

**How can I interpret my ML  
system outputs?**

**Are my interpretations  
reasonable?**

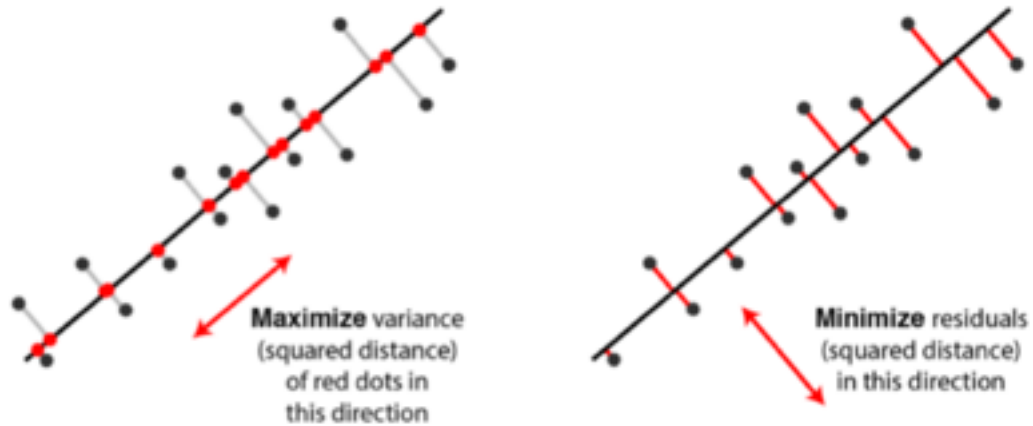
---



# This course so far

---

How can one predict some value, find structure from data?



Two equivalent views of principal component analysis.

# Let's consider LASSO

---

We have come up with some sparse  $w$  for predicting  $y$  from  $x$

# How can we interpret the coefficients? Relative

---

# How can we interpret the coefficients? Relative

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

# How can we interpret the coefficients? Relative

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

# How can we interpret the coefficients? Relative

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

**FALSE**

# How can we interpret the coefficients? Relative

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

**FALSE**

What if features have different scales?

# How can we interpret the coefficients? Relative

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

**FALSE**

What if features have different scales?

- Sqft have a mean of 2500  $\rightarrow w_i$  is

# How can we interpret the coefficients? Relative

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

**FALSE**

What if features have different scales?

- Sqft have a mean of 2500  $\rightarrow w_i$  is
- # baths have a mean of 2  $\rightarrow w_i$  is

# How can we interpret the coefficients? Relative

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

**FALSE**

What if features have different scales?

- Sqft have a mean of 2500  $\rightarrow w_i$  is
- # baths have a mean of 2  $\rightarrow w_i$  is

How to fix?

# How can we interpret the coefficients? Relative

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

**FALSE**

What if features have different scales?

- Sqft have a mean of 2500  $\rightarrow w_i$  is
- # baths have a mean of 2  $\rightarrow w_i$  is

How to fix?

Normalize features to have mean and unit variance (z-score)

# How can we interpret the coefficients? Relative

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

**FALSE**

What if features have different scales?

- Sqft have a mean of 2500  $\rightarrow w_i$  is \$/sqft
- # baths have a mean of 2  $\rightarrow w_i$  is

How to fix?

Normalize features to have mean and unit variance (z-score)

# How can we interpret the coefficients? Relative

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i > \hat{w}_j$  means feature  $i$  is more important than feature  $j$

**FALSE**

What if features have different scales?

- Sqft have a mean of 2500  $\rightarrow w_i$  is \$/sqft
- # baths have a mean of 2  $\rightarrow w_i$  is \$/bath

How to fix?

Normalize features to have mean and unit variance (z-score)

# How can we interpret the coefficients? Zeros

---

# How can we interpret the coefficients? Zeros

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2 + \lambda \|w\|_1 \quad \text{with}$$

normalized data

Claim:  $\hat{w}_i = 0$  means feature  $i$  has no predictive power for  $y$

# How can we interpret the coefficients? Zeros

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2 + \lambda \|w\|_1 \quad \text{with}$$

normalized data

Claim:  $\hat{w}_i = 0$  means feature  $i$  has no predictive power for  $y$

**FALSE**

# How can we interpret the coefficients? Zeros

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2 + \lambda \|w\|_1 \quad \text{with}$$

normalized data

Claim:  $\hat{w}_i = 0$  means feature  $i$  has no predictive power for  $y$

**FALSE**

Feature  $i$  might be correlated with  $y$  but redundant with feature  $j$

# How can we interpret the coefficients? Zeros

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2 + \lambda \|w\|_1 \quad \text{with}$$

normalized data

Claim:  $\hat{w}_i = 0$  means feature  $i$  has no predictive power for  $y$

**FALSE**

Feature  $i$  might be correlated with  $y$  but redundant with feature  $j$   
e.g. area in square feet and area in square meters

# How can we interpret the coefficients?

## Distributional changes

---

# How can we interpret the coefficients?

## Distributional changes

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i = 90,000$  and the  $i$ th feature = #fireplaces. If I add 10 more fireplaces, I can expect to sell my house for \$900,000 more!

# How can we interpret the coefficients?

## Distributional changes

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i = 90,000$  and the  $i$ th feature = #fireplaces. If I add 10 more fireplaces, I can expect to sell my house for \$900,000 more!

**FALSE**

# How can we interpret the coefficients?

## Distributional changes

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i = 90,000$  and the  $i$ th feature = #fireplaces. If I add 10 more fireplaces, I can expect to sell my house for \$900,000 more!

**FALSE**

(1) Distribution shift: we've never seen a house with 10 fireplaces in the training data, so we don't have accurate predictions for this range. It's OOD: out of distribution

# How can we interpret the coefficients?

## Distributional changes

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i = 90,000$  and the  $i$ th feature = #fireplaces. If I add 10 more fireplaces, I can expect to sell my house for \$900,000 more!

**FALSE**

(1) Distribution shift: we've never seen a house with 10 fireplaces in the training data, so we don't have accurate predictions for this range. It's OOD: out of distribution

(2) Correlation  $\neq$  causation

# How can we interpret the coefficients?

## Distributional changes

---

Consider a linear model

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Claim:  $\hat{w}_i = 90,000$  and the  $i$ th feature = #fireplaces. If I add 10 more fireplaces, I can expect to sell my house for \$900,000 more!

**FALSE**

(1) Distribution shift: we've never seen a house with 10 fireplaces in the training data, so we don't have accurate predictions for this range. It's OOD: out of distribution

(2) Correlation  $\neq$  causation

(3) Linear model, non-linear relationship

# Generalization

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Train on hospitals 1, 2, 3 → Test on hospital 4

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Train on hospitals 1, 2, 3 → Test on hospital 4

Claim: The test set performance is always a good indicator of how our model will do if we deploy this model in a new hospital.

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Train on hospitals 1, 2, 3 → Test on hospital 4

Claim: The test set performance is always a good indicator of how our model will do if we deploy this model in a new hospital.

**FALSE**

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Train on hospitals 1, 2, 3 → Test on hospital 4

Claim: The test set performance is always a good indicator of how our model will do if we deploy this model in a new hospital.

FALSE

**Distribution shift:** new hospital could have different equipment, demographics, staff training, location...

# Generalization

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

Claim: The test set performance is always a good indicator of how our model will do if we deploy this model in the **same hospitals**.

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Train on hospitals 1, 2, 3 → Test on  
hospitals 1, 2, 3

Claim: The test set performance is always a good indicator of how our model will do if we deploy this model in the **same hospitals**.

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Train on hospitals 1, 2, 3 → Test on  
hospitals 1, 2, 3

Claim: The test set performance is always a good indicator of how our model will do if we deploy this model in the **same hospitals**.

**FALSE**

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Train on hospitals 1, 2, 3 → Test on hospitals 1, 2, 3

Claim: The test set performance is always a good indicator of how our model will do if we deploy this model in the **same hospitals**.

**FALSE**

**Distribution shift:** temporal shift.

Conditions can change in our hospital by the time we deploy our system! e.g. new training procedures, new staff, new equipment, etc.

# Generalization

Say we've trained a model to interpret medical x-rays using data from a few hospitals. We randomly split the data into train/validation/test splits.

# Train on hospitals 1, 2, 3 → Test on hospitals 1, 2, 3

Claim: The test set performance is always a good indicator of how our model will do if we deploy this model in the **same hospitals**.

**FALSE**

**Distribution shift:** temporal shift.

Conditions can change in our hospital by the time we deploy our system! e.g. new training procedures, new staff, new equipment, etc.

# My particular hobby horse: the world is non-stationary!

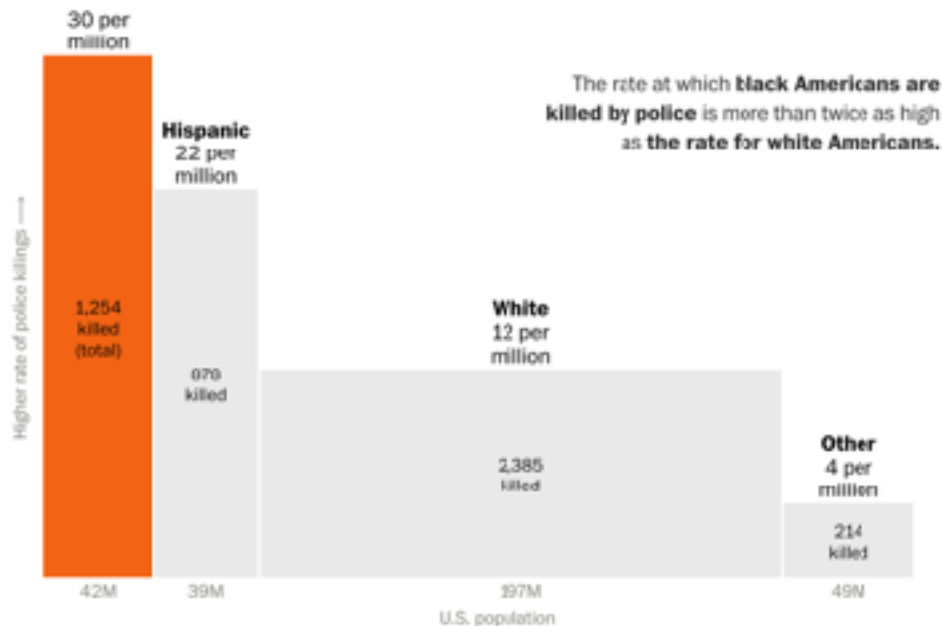
**Now, let's talk about high-stakes decision making**

---

# Some sobering statistics

Many, many people are killed by police every year (10,429 in the past 9 years)

> half of which aren't reported to the FBI (reporting is voluntary)



Source: an excellent Washington Post [interactive database](#)

# Design choices

---

Model Class

Features to collect

Loss function

Regularization

Optimization Method

Constraints

...

# Design choices affect your model

---

Model Class

Features to collect

Loss function

Regularization

Optimization Method

Constraints

...

These will all affect what model you find, and how well and when that structure will generalize

# ... and that can have consequences “beyond” ML

---

Model Class

Features to collect

Loss function

Regularization

Optimization Method

Constraints

...

These will all affect what model you find, and how well and when that structure will generalize

Different models will have different kinds of errors, predictions, and failure modes

# Consequence #1: Good predictions for whom?

---



# The US Criminal Justice System

---

# The US Criminal Justice System

---

This is a description of how things (largely) are and have been, not how they ought to be.

# The US Criminal Justice System

---

This is a description of how things (largely) are and have been, not how they ought to be.

Any use of ML in the criminal justice system is just a tiny part of the larger CJ ecosystem.

# The US Criminal Justice System

---

**This is a description of how things (largely) are and have been, not how they ought to be.**

**Any use of ML in the criminal justice system is just a tiny part of the larger CJ ecosystem.**

**Moreover, “reoffend” or “recidivate” are very, very, very far from precise terms in their use (in this lecture, and in industry).**

# The US Criminal Justice System

---

# The US Criminal Justice System

---

(Lots and lots and lots and lots of possibly biased processes, including racialized choices of where to police, what behaviors are most dangerous, ...)

# The US Criminal Justice System

---

(Lots and lots and lots and lots of possibly biased processes, including racialized choices of where to police, what behaviors are most dangerous, ...)

-> A person is charged with a crime

# The US Criminal Justice System

---

**(Lots and lots and lots and lots of possibly biased processes, including racialized choices of where to police, what behaviors are most dangerous, ...)**

-> A person is charged with a crime

-> A judge determines whether to detain them or release them on bail until their trial

# The US Criminal Justice System

---

(Lots and lots and lots and lots of possibly biased processes, including racialized choices of where to police, what behaviors are most dangerous, ...)

-> A person is charged with a crime

-> A judge determines whether to detain them or release them on bail until their trial

-> (Many, many months or years pass)

# The US Criminal Justice System

---

**(Lots and lots and lots and lots of possibly biased processes, including racialized choices of where to police, what behaviors are most dangerous, ...)**

-> A person is charged with a crime

-> A judge determines whether to detain them or release them on bail until their trial

-> (Many, many months or years pass)

-> The trial determines whether the justice system considers a person guilty or innocent of charges

# The US Criminal Justice System

---

**(Lots and lots and lots and lots of possibly biased processes, including racialized choices of where to police, what behaviors are most dangerous, ...)**

-> A person is charged with a crime

-> A judge determines whether to detain them or release them on bail until their trial

-> (Many, many months or years pass)

-> The trial determines whether the justice system considers a person guilty or innocent of charges

-> A sentencing hearing determines the length of prison time, fines, ...

# The US Criminal Justice System

---

**(Lots and lots and lots and lots of possibly biased processes, including racialized choices of where to police, what behaviors are most dangerous, ...)**

-> A person is charged with a crime

-> A judge determines whether to detain them or release them on bail until their trial

-> (Many, many months or years pass)

-> The trial determines whether the justice system considers a person guilty or innocent of charges

-> A sentencing hearing determines the length of prison time, fines, ...

-> With “good behavior”, people in prison can be considered for parole

# What does the CJ process have to do with ML?

---

(Lots and lots and lots and lots of possibly biased processes, including racialized choices of where to police, what behaviors are most dangerous, ...)

