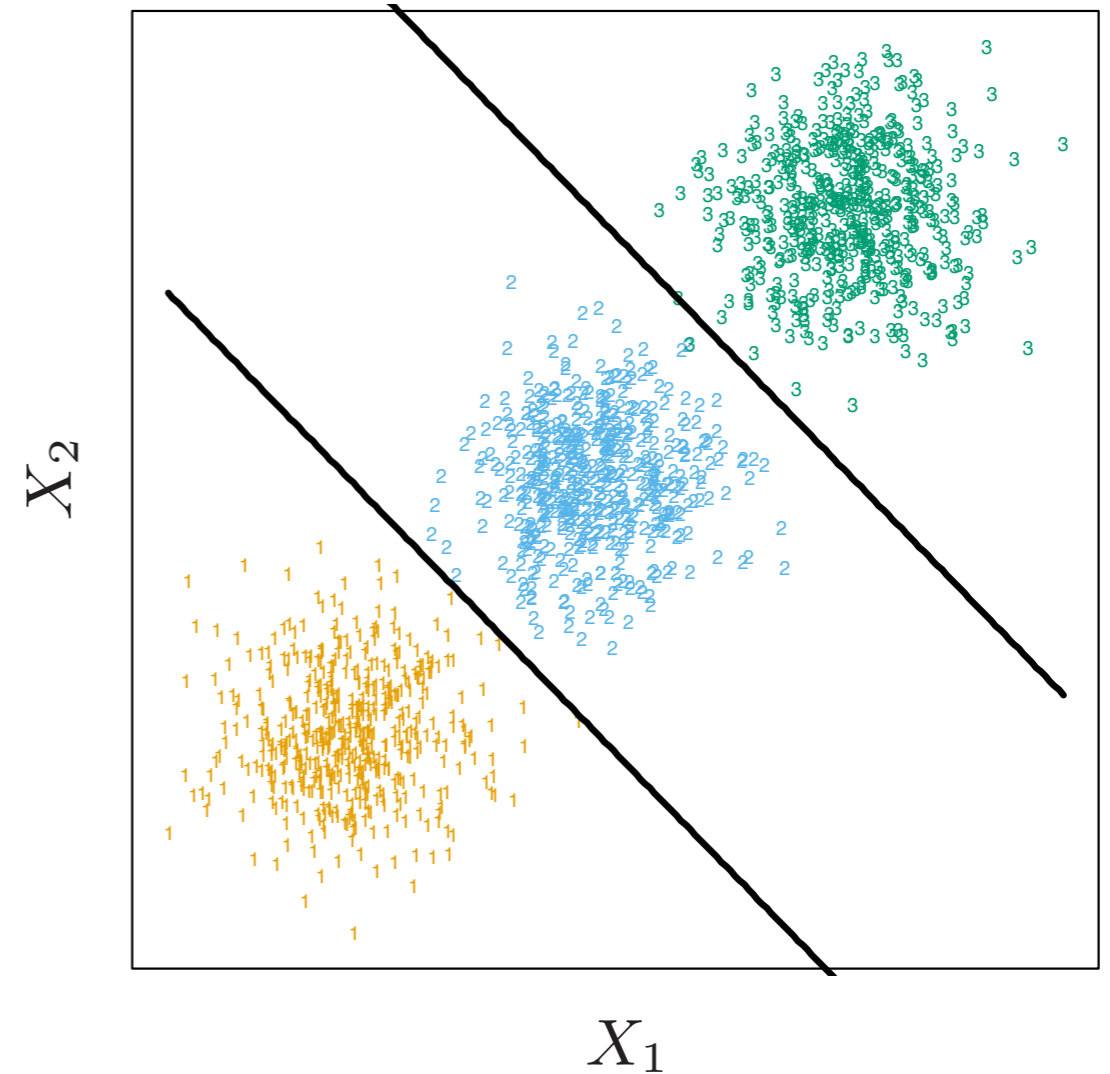


# Linear Discriminant Analysis

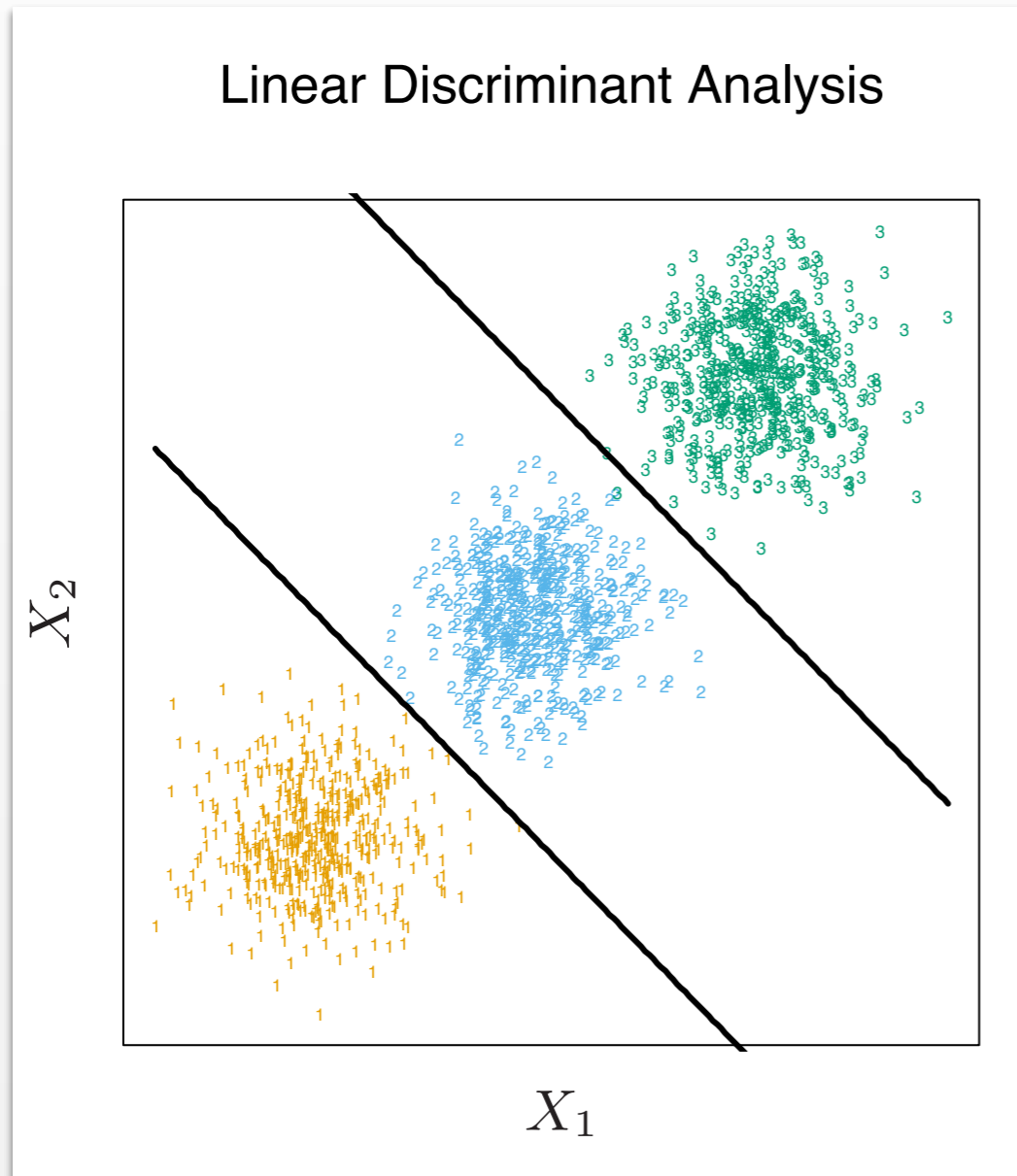
Linear Regression



Linear Discriminant Analysis



# Linear Discriminant Analysis



## Algorithm

- Mean for each class

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n:y_n=k} \mathbf{x}_n$$

- Covariance for each class

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n:y_n=k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

- Average covariance

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_k N_k \boldsymbol{\Sigma}_k$$

# Linear Discriminant Analysis

## Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

## Linear Discriminant Analysis

$$p(\mathbf{x} | y = k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

## Quadratic Discriminant Analysis

$$p(\mathbf{x} | y = k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

## Algorithm

- Mean for each class

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n:y_n=k} \mathbf{x}_n$$

- Covariance for each class

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n:y_n=k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

- Average covariance

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_k N_k \boldsymbol{\Sigma}_k$$

# Linear Discriminant Analysis

## Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

## Linear Discriminant Analysis

$$p(\mathbf{x} | y = k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

## Quadratic Discriminant Analysis

$$p(\mathbf{x} | y = k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

**What is “linear” or “quadratic”  
about these methods?**

# Linear Discriminant Analysis

$$\log \frac{p(\mathbf{x} | y = k)}{p(\mathbf{x} | y = l)} = \frac{1}{2}(\mathbf{x} - \mu_l)^\top \Sigma^{-1}(\mathbf{x} - \mu_l) - \frac{1}{2}(\mathbf{x} - \mu_k)^\top \Sigma^{-1}(\mathbf{x} - \mu_k)$$