-> A person is charged with a crime

-> A judge determines whether to detain them or release them on bail until their trial

-> (Many, many months or years pass)

-> The trial determines whether the justice system considers a person guilty or innocent of charges

-> A sentencing hearing determines the length of prison time, fines, ...

-> With “good behavior”, people in prison can be considered for parole

# What does the CJ process have to do with ML?

---

Predictive policing: determine (statistically) where to send police

- > A person is charged with a crime
- > A judge determines whether to detain them or release them on bail until their trial
- > (Many, many months or years pass)
- > The trial determines whether the justice system considers a person guilty or innocent of charges
- > A sentencing hearing determines the length of prison time, fines, ...
- > With “good behavior”, people in prison can be considered for parole

# What does the CJ process have to do with ML?

---

Predictive policing: determine (statistically) where to send police

-> A person is charged with a crime

Predict probability of reappearing, being charged with another crime if on parole, ...

-> (Many, many months or years pass)

-> The trial determines whether the justice system considers a person guilty or innocent of charges

-> A sentencing hearing determines the length of prison time, fines, ...

-> With “good behavior”, people in prison can be considered for parole

# What does the CJ process have to do with ML?

---

Predictive policing: determine (statistically) where to send police

-> A person is charged with a crime

Predict probability of reappearing, being charged with another crime if on parole, ...

-> (Many, many months or years pass)

-> The trial determines whether the justice system considers a person guilty or innocent of charges

Predict risk of being charged with another crime, use in sentencing

-> With “good behavior”, people in prison can be considered for parole

# What does the CJ process have to do with ML?

---

Predictive policing: determine (statistically) where to send police

-> A person is charged with a crime

Predict probability of reappearing, being charged with another crime if on parole, ...

-> (Many, many months or years pass)

-> The trial determines whether the justice system considers a person guilty or innocent of charges

Predict risk of being charged with another crime, use in sentencing

Predict risk of being charged with another crime, use in parole decisions

# The US Criminal Justice System

---

# The US Criminal Justice System

---

**“reoffend” or “recidivate” are very, very, very far from precise terms in their use (in this lecture, and in industry).**

# The US Criminal Justice System

---

**“reoffend” or “recidivate” are very, very, very far from precise terms in their use (in this lecture, and in industry).**

**Largely, it corresponds to either being:**

# The US Criminal Justice System

---

“reoffend” or “recidivate” are very, very, very far from precise terms in their use (in this lecture, and in industry).

Largely, it corresponds to either being:

- arrested

# The US Criminal Justice System

---

“reoffend” or “recidivate” are very, very, very far from precise terms in their use (in this lecture, and in industry).

Largely, it corresponds to either being:

- arrested
- charged

# The US Criminal Justice System

---

“reoffend” or “recidivate” are very, very, very far from precise terms in their use (in this lecture, and in industry).

Largely, it corresponds to either being:

- arrested
- charged
- or being found guilty of a crime

# The US Criminal Justice System

---

“reoffend” or “recidivate” are very, very, very far from precise terms in their use (in this lecture, and in industry).

Largely, it corresponds to either being:

- arrested
- charged
- or being found guilty of a crime

usually one which is violent.

**Can such predictions be fair with respect to race?**

---

# Can such predictions be fair with respect to race?

---

Attempt #1: don't use race as a feature

Issue: race is strongly correlated with other features

Attempt #2: Remove any feature correlated with race

Issue: (nearly) all features correlated with race

“Fairness through unawareness”

What exactly is the goal of this?

# Can such predictions be fair with respect to race?

---

Attempt #1: don't use race as a feature

Issue: race is strongly correlated with other features

Attempt #2: Remove any feature correlated with race

Issue: (nearly) all features correlated with race

“Fairness through unawareness”

What exactly is the goal of this?

To accurately predict which people are most likely commit a crime if released, so they can be released

# Can such predictions be fair with respect to race?

---

Attempt #1: don't use race as a feature

Issue: race is strongly correlated with other features

Attempt #2: Remove any feature correlated with race

Issue: (nearly) all features correlated with race

“Fairness through unawareness”

What exactly is the goal of this?

To accurately predict which people are most likely commit a crime if released, so they can be released

# Can such predictions be fair with respect to race?

---

Attempt #1: don't use race as a feature

Issue: race is strongly correlated with other features

Attempt #2: Remove any feature correlated with race

Issue: (nearly) all features correlated with race

“Fairness through unawareness”

What exactly is the goal of this?

To accurately predict which people are most likely commit a crime if released, so they can be released

... where these predictions should be accurate? similar?  
for people of all races.

# Can such predictions be fair with respect to race?

---

Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale

Attempt #1: don't use race as a feature

Issue: race is strongly correlated with other features

Attempt #2: Remove any feature correlated with race

Issue: (nearly) all features correlated with race

“Fairness through unawareness”

What exactly is the goal of this?

To accurately predict which people are most likely commit a crime if released, so they can be released

... where these predictions should be accurate? similar?  
for people of all races.

# Can such predictions be fair with respect to race?

Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale

Attempt #1: don't use race as a feature

Issue: race is strongly correlated with other features

Attempt #2: Remove any feature correlated with race

Issue: (nearly) all features correlated with race

“Fairness through unawareness”

What exactly is the goal of this?

To accurately predict which people are most likely commit a crime if released, so they can be released

Concern: not all errors are equally costly.

... where these predictions should be accurate? similar?  
for people of all races.

# Can such predictions be fair with respect to race?

---



# Can such predictions be fair with respect to race?

---



Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale

# Can such predictions be fair with respect to race?

---

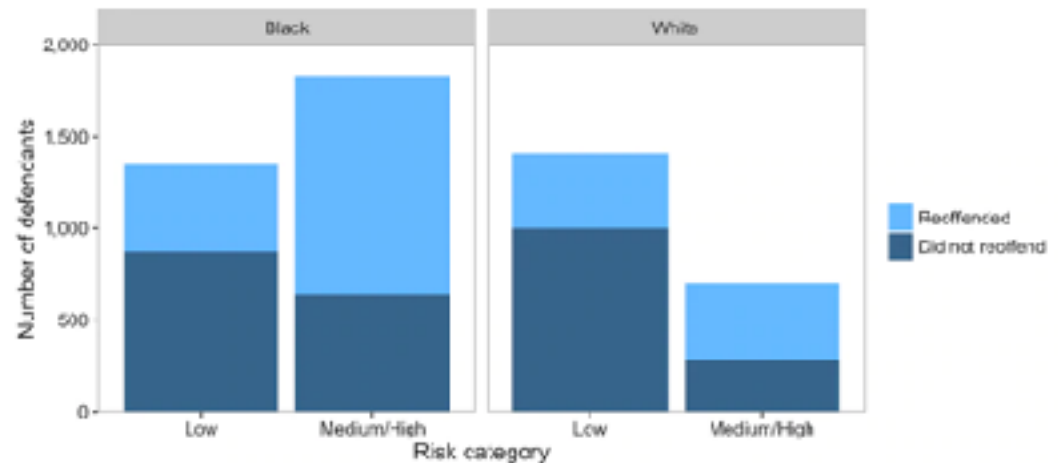


Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale

# Can such predictions be fair with respect to race?



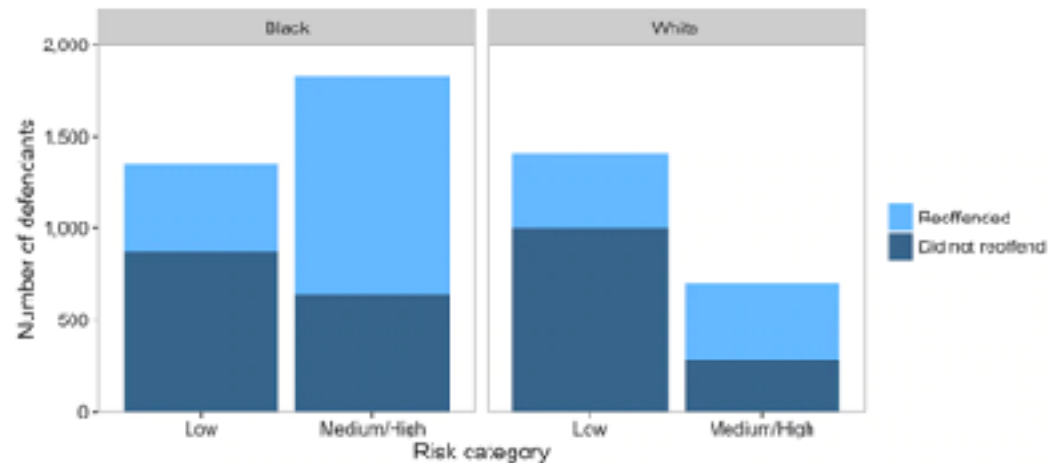
Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale



# Can such predictions be fair with respect to race?



Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale

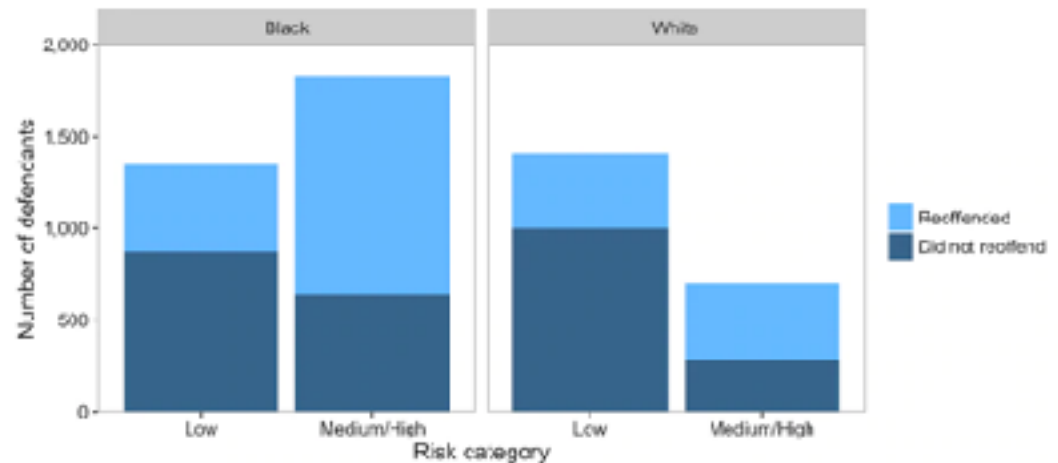


$\mathbb{P}[\text{does not reoffend} | \text{predicted low risk, white}] <$   
 $\mathbb{P}[\text{does not reoffend} | \text{predicted low risk, black}]$

# Can such predictions be fair with respect to race?



Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale



$\mathbb{P}[\text{does not reoffend} | \text{predicted low risk, white}] <$   
 $\mathbb{P}[\text{does not reoffend} | \text{predicted low risk, black}]$

False positive rate for black defendants higher than for white defendants.

# Can such predictions be fair with respect to race?

---



Calibration of scores w.r.t. race

# Can such predictions be fair with respect to race?

---



Calibration of scores w.r.t. race

[Source: Washington Post](#)

# Can such predictions be fair with respect to race?

---



Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale

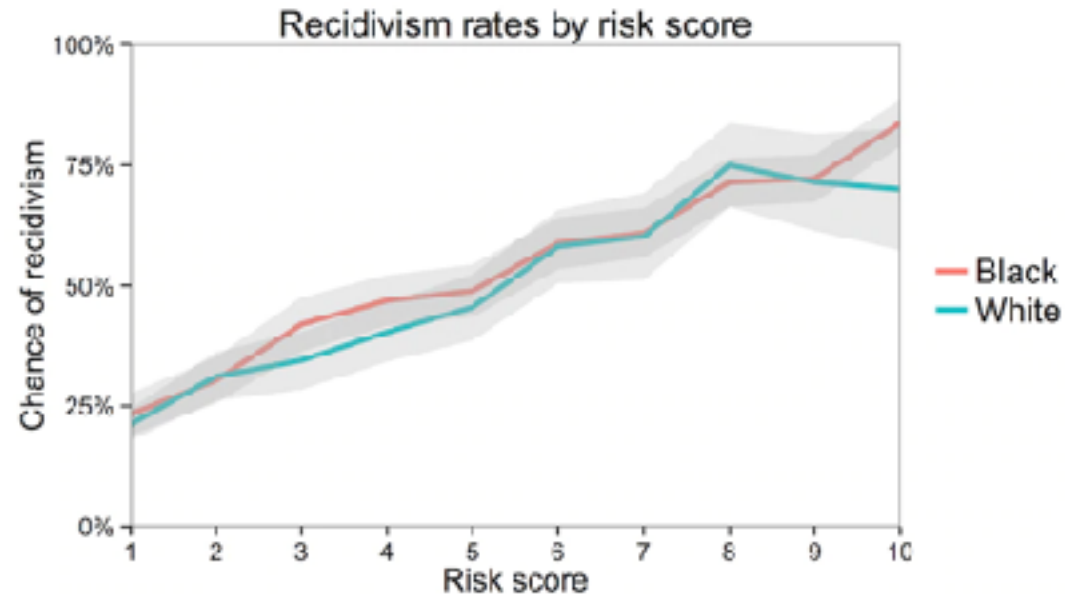
Calibration of scores w.r.t. race

[Source: Washington Post](#)

# Can such predictions be fair with respect to race?



Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale



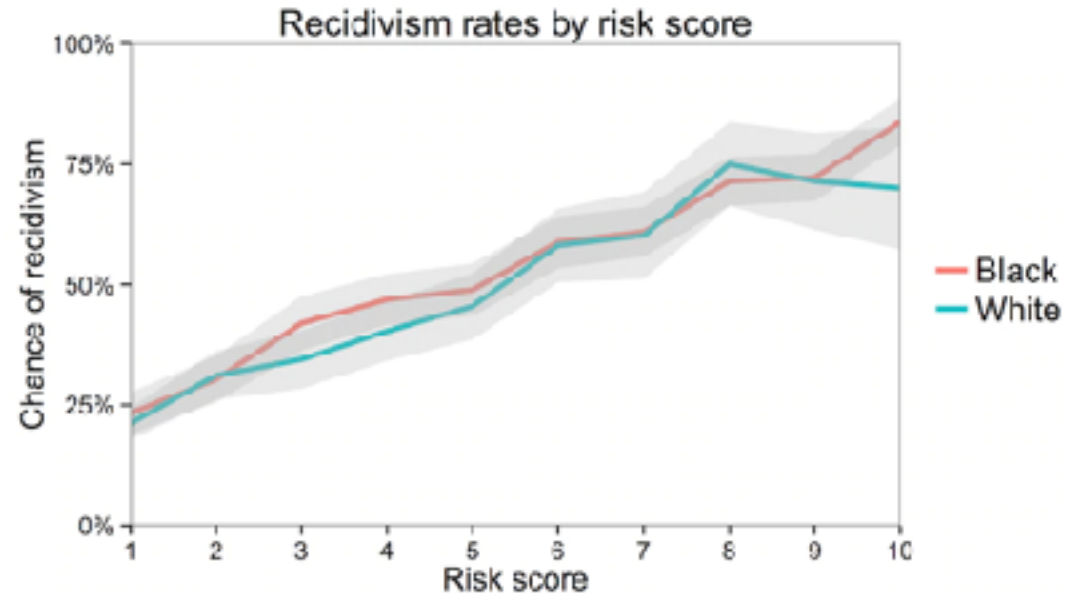
Calibration of scores w.r.t. race

[Source: Washington Post](#)

# Can such predictions be fair with respect to race?



Pretrial release risk scale: 1-10  
General recidivism scale  
Violent recidivism scale



$$\mathbb{P}[\text{reoffend} | \text{prediction, white}] \approx \mathbb{P}[\text{reoffend} | \text{prediction, black}]$$

Calibration of scores w.r.t. race

# Can any predictions be fair with respect to race?

---



$\mathbb{P}[\text{does not reoffend} | \text{predicted low risk, white}] <$   
 $\mathbb{P}[\text{does not reoffend} | \text{predicted low risk, black}]$

while

$\mathbb{P}[\text{reoffend} | \text{prediction, white}] \approx \mathbb{P}[\text{reoffend} | \text{prediction, black}]$

How is this possible?

Is this avoidable?

# Can any predictions be fair with respect to race?

---

False positives approximately equal and calibration are mutually exclusive, unless we have perfect predictions or the rate of what we are predicting is equal in every racial group.

... and this is true for more settings than just criminal justice.

Lending

Advertising

...

**We are making a choice about how we allocate predictions.**

Classification! Logistic  
“Regression” (linear  
classification), oh my...

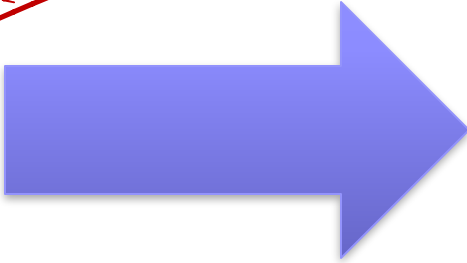
---

# **Thus far, regression:**

---

**predict a continuous value given some inputs**

# Weather prediction revisited



# Recent Kaggle competitions

The screenshot displays the Kaggle website interface. On the left is a navigation sidebar with the Kaggle logo and menu items: Home, Competitions (selected), Benchmarks, Game Assets, Data Hub, and More. The main content area features a search bar for competitions, a 'Sign In' button, and a 'Register' button. Below the search bar is a 'Completed' filter and a 'Results' section. The results are presented in a table with columns for competition details and 'Recently Completed' status.

| Competition  | Recently Completed |
|--|--------------------|
| <b>Predicting Inflation Need</b><br>Program on Series - Season 4 Episode 4<br>Playground - 4285 Teams - 44 minutes ago   | 5wag<br>--         |
| <b>Record aJLUC - Scientific Image Forgery Detection</b><br>Develop methods that can accurately detect and segment copy-move forgeries across biomedical research images.<br>Research - Code Competition - 1188 Teams - 8 days ago | \$25,000<br>--     |
| <b>Measuring Progress Toward AGI - Cognitive Abilities</b><br>Design high-quality benchmarks that go beyond recall to evaluate how frontier models truly reason, act, and judge.<br>Featured - 87 Teams - 18 days ago              | \$200,000<br>--    |
| <b>AI Mathematical Olympiad - Progress Prize 3</b><br>Solve international-level math challenges using artificial intelligence models.<br>Featured - Code Competition - 4708 Teams - 10 days ago                                    | \$2,207,522<br>--  |
| <b>Meta Machine Learning Memia 2025</b><br>Present the 150th ICLR & Funded by Teamwork<br>Featured - 1482 Teams - 29 days ago  | \$50,000<br>--     |

# Classification

---

- Learn  $f: X \rightarrow Y$ 
  - $X$  - features
  - $Y$  - target classes
- Loss Function
- Expected loss of  $f$ :
  
- Suppose you knew  $P(Y|X)$  exactly, how should you classify?
  - Bayes-Optimal classifier:

# Classification

- Learn  $f: X \rightarrow Y$ 
  - $X$  - features
  - $Y$  - target classes

- **Loss Function**

$$\ell(f(x), y) = \mathbf{1}\{f(x) \neq y\}$$

- **Expected loss of  $f$ :**

$$\mathbb{E}_{XY}[\mathbf{1}\{f(X) \neq Y\}] = \mathbb{E}_X[\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x]]$$

$$\begin{aligned}\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x] &= \sum_i P(Y = i|X = x)\mathbf{1}\{f(x) \neq i\} = \sum_{i \neq f(x)} P(Y = i|X = x) \\ &= 1 - P(Y = f(x)|X = x)\end{aligned}$$

- **Suppose you knew  $P(Y|X)$  exactly, how should you classify?**
  - **Bayes-Optimal classifier:**

$$f(x) = \arg \max_y \mathbb{P}(Y = y|X = x)$$

# Binary Classification

- Learn  $f: X \rightarrow Y$

- $X$  - features

- $Y$  - target classes

$$Y \in \{0, 1\}$$

- Loss Function

$$\ell(f(x), y) = \mathbf{1}\{f(x) \neq y\}$$

- Expected loss of  $f$ :

$$\mathbb{E}_{XY}[\mathbf{1}\{f(X) \neq Y\}] = \mathbb{E}_X[\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x]]$$

$$\begin{aligned}\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x] &= \sum_i P(Y = i|X = x)\mathbf{1}\{f(x) \neq i\} = \sum_{i \neq f(x)} P(Y = i|X = x) \\ &= 1 - P(Y = f(x)|X = x)\end{aligned}$$

- Suppose you knew  $P(Y|X)$  exactly, how should you classify?

- Bayes-Optimal classifier:

$$f(x) = \arg \max_y \mathbb{P}(Y = y|X = x)$$

# Bayes Optimal Binary Classifier

---

$$Y \in \{0, 1\}$$

- Suppose you knew  $P(Y|X)$  exactly, how should you classify?
- Bayes-Optimal classifier:

$$f(x) = \arg \max_y \mathbb{P}(Y = y | X = x)$$

- Suppose we don't know  $P(Y|X)$ , but have  $n$  iid examples

$$\{(x_i, y_i)\}_{i=1}^n$$

- What is a natural estimator for  $P(Y | X)$ ?

# Bayes Optimal Binary Classifier

- Suppose we don't know  $P(Y|X)$ , but have  $n$  iid examples

$$\{(x_i, y_i)\}_{i=1}^n \quad Y \in \{0, 1\}$$

- What is a natural estimator for  $P(Y | X)$ ?

Fix some  $\tilde{x} \in X$

Suppose  $x_i = \tilde{x}$  for  $m \leq n$  samples

What is a natural estimator for  $\theta_* := \mathbb{P}(Y = 1 | X = \tilde{x})$ ?

If  $k$  of the  $m$  labels are equal to  $Y = 1$  then

# Bayes Optimal Binary Classifier

- Suppose we don't know  $P(Y|X)$ , but have  $n$  iid examples

$$\{(x_i, y_i)\}_{i=1}^n$$

$$Y \in \{0, 1\}$$

- What is a natural estimator for  $\operatorname{argmax}_y P(Y = y | X)$ ?

If  $X = \{0, 1\}^d$ , or is generally discrete

$$\hat{f}(x) = \operatorname{argmax}_{y \in \{0,1\}} \frac{\sum_{i=1}^n \mathbf{1}[x_i = x, y_i = y]}{\sum_{i=1}^n \mathbf{1}[x_i = x]}$$

Issues?

# Bayes Optimal Binary Classifier

- What is a natural estimator for  $\operatorname{argmax}_y \mathbb{P}(Y = y | X)$ ?

If  $X = \{0, 1\}^d$ , or is generally discrete  $Y \in \{0, 1\}$

$$\hat{f}(x) = \operatorname{argmax}_{y \in \{0, 1\}} \frac{\sum_{i=1}^n \mathbf{1}[\mathbf{x}_i = \mathbf{x}, \mathbf{y}_i = y]}{\sum_{i=1}^n \mathbf{1}[\mathbf{x}_i = \mathbf{x}]}$$

Issues?

$2^d$  possible inputs, for small  $d$  requires huge  $n$

To make predictions for unseen inputs ( $x$ s),

need a **general** model for  $\mathbb{P}(Y = 1 | X = x)$

# Process

---

Decide on a **model**

Find the function which fits the data best

**Choose a loss function**

**Pick the function which minimizes loss  
on data**

Use function to make prediction on new  
examples

# Decide on a model, Binary Classification

---

To make predictions for unseen inputs ( $x$ s),

need a **general** model for  $\mathbb{P}(Y = 1|X = x)$

- **What about standard linear regression model?**

- **Need to map real values to  $[0,1]$** 
  - **We call such maps “link functions”**

# Logistic Regression

Actually classification, not regression :)

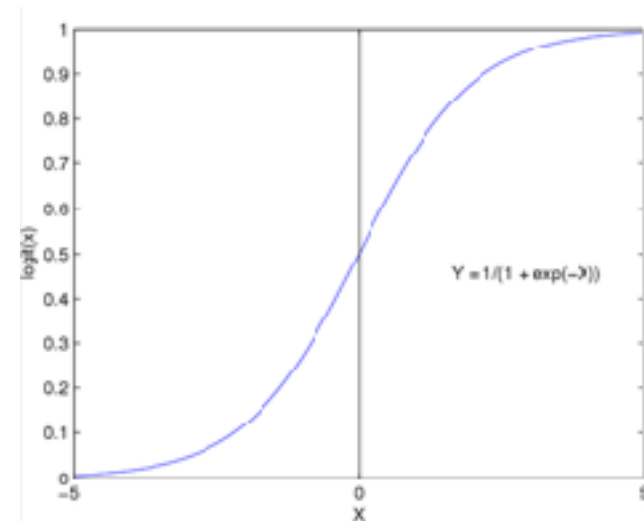
Learn  $\mathbb{P}(Y = 1|X = x)$  using  $\sigma(w^T x)$ , for link function  $\sigma =$

Logistic function(or Sigmoid):

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

$$\mathbb{P}[Y = 1|X = x, w] = \sigma(w^T x) = \frac{1}{1 + \exp(-w^T x)}$$

$$\begin{aligned}\mathbb{P}[Y = 0|X = x, w] &= 1 - \sigma(w^T x) = \frac{\exp(-w^T x)}{1 + \exp(-w^T x)} \\ &= \frac{1}{1 + \exp(w^T x)}\end{aligned}$$

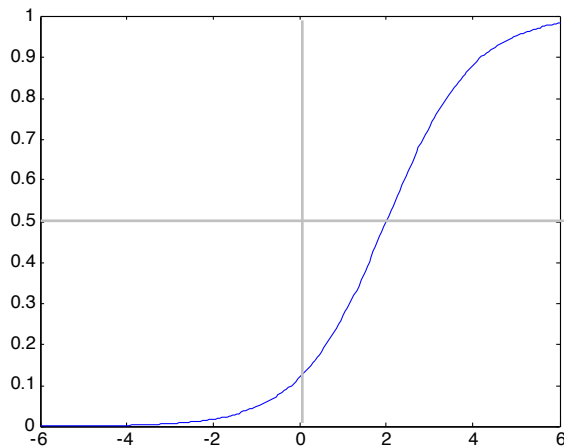


Features can be discrete or continuous!

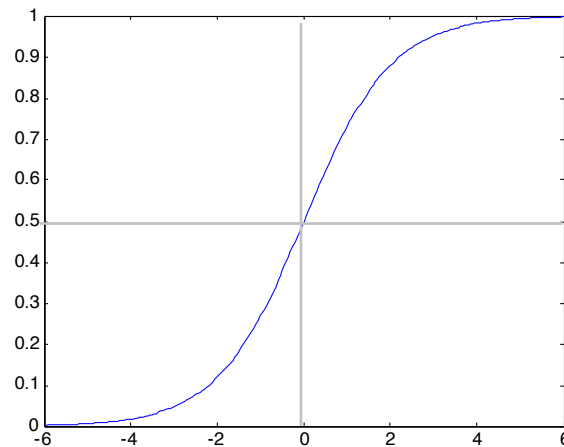
# Understanding the sigmoid

$$\sigma(w_0 + \sum_k w_k x_k) = \frac{1}{1 + e^{w_0 + \sum_k w_k x_k}}$$

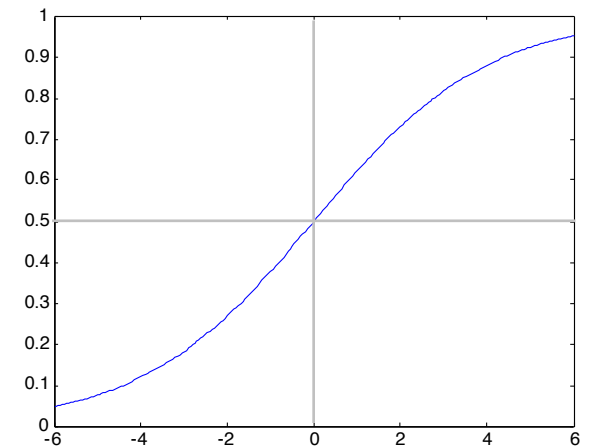
$$w_0 = -2, w_1 = -1$$



$$w_0 = 0, w_1 = -1$$



$$w_0 = 0, w_1 = -0.5$$



# Sigmoid for binary classes

---

$$\mathbb{P}(Y = 0|w, X) = \frac{1}{1 + \exp(w_0 + \sum_k w_k X_k)}$$

$$\mathbb{P}(Y = 1|w, X) = 1 - \mathbb{P}(Y = 0|w, X) = \frac{\exp(w_0 + \sum_k w_k X_k)}{1 + \exp(w_0 + \sum_k w_k X_k)}$$

$$\frac{\mathbb{P}(Y = 1|w, X)}{\mathbb{P}(Y = 0|w, X)} =$$

# Sigmoid for binary classes

$$\mathbb{P}(Y = 0|w, X) = \frac{1}{1 + \exp(w_0 + \sum_k w_k X_k)}$$

$$\mathbb{P}(Y = 1|w, X) = 1 - \mathbb{P}(Y = 0|w, X) = \frac{\exp(w_0 + \sum_k w_k X_k)}{1 + \exp(w_0 + \sum_k w_k X_k)}$$