# Linear Discriminant Analysis

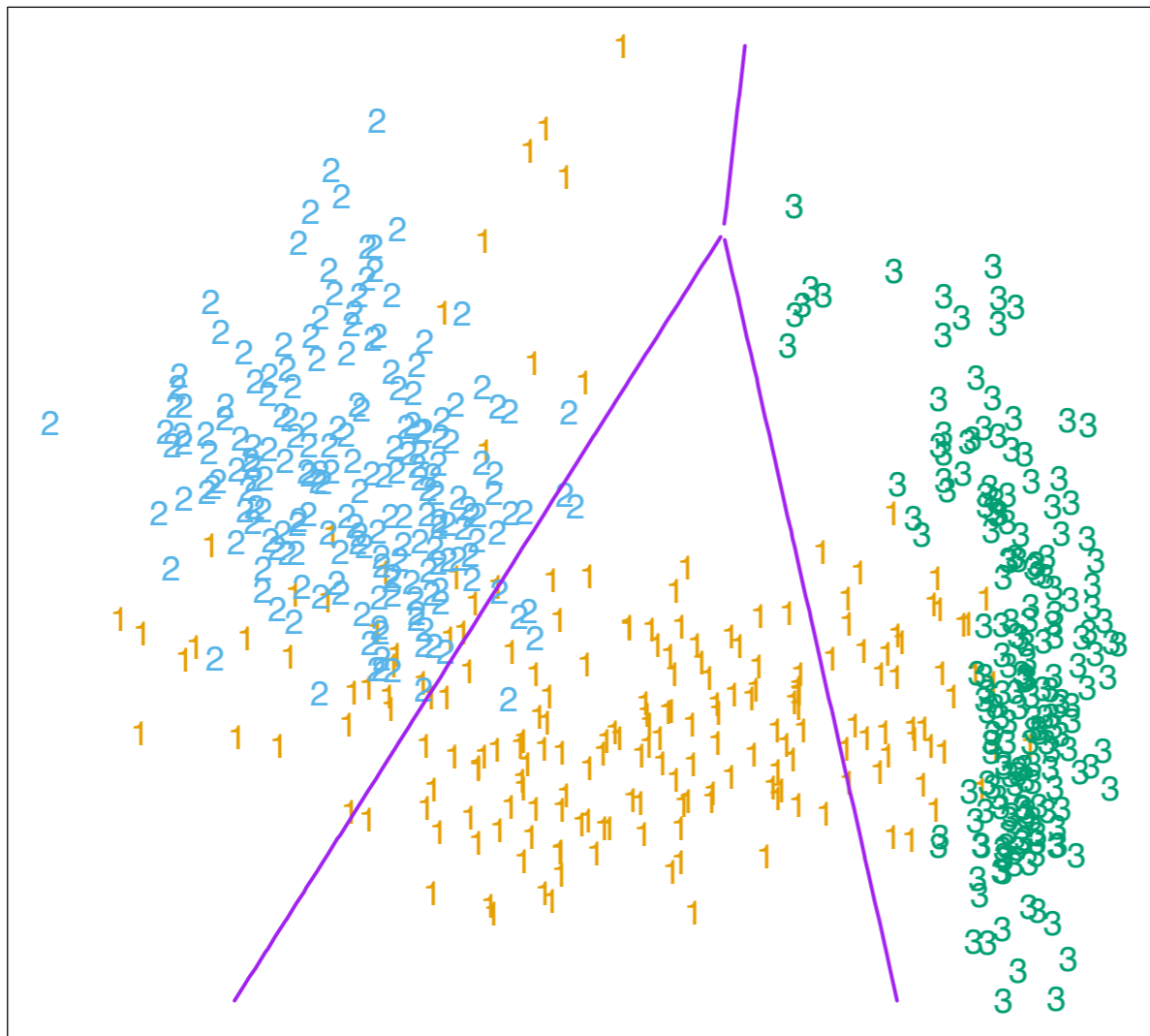
$$\begin{aligned}\log \frac{p(\mathbf{x} | y = k)}{p(\mathbf{x} | y = l)} &= \frac{1}{2}(\mathbf{x} - \mu_l)^\top \Sigma^{-1}(\mathbf{x} - \mu_l) \\ &\quad - \frac{1}{2}(\mathbf{x} - \mu_k)^\top \Sigma^{-1}(\mathbf{x} - \mu_k) \\ &= -\frac{1}{2}(\mu_k + \mu_l)^\top \Sigma^{-1}(\mu_k - \mu_l) \\ &\quad + \mathbf{x}^\top \Sigma^{-1}(\mu_k - \mu_l)\end{aligned}$$

# Logistic Regression vs LDA

- Both methods use linear log odds
- *LDA*: Calculate weights directly using mean and covariance for each class
  - Easy to calculate
  - Can be more sensitive to outliers
- *Logistic Regression*: Perform maximum likelihood estimation with logistic activation function
  - Harder to solve, but more expressive.

# Linear Discriminant Analysis

LDA



QDA