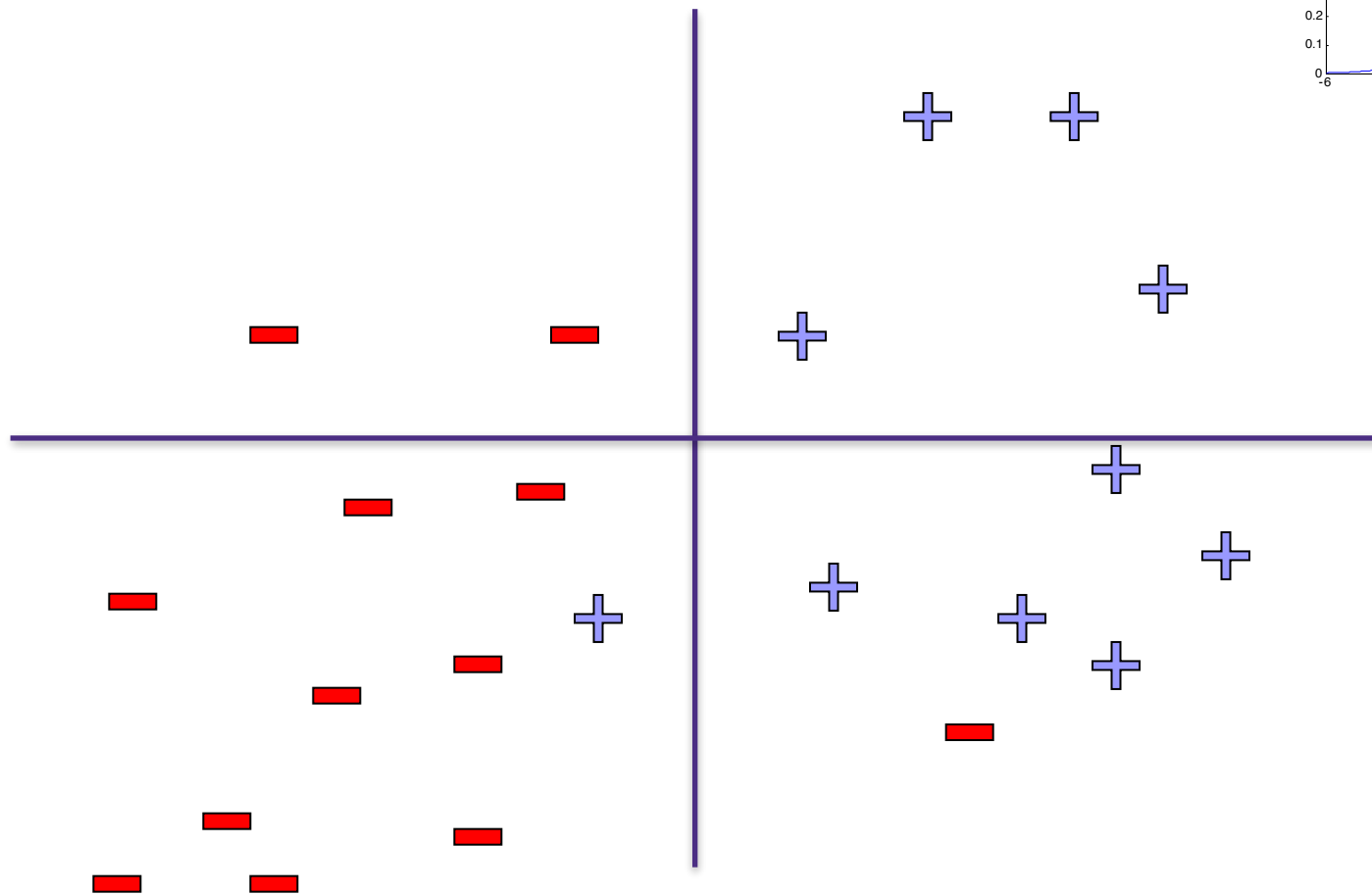
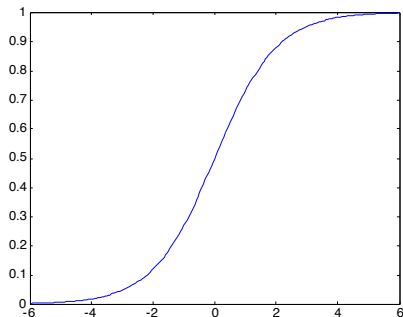
$$\frac{\mathbb{P}(Y = 1|w, X)}{\mathbb{P}(Y = 0|w, X)} = \exp(w_0 + \sum_k w_k X_k)$$

$$\log \frac{\mathbb{P}(Y = 1|w, X)}{\mathbb{P}(Y = 0|w, X)} = w_0 + \sum_k w_k X_k$$

**Linear Decision Rule!**

# Logistic Regression – a Linear classifier

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



$$\ln \frac{P(Y = 0|X)}{P(Y = 1|X)} = w_0 + \sum_i w_i X_i$$

# Process

---

Decide on a **model**

Find the function which fits the data best

**Choose a loss function**

**Pick the function which minimizes loss  
on data**

Use function to make prediction on new  
examples

# Loss function: Conditional Likelihood

- **Have a bunch of iid data:**  $\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$

$$P(Y = -1|x, w) = \frac{1}{1 + \exp(w^T x)}$$

$$P(Y = 1|x, w) = \frac{\exp(w^T x)}{1 + \exp(w^T x)}$$

- **This is equivalent to:**

$$P(Y = y|x, w) = \frac{1}{1 + \exp(-y w^T x)}$$

- **So we can compute the maximum likelihood estimator:**

$$\hat{w}_{MLE} = \arg \max_w \prod_{i=1}^n P(y_i|x_i, w)$$

# Loss function: Conditional Likelihood

- **Have a bunch of iid data:**  $\{(x_i, y_i)\}_{i=1}^n$   $x_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$

$$\begin{aligned}\hat{w}_{MLE} &= \arg \max_w \prod_{i=1}^n P(y_i | x_i, w) & P(Y = y | x, w) &= \frac{1}{1 + \exp(-y w^T x)} \\ &= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w))\end{aligned}$$

# Loss function: Conditional Likelihood

- **Have a bunch of iid data:**  $\{(x_i, y_i)\}_{i=1}^n$   $x_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$

$$\begin{aligned}\hat{w}_{MLE} &= \arg \max_w \prod_{i=1}^n P(y_i | x_i, w) & P(Y = y | x, w) &= \frac{1}{1 + \exp(-y w^T x)} \\ &= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w))\end{aligned}$$

Logistic Loss:  $\ell_i(w) = \log(1 + \exp(-y_i x_i^T w))$

Squared error Loss:  $\ell_i(w) = (y_i - x_i^T w)^2$  (MLE for Gaussian noise)

# Process

---

Decide on a **model**

Find the function which fits the data best

**Choose a loss function**

**Pick the function which minimizes loss  
on data**

Use function to make prediction on new  
examples

# Loss function: Conditional Likelihood

- **Have a bunch of iid data:**  $\{(x_i, y_i)\}_{i=1}^n$   $x_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$

$$\begin{aligned}\hat{w}_{MLE} &= \arg \max_w \prod_{i=1}^n P(y_i | x_i, w) & P(Y = y | x, w) &= \frac{1}{1 + \exp(-y w^T x)} \\ &= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w)) = J(w)\end{aligned}$$

What does  $J(w)$  look like? Is it convex?

# Loss function: Conditional Likelihood

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$$

$$\begin{aligned} \hat{w}_{MLE} &= \arg \max_w \prod_{i=1}^n P(y_i | x_i, w) & P(Y = y | x, w) &= \frac{1}{1 + \exp(-y w^T x)} \\ &= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w)) = J(w) \end{aligned}$$

**Good news:**  $J(\mathbf{w})$  is convex function of  $\mathbf{w}$ , no local optima problems

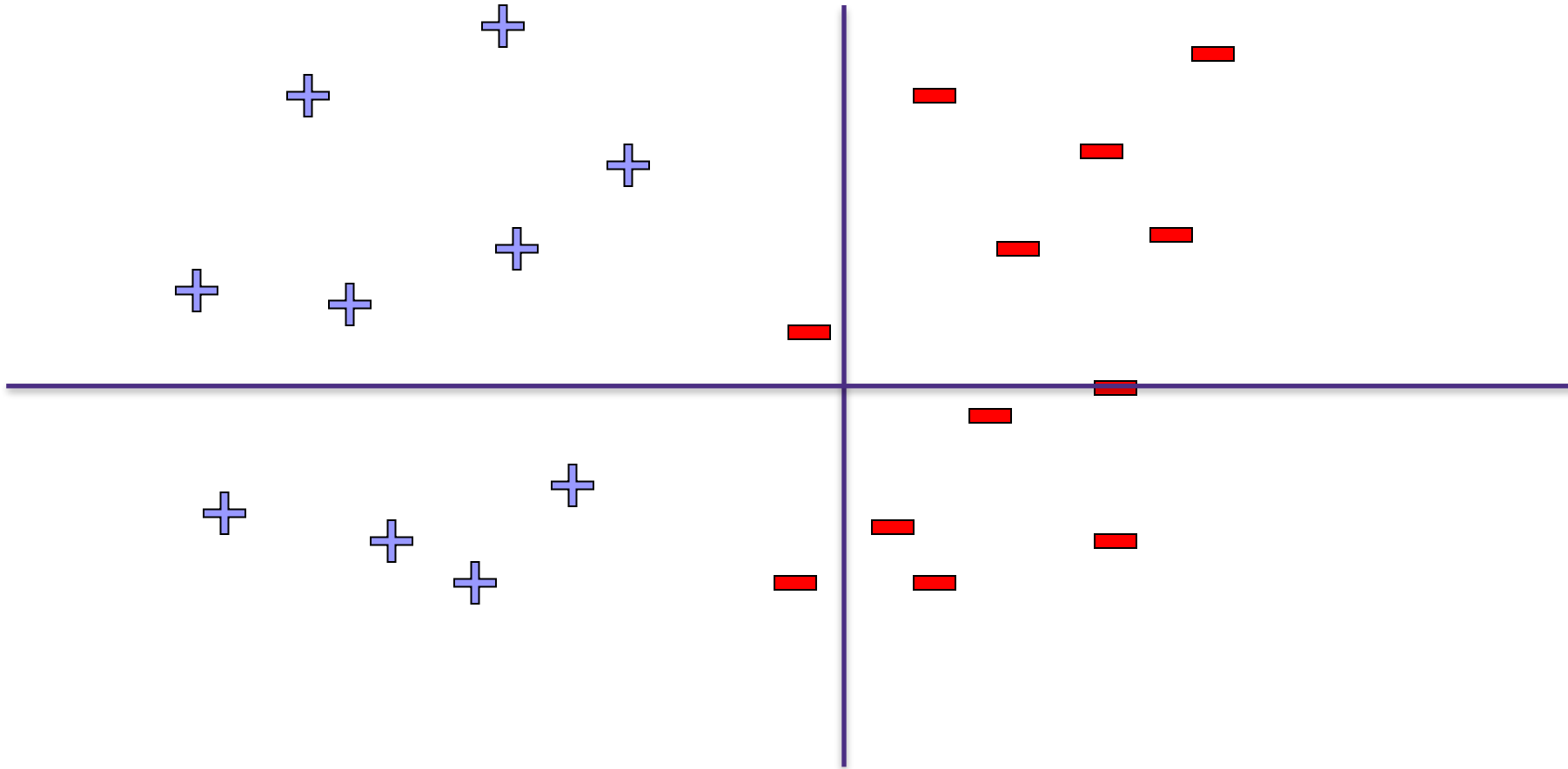
**Bad news:** no closed-form solution to maximize  $J(\mathbf{w})$

**Good news:** convex functions easy to optimize

# Linear Separability

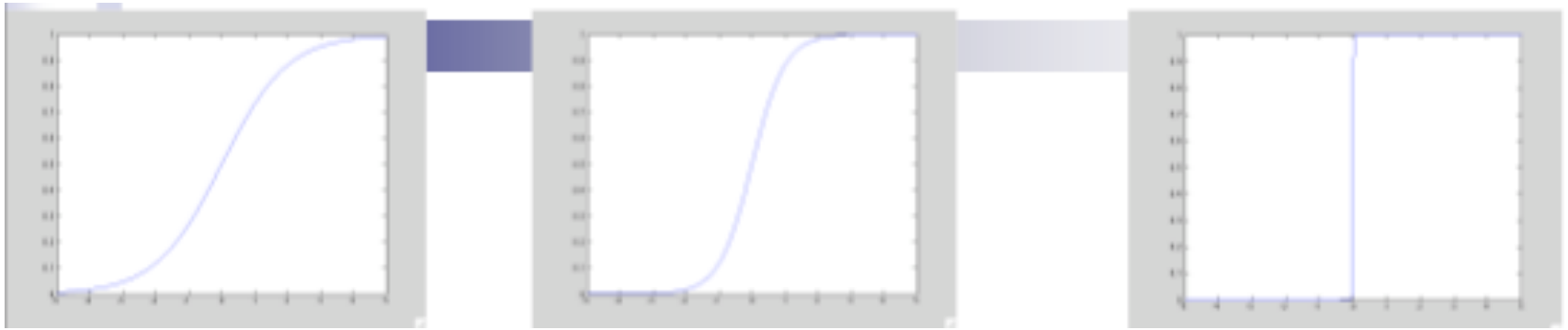
$$\arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w))$$

When is this loss small?



# Large parameters $\rightarrow$ Overfitting

When data is linearly separable, weights  $\Rightarrow \infty$



$$\frac{1}{1 + e^{-x}}$$

$$\frac{1}{1 + e^{-2x}}$$

$$\frac{1}{1 + e^{-100x}}$$

Overfitting

Penalize high weights to prevent overfitting?

# Regularized Conditional Log Likelihood

---

Add a penalty to avoid high weights/overfitting?:

$$\arg \min_{w,b} \sum_{i=1}^n \log (1 + \exp(-y_i (x_i^T w + b))) + \lambda \|w\|_2^2$$

Be sure to not regularize the offset  $b$ !

# Multi-class classification

---

# Multi-class classification

---

- So far: binary  $y \in \{-1, +1\}$

# Multi-class classification

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$

# Multi-class classification

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

# Multi-class classification

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



# Multi-class classification

---

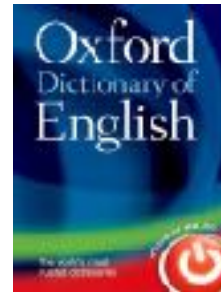
- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



# Multi-class classification

---

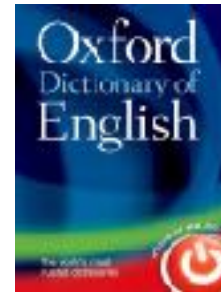
- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



# Multi-class classification

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



- A k-class classifier predicts  $y$  given  $x$

# How to handle categorical labels

---

# How to handle categorical labels

---

- So far: binary  $y \in \{-1, +1\}$

# How to handle categorical labels

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$

# How to handle categorical labels

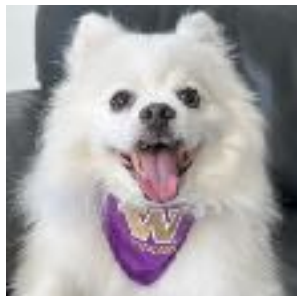
---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

# How to handle categorical labels

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



# How to handle categorical labels

---

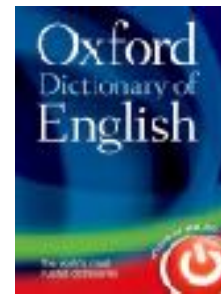
- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



# How to handle categorical labels

---

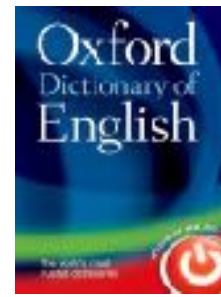
- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



# How to handle categorical labels

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

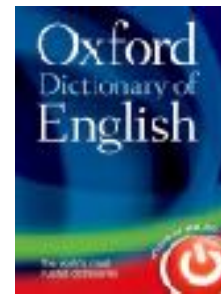


- A k-class classifier predicts  $y$  given  $x$

# How to handle categorical labels

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

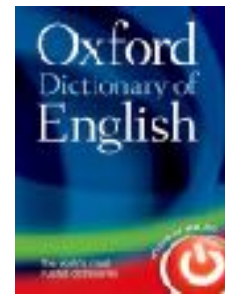


- A k-class classifier predicts  $y$  given  $x$

Ideas for how to use logistic regression for multi class prediction:

# How to handle categorical labels

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



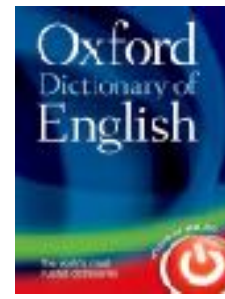
- A k-class classifier predicts  $y$  given  $x$

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$

# How to handle categorical labels

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels



- A k-class classifier predicts  $y$  given  $x$

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$

# One-HOT encoding labels + softmax

---

# One-HOT encoding labels + softmax

---

- So far: binary  $y \in \{-1, +1\}$

# One-HOT encoding labels + softmax

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$

# One-HOT encoding labels + softmax

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

# One-HOT encoding labels + softmax

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

# One-HOT encoding labels + softmax

---

- So far: binary  $y \in \{-1, +1\}$

- In general:  $y \in \{c_1, c_2, \dots, c_k\}$

- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$

# One-HOT encoding labels + softmax

---

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

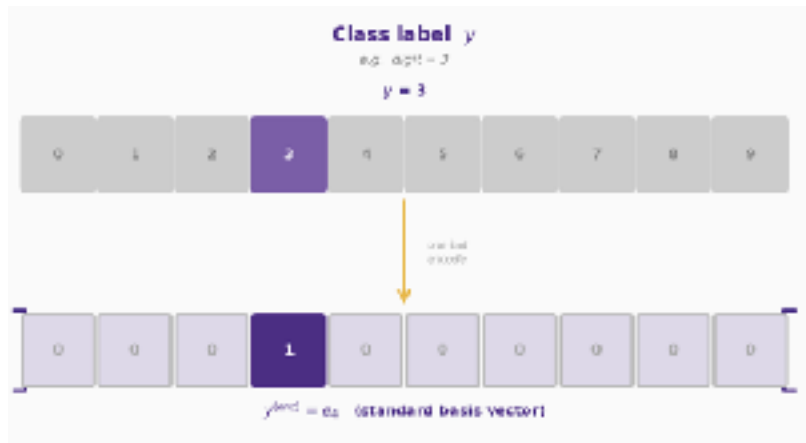
- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$

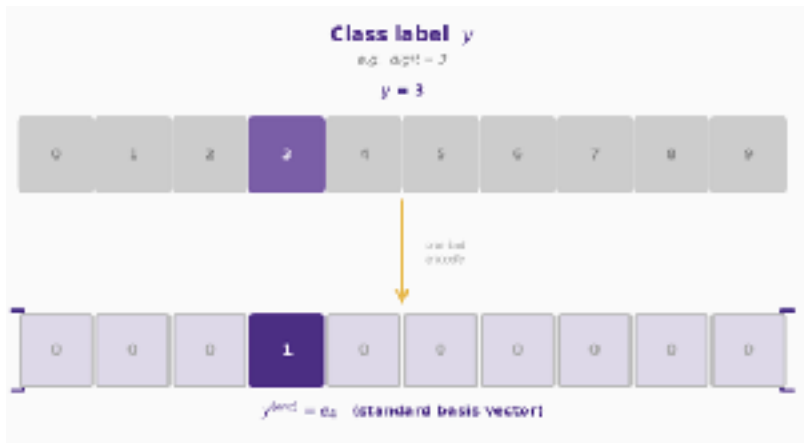


# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



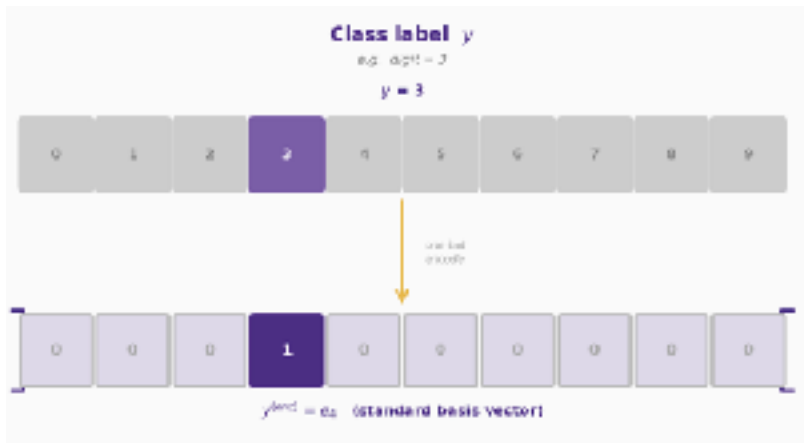
$$\mathbb{P}[c_1 | x]$$

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



$$\mathbb{P}[c_1 | x]$$

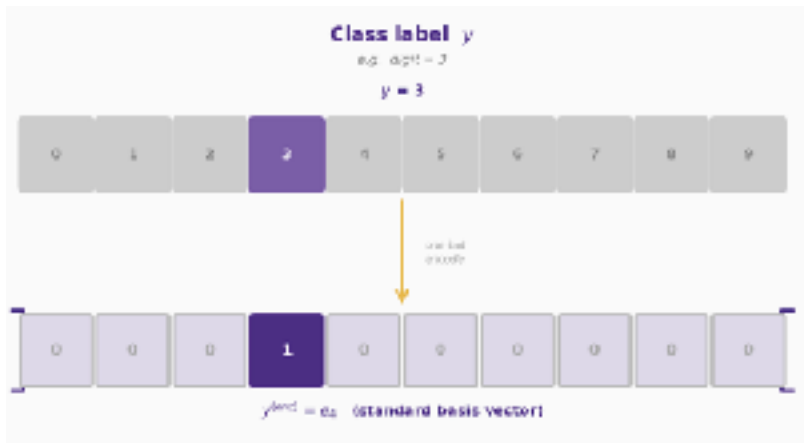
$$\mathbb{P}[c_2 | x]$$

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



$$\mathbb{P}[c_1 | x]$$

$$\mathbb{P}[c_2 | x]$$

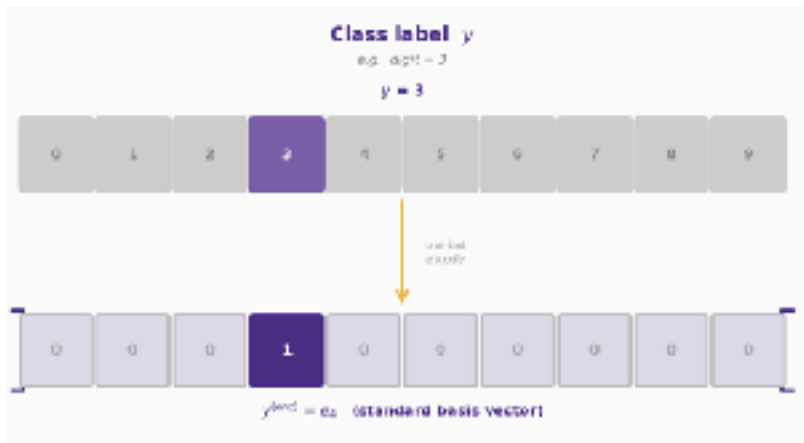
⋮

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



$$\mathbb{P}[c_1 | x]$$

$$\mathbb{P}[c_2 | x]$$

⋮

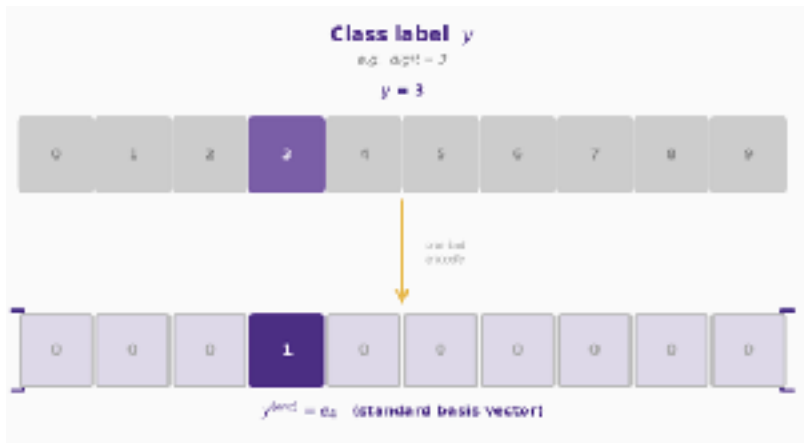
$$\mathbb{P}[c_k | x]$$

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



$$w_1 \cdot x$$

$$\mathbb{P}[c_1 | x]$$

$$\mathbb{P}[c_2 | x]$$

⋮

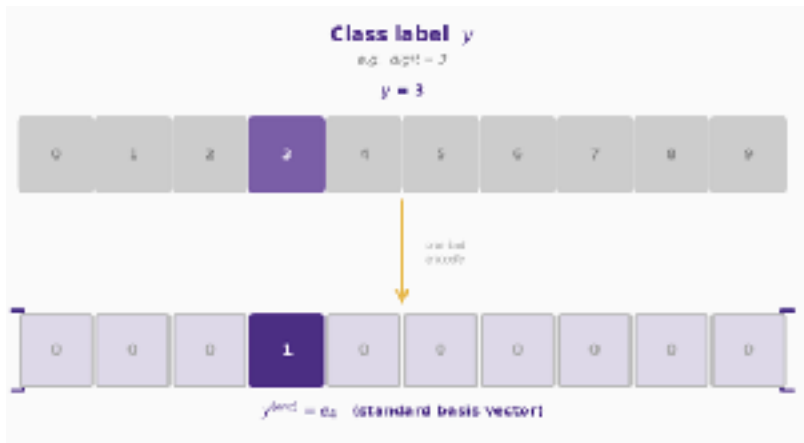
$$\mathbb{P}[c_k | x]$$

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



$$w_1 \cdot x$$

$$\mathbb{P}[c_1 | x]$$

$$w_2 \cdot x$$

$$\mathbb{P}[c_2 | x]$$

⋮

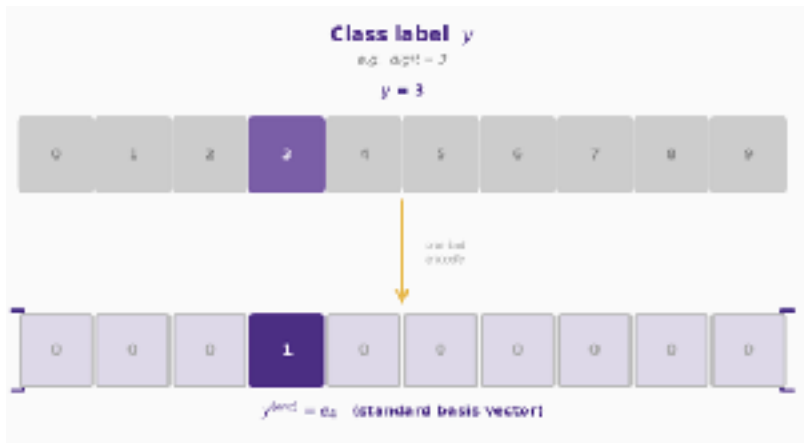
$$\mathbb{P}[c_k | x]$$

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



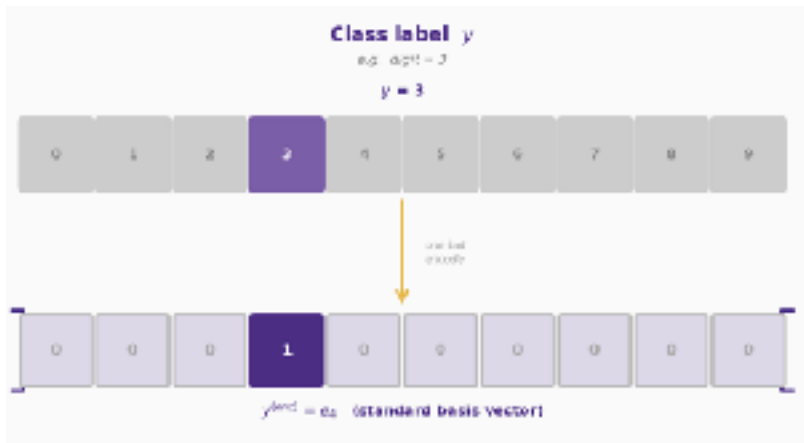
$$\begin{array}{ll} w_1 \cdot x & \mathbb{P}[c_1 | x] \\ w_2 \cdot x & \mathbb{P}[c_2 | x] \\ \vdots & \vdots \\ & \mathbb{P}[c_k | x] \end{array}$$

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



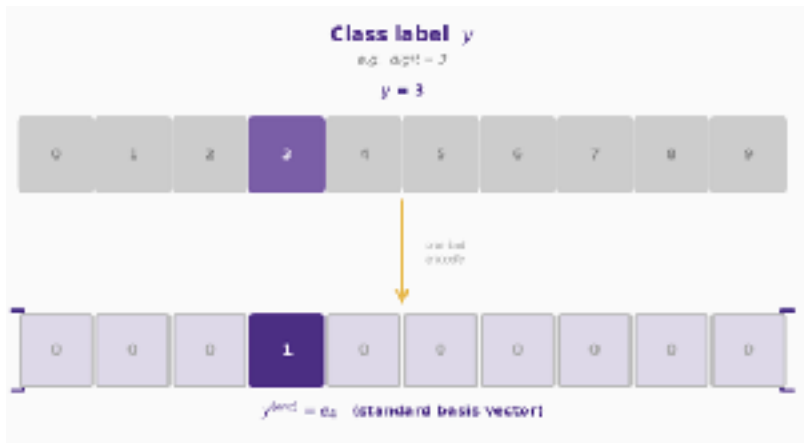
$$\begin{array}{ll} w_1 \cdot x & \mathbb{P}[c_1 | x] \\ w_2 \cdot x & \mathbb{P}[c_2 | x] \\ \vdots & \vdots \\ w_k \cdot x & \mathbb{P}[c_k | x] \end{array}$$

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



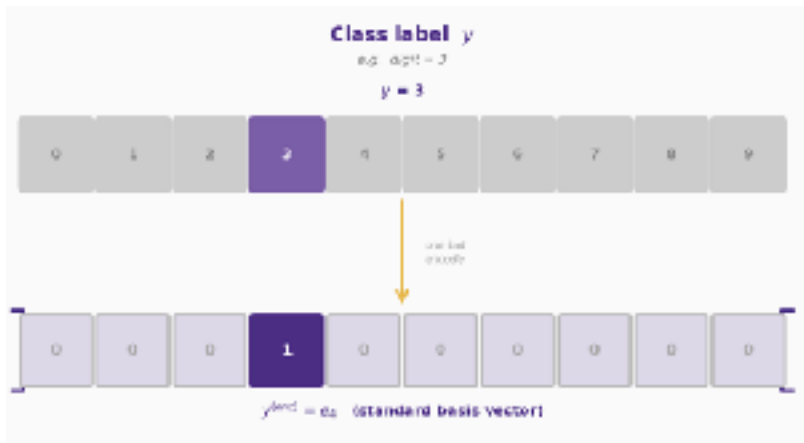
|               |                       |                       |
|---------------|-----------------------|-----------------------|
| $w_1 \cdot x$ | $\mathbb{P}[c_1   x]$ | $e^{w_1 \cdot x} / N$ |
| $w_2 \cdot x$ | $\mathbb{P}[c_2   x]$ |                       |
| $\vdots$      | $\vdots$              |                       |
| $w_k \cdot x$ | $\mathbb{P}[c_k   x]$ |                       |

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



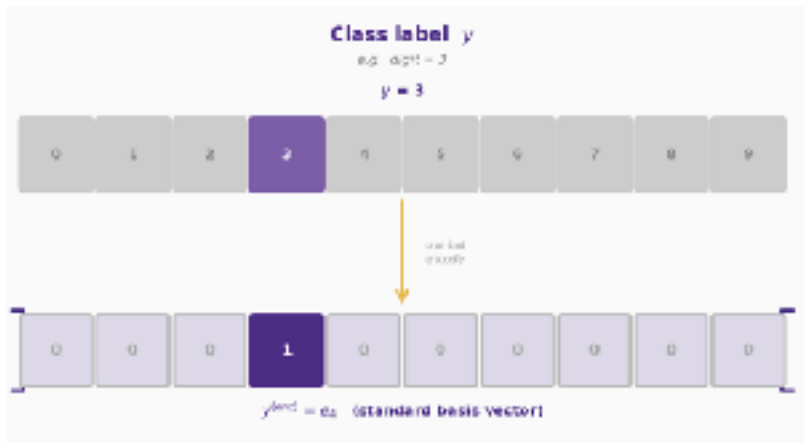
|               |                       |                       |
|---------------|-----------------------|-----------------------|
| $w_1 \cdot x$ | $\mathbb{P}[c_1   x]$ | $e^{w_1 \cdot x} / N$ |
| $w_2 \cdot x$ | $\mathbb{P}[c_2   x]$ | $e^{w_2 \cdot x} / N$ |
| $\vdots$      | $\vdots$              |                       |
| $w_k \cdot x$ | $\mathbb{P}[c_k   x]$ |                       |

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



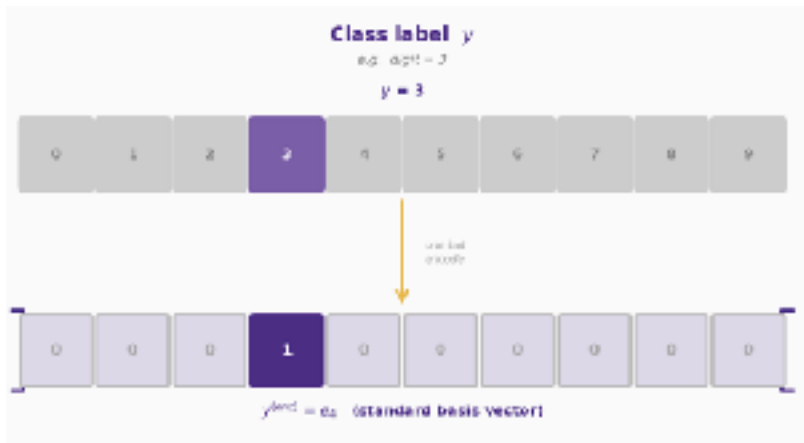
|               |                       |                       |
|---------------|-----------------------|-----------------------|
| $w_1 \cdot x$ | $\mathbb{P}[c_1   x]$ | $e^{w_1 \cdot x} / N$ |
| $w_2 \cdot x$ | $\mathbb{P}[c_2   x]$ | $e^{w_2 \cdot x} / N$ |
| $\vdots$      | $\vdots$              | $\vdots$              |
| $w_k \cdot x$ | $\mathbb{P}[c_k   x]$ |                       |

# One-HOT encoding labels + softmax

- So far: binary  $y \in \{-1, +1\}$
- In general:  $y \in \{c_1, c_2, \dots, c_k\}$
- $c_j$ 's are called classes or labels

Ideas for how to use logistic regression for multi class prediction:

- Try to predict  $\mathbb{P}[c_i | x]$  versus  $\mathbb{P}[\text{not } c_i | x]$  for each  $i$
- Try to predict  $\mathbb{P}[c_i | x]$  as a linear function  $w_i \cdot x$  for each  $i$



|               |                       |                       |
|---------------|-----------------------|-----------------------|
| $w_1 \cdot x$ | $\mathbb{P}[c_1   x]$ | $e^{w_1 \cdot x} / N$ |
| $w_2 \cdot x$ | $\mathbb{P}[c_2   x]$ | $e^{w_2 \cdot x} / N$ |
| $\vdots$      | $\vdots$              | $\vdots$              |
| $w_k \cdot x$ | $\mathbb{P}[c_k   x]$ | $e^{w_k \cdot x} / N$ |

# Softmax classification

# Softmax classification

## 2 classes

Conditional probabilities

$$P(y = 1 \mid x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 \mid x) = \frac{1}{1 + e^{w^\top x}}$$

# Softmax classification

2 classes

# Sigmoid

Conditional probabilities

$$P(y = 1 \mid x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 \mid x) = \frac{1}{1 + e^{w^\top x}}$$

# Softmax classification

2 classes

# Sigmoid

Conditional probabilities

$$P(y = 1 \mid x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 \mid x) = \frac{1}{1 + e^{w^\top x}}$$

$k$  classes

$$P(y = c_j \mid x) = \frac{e^{w_j^\top x}}{\sum_{j'} e^{w_{j'}^\top x}}$$

# Softmax classification

2 classes

# Sigmoid

Conditional probabilities

$k$  classes

# Softmax

$$P(y = 1 \mid x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 \mid x) = \frac{1}{1 + e^{w^\top x}}$$

$$P(y = c_j \mid x) = \frac{e^{w_j^\top x}}{\sum_{j'} e^{w_{j'}^\top x}}$$

# Softmax classification

2 classes

# Sigmoid

Conditional probabilities

$k$  classes

# Softmax

$$P(y = 1 \mid x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 \mid x) = \frac{1}{1 + e^{w^\top x}}$$

$$P(y = c_j \mid x) = \frac{e^{w_j^\top x}}{\sum_{j'} e^{w_{j'}^\top x}}$$

MLE

$$\operatorname{argmax}_w \sum_{i=1}^n \log \left( \frac{1}{1 + e^{-y_i w^\top x_i}} \right)$$

# Softmax classification

2 classes

# Sigmoid

Conditional probabilities

$$P(y = 1 | x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 | x) = \frac{1}{1 + e^{w^\top x}}$$

$k$  classes

# Softmax

$$P(y = c_j | x) = \frac{e^{w_j^\top x}}{\sum_{j'} e^{w_{j'}^\top x}}$$

MLE

$$\operatorname{argmax}_w \sum_{i=1}^n \log \left( \frac{1}{1 + e^{-y_i w^\top x_i}} \right)$$

# Logistic loss

# Softmax classification

2 classes

# Sigmoid

Conditional probabilities

$$P(y = 1 | x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 | x) = \frac{1}{1 + e^{w^\top x}}$$

$k$  classes

# Softmax

$$P(y = c_j | x) = \frac{e^{w_j^\top x}}{\sum_{j'} e^{w_{j'}^\top x}}$$

MLE

$$\operatorname{argmax}_w \sum_{i=1}^n \log \left( \frac{1}{1 + e^{-y_i w^\top x_i}} \right)$$

# Logistic loss

$$\operatorname{argmax}_w \sum_{i=1}^n \sum_{j=1}^k 1\{y_i = c_j\} \log \left( \frac{e^{w_j^\top x_i}}{\sum_{j'=1}^k e^{w_{j'}^\top x_i}} \right)$$

# Softmax classification

2 classes

# Sigmoid

Conditional probabilities

$$P(y = 1 | x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 | x) = \frac{1}{1 + e^{w^\top x}}$$

$k$  classes

# Softmax

$$P(y = c_j | x) = \frac{e^{w_j^\top x}}{\sum_{j'} e^{w_{j'}^\top x}}$$

MLE

$$\operatorname{argmax}_w \sum_{i=1}^n \log \left( \frac{1}{1 + e^{-y_i w^\top x_i}} \right)$$

# Logistic loss

$$\operatorname{argmax}_w \sum_{i=1}^n \sum_{j=1}^k 1\{y_i = c_j\} \log \left( \frac{e^{w_j^\top x_i}}{\sum_{j'=1}^k e^{w_{j'}^\top x_i}} \right)$$

# Binary cross entropy

# Softmax classification

without loss of generality, set  $k = 2$ ,  $w_1 = 0$

2 classes

$k$  classes

# Sigmoid

Conditional probabilities

# Softmax

$$P(y = 1 | x) = \frac{1}{1 + e^{-w^\top x}} = \frac{e^{w^\top x}}{1 + e^{w^\top x}}$$

$$P(y = -1 | x) = \frac{1}{1 + e^{w^\top x}}$$

$$P(y = c_j | x) = \frac{e^{w_j^\top x}}{\sum_{j'} e^{w_{j'}^\top x}}$$

MLE

$$\operatorname{argmax}_w \sum_{i=1}^n \log \left( \frac{1}{1 + e^{-y_i w^\top x_i}} \right)$$

# Logistic loss

$$\operatorname{argmax}_w \sum_{i=1}^n \sum_{j=1}^k 1\{y_i = c_j\} \log \left( \frac{e^{w_j^\top x_i}}{\sum_{j'=1}^k e^{w_{j'}^\top x_i}} \right)$$

# Binary cross entropy

# Regression and classification

# Regression and classification

- ML paradigm: define prediction  $f_w(x)$  and loss  $\ell(f_w(x), y)$

# Regression and classification

- ML paradigm: define prediction  $f_w(x)$  and loss  $\ell(f_w(x), y)$
- Then optimize:

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n \ell(f_w(x_i), y_i)$$

# Regression and classification

- ML paradigm: define prediction  $f_w(x)$  and loss  $\ell(f_w(x), y)$

- Then optimize:

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n \ell(f_w(x_i), y_i)$$

- Squared error loss:  $\ell(f_w(x), y) = (y - f_w(x))^2$

# Regression and classification

- ML paradigm: define prediction  $f_w(x)$  and loss  $\ell(f_w(x), y)$
- Then optimize:

$$\hat{w} = \operatorname{argmin}_w \sum_{i=1}^n \ell(f_w(x_i), y_i)$$

- Squared error loss:  $\ell(f_w(x), y) = (y - f_w(x))^2$
- Logistic loss:  $\ell(f_w(x), y) = \log(1 + \exp(-yf_w(x)))$



# Regression and classification

- Can we treat classification as a regression problem?
  - # Yes: could regress to a continuous output  $y$  and round to a valid label
- Can we treat regression as a classification problem?

# Regression and classification

- Can we treat classification as a regression problem?

# Yes: could regress to a continuous output  $y$  and round to a valid label

# Is this a good idea?

- Can we treat regression as a classification problem?

# Regression and classification

- Can we treat classification as a regression problem?

# Yes: could regress to a continuous output  $y$  and round to a valid label

# Is this a good idea?

No, no smooth relationship between label values

- Can we treat regression as a classification problem?

# Regression and classification

- Can we treat classification as a regression problem?

# Yes: could regress to a continuous output  $y$  and round to a valid label

# Is this a good idea?

No, no smooth relationship between label values

- Can we treat regression as a classification problem?

# Yes: classify which “bin” the true number falls into:  $c_1 = [0-10]$ ,  $c_2 = [11-20]$ , ...

# Regression and classification

- Can we treat classification as a regression problem?

# Yes: could regress to a continuous output  $y$  and round to a valid label

# Is this a good idea?

No, no smooth relationship between label values

- Can we treat regression as a classification problem?

# Yes: classify which “bin” the true number falls into:  $c_1 = [0-10]$ ,  $c_2 = [11-20]$ , ...

# Is this a good idea?

# Regression and classification

- Can we treat classification as a regression problem?

# Yes: could regress to a continuous output  $y$  and round to a valid label

# Is this a good idea?

No, no smooth relationship between label values

- Can we treat regression as a classification problem?

# Yes: classify which “bin” the true number falls into:  $c_1 = [0-10]$ ,  $c_2 = [11-20]$ , ...

# Is this a good idea?

You lose some precision, but it actually can be. Sometimes easier to optimize. Can win a kaggle competition this way. Look into “distributional RL”

# The ML pipeline

1. Task
2. Data
3. Model family
4. Training loss
4. Optimize
5. Pick hyperparameters
6. Evaluate

# Non-parametric methods

## Nearest Neighbours

---

Natasha Jaques



# Parametric vs non-parametric

# Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples

# Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples

# Learning a bunch of parameters / weights. Like linear regression, neural networks

# Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples

# Learning a bunch of parameters / weights. Like linear regression, neural networks

# Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples
  - # Learning a bunch of parameters / weights. Like linear regression, neural networks
- A model is non-parametric if # parameters increases with # samples
  - Does not mean absence of parameters!

# Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples

# Learning a bunch of parameters / weights. Like linear regression, neural networks

- A model is non-parametric if # parameters increases with # samples

- Does not mean absence of parameters!

# Today's class

This lecture:  $k$  nearest neighbors

# This lecture: $k$ nearest neighbors

# Simple

# This lecture: $k$ nearest neighbors

- Assume we have a classification task

# Simple

# This lecture: $k$ nearest neighbors

# Simple

- Assume we have a classification task
- To classify a new point  $x$ :
  - Find its  $k$  nearest neighbors in the training data
  - Set  $y$  to be the majority vote of the labels of these nearest neighbors

# This lecture: $k$ nearest neighbors

# Simple

- Assume we have a classification task
  - To classify a new point  $x$ :
    - Find its  $k$  nearest neighbors in the training data
    - Set  $y$  to be the majority vote of the labels of these nearest neighbors
- # What do we mean by nearest neighbors?

# This lecture: $k$ nearest neighbors

# Simple

- Assume we have a classification task
- To classify a new point  $x$ :
  - Find its  $k$  nearest neighbors in the training data
  - Set  $y$  to be the majority vote of the labels of these nearest neighbors

# What do we mean by nearest neighbors?

# Smallest distance in feature space (with  $d$  features)

# This lecture: $k$ nearest neighbors

# Simple

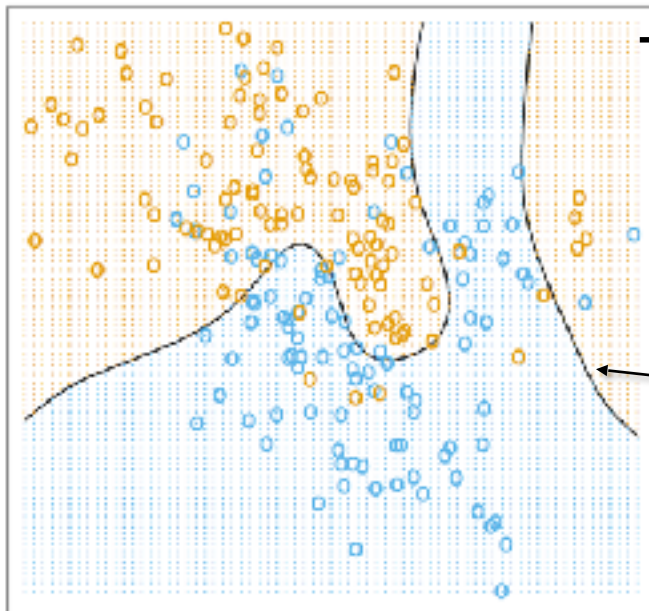
- Assume we have a classification task
- To classify a new point  $x$ :
  - Find its  $k$  nearest neighbors in the training data
  - Set  $y$  to be the majority vote of the labels of these nearest neighbors

# What do we mean by nearest neighbors?

# Smallest distance in feature space (with  $d$  features)

- Design choices / hyperparameters:
  - Number of nearest neighbors  $k$
  - Distance metric
  - Aggregation method

# Example: Bayes classifier



Training data:

- True label: +1
- True label: -1

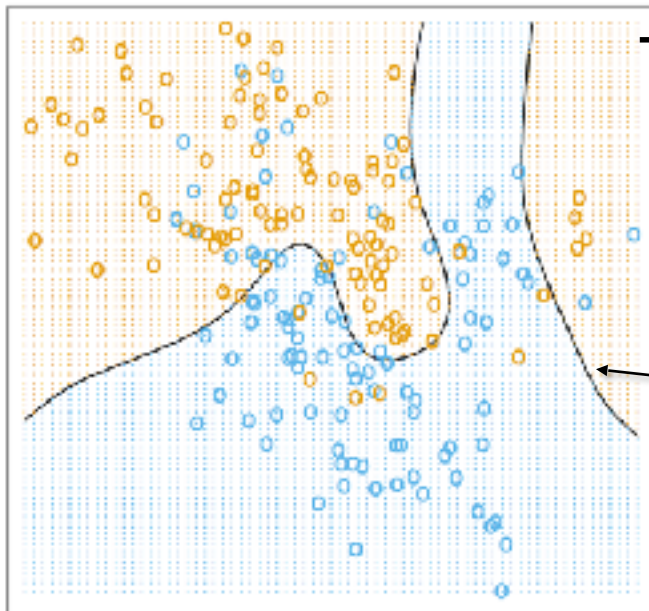
Optimal Bayes classifier:

$$\mathbb{P}(Y = 1|X = x) = \frac{1}{2}$$

▣ Predicted label: +1

▣ Predicted label: -1

# Example: Bayes classifier



Training data:

- True label: +1
- True label: -1

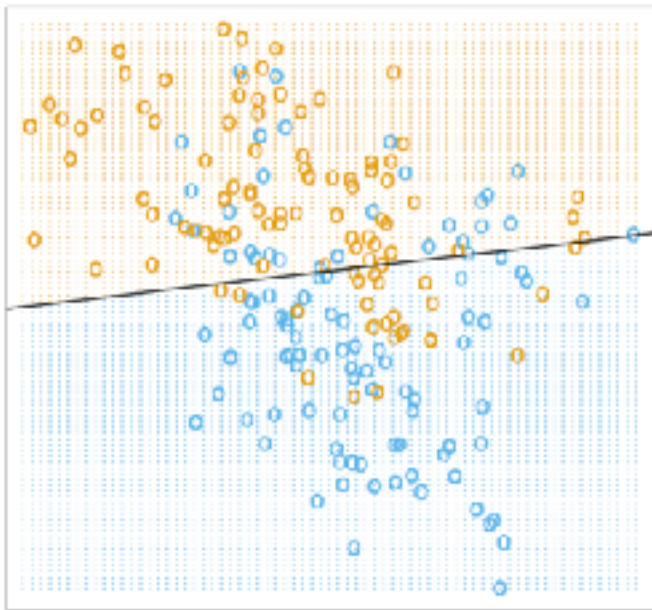
Optimal Bayes classifier:

$$\mathbb{P}(Y = 1 | X = x) = \frac{1}{2} \quad \# \text{ Ground truth}$$

 Predicted label: +1

 Predicted label: -1

# Linear decision boundary



Training data:

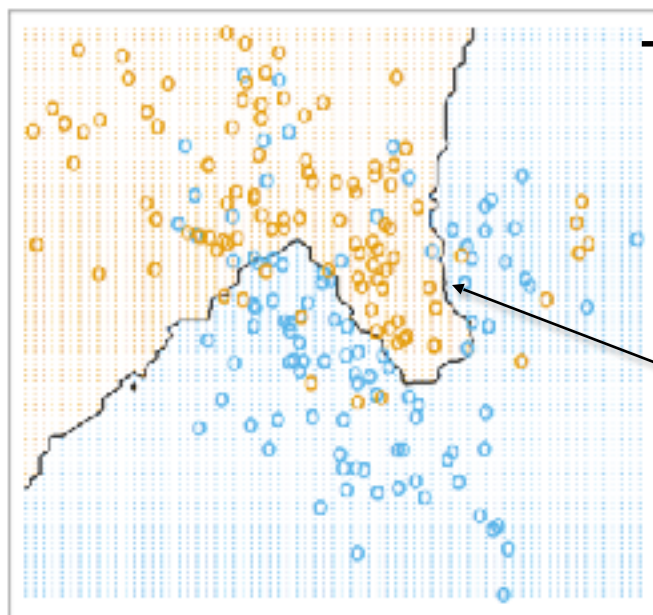
- True label: +1
- True label: -1

Learned linear decision bound  
 $x^T w + b = 0$

 Predicted label: +1

 Predicted label: -1

# $k = 15$ nearest neighbors boundary



Training data:

○ True label: +1

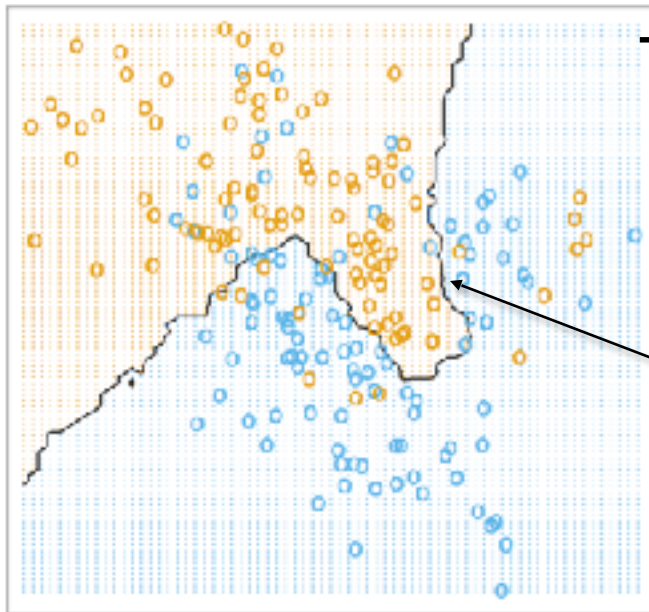
○ True label: -1

15 nearest neighbors  
decision boundary

(Priority vote)

○ Predicted label: +1  
○ Predicted label: -1

# $k = 15$ nearest neighbors boundary



Training data:

- True label: +1
- True label: -1

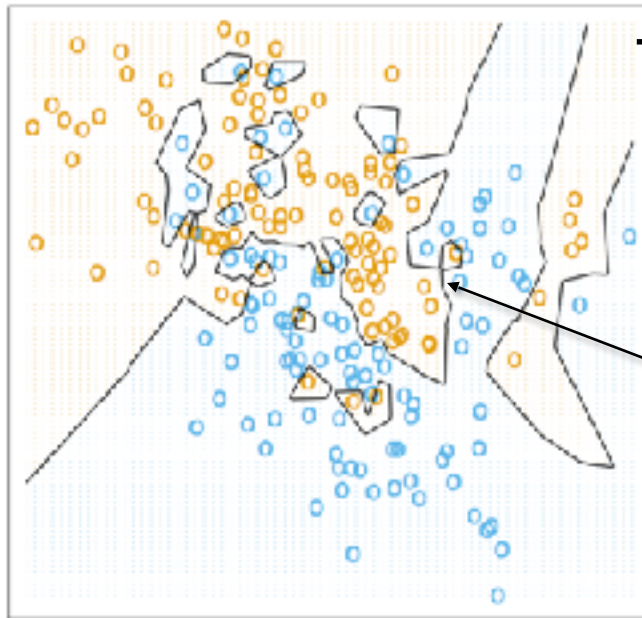
15 nearest neighbors  
decision boundary

(Priority vote)

○ Predicted label: +1  
○ Predicted label: -1

# Boundary depends on  
the training data points

# $k = 1$ nearest neighbor boundary



Training data:

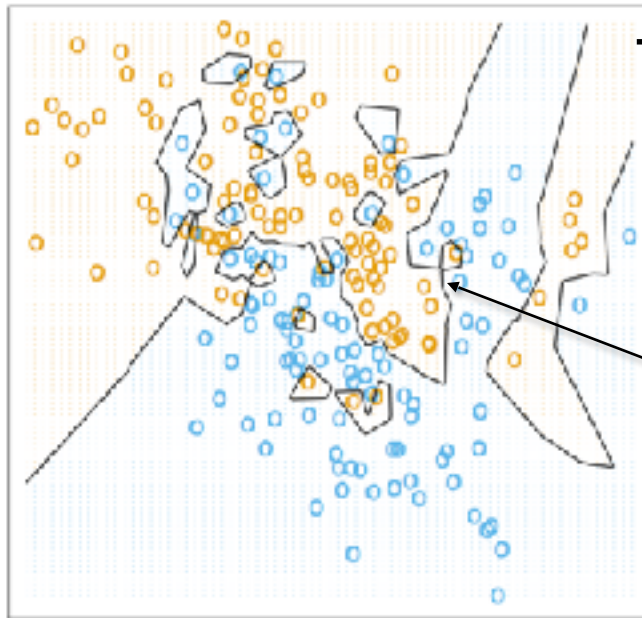
- True label: +1
- True label: -1

1 nearest neighbor  
decision boundary  
(majority vote)

Predicted label: -1

Predicted label: +1

# $k = 1$ nearest neighbor boundary



# Overfitting

Training data:

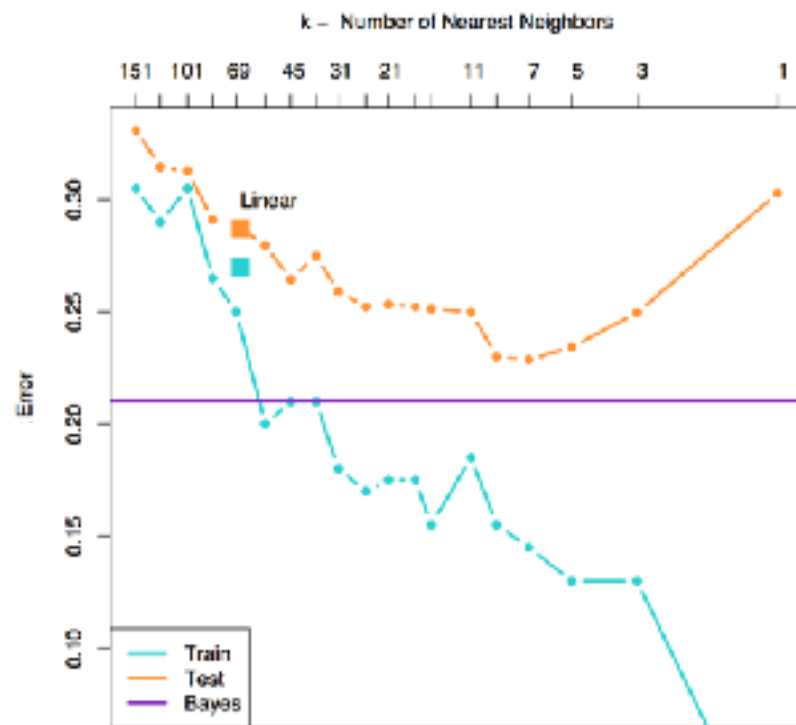
- True label: +1
- True label: -1

1 nearest neighbor  
decision boundary  
(majority vote)

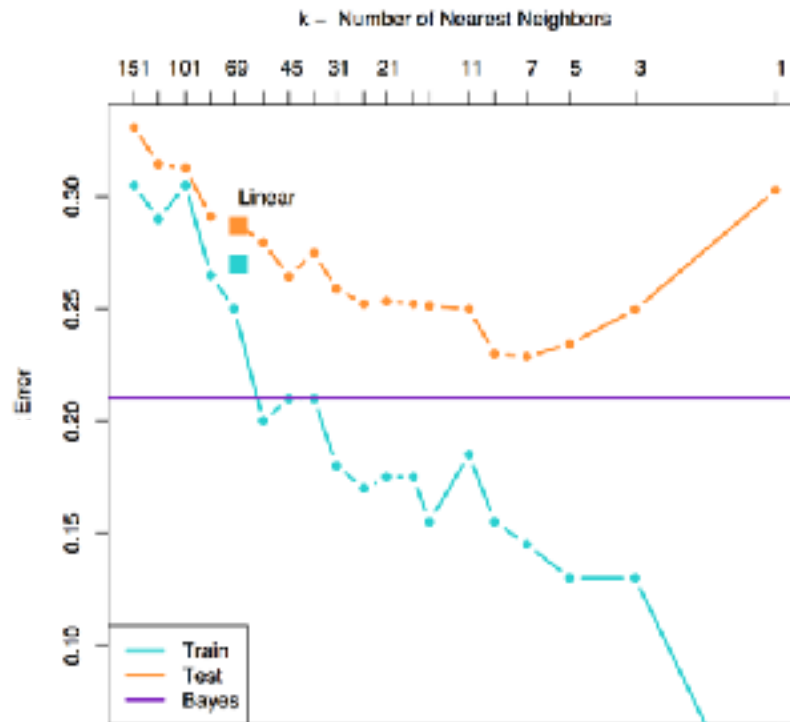
Predicted label: -1

Predicted label: +1

# $k$ nearest neighbors error

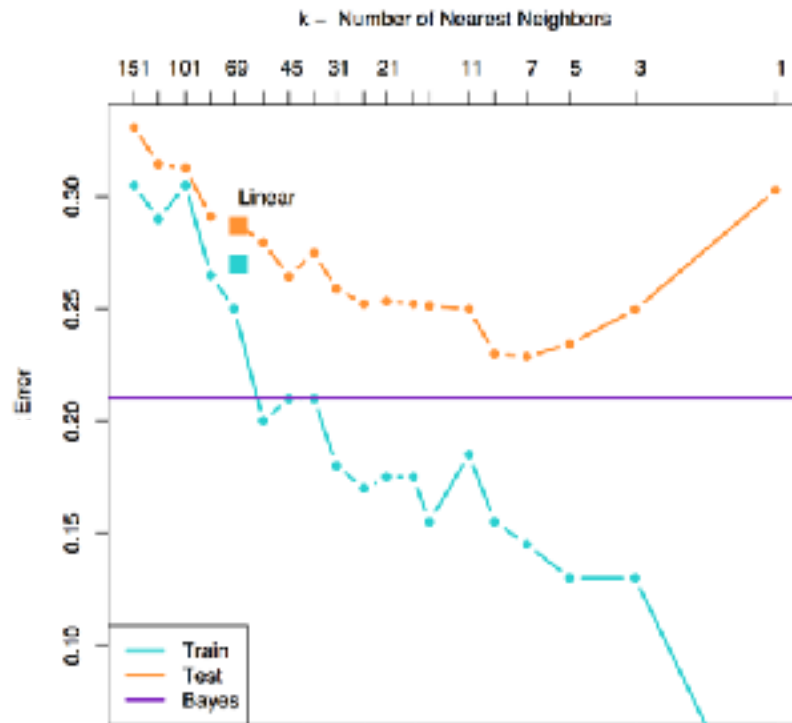


# $k$ nearest neighbors error



# Overfitting

# $k$ nearest neighbors error



# Underfitting

# Overfitting

# Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples
- A model is non-parametric if # parameters increases with # samples

# Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples
  - # After training, you can discard your data
- A model is non-parametric if # parameters increases with # samples

# Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples
  - # After training, you can discard your data
- A model is non-parametric if # parameters increases with # samples
  - # Keep training data around for inference

# Notable distance metrics & level sets

# Notable distance metrics & level sets

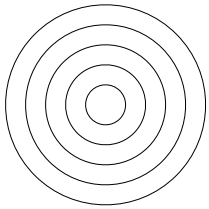
$\ell_2$  norm (Euclidean)

# Notable distance metrics & level sets

$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$

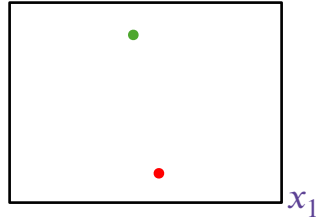
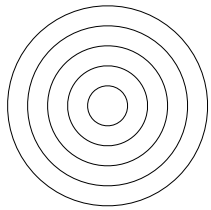
# Notable distance metrics & level sets

$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



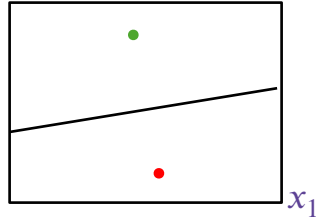
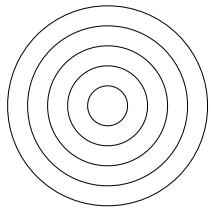
# Notable distance metrics & level sets

$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



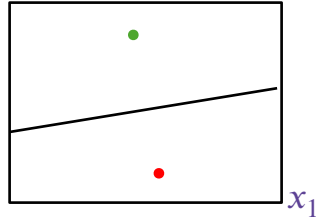
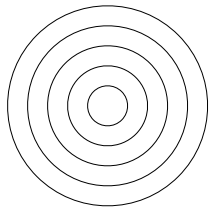
# Notable distance metrics & level sets

$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



# Notable distance metrics & level sets

$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$

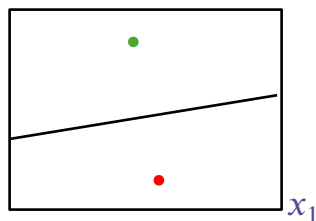
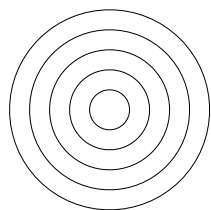


$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$

# Notable distance metrics & level sets

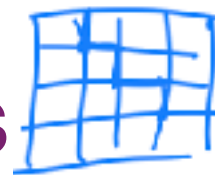


$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$

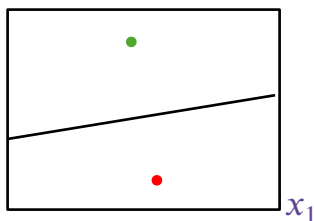
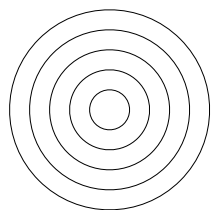


$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$

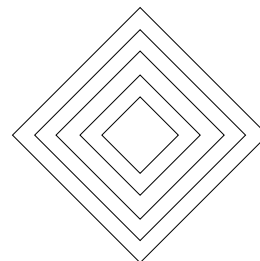
# Notable distance metrics & level sets



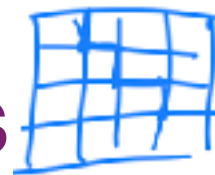
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



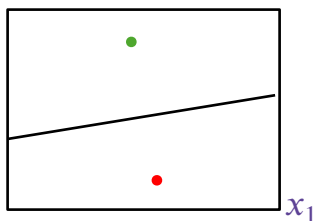
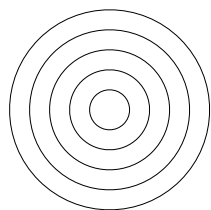
$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



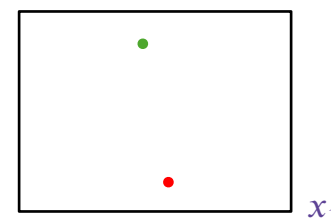
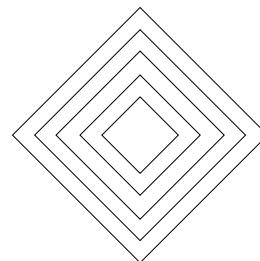
# Notable distance metrics & level sets



$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



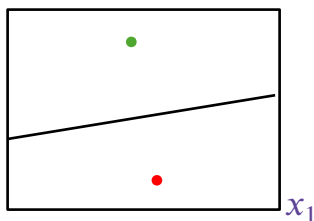
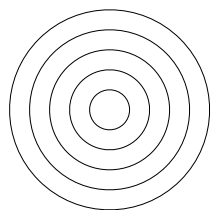
$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



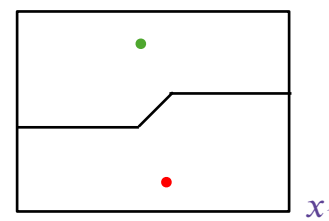
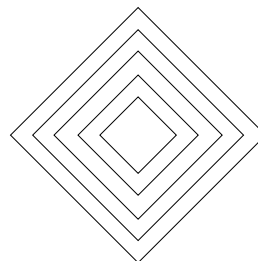
# Notable distance metrics & level sets



$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



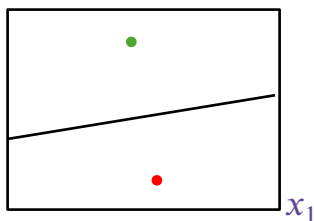
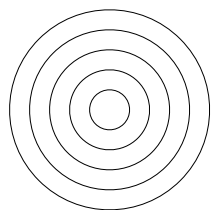
$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



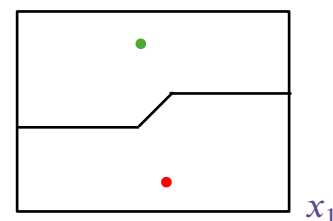
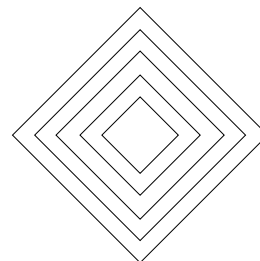
# Notable distance metrics & level sets



$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$

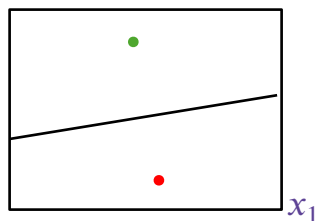
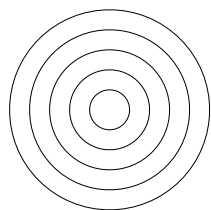


Mahalanobis norm

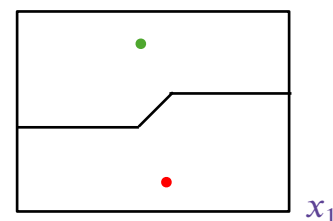
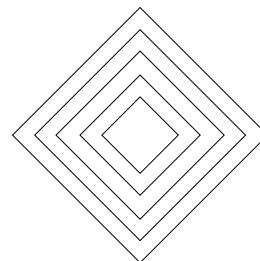
# Notable distance metrics & level sets



$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$

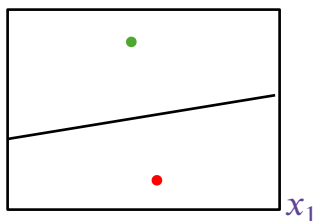
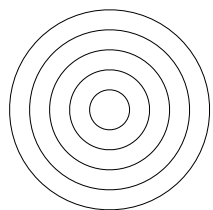


Mahalanobis norm  $(u - v)^T M (u - v)$

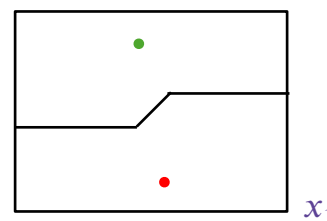
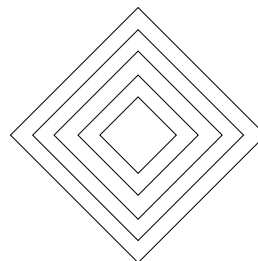
# Notable distance metrics & level sets



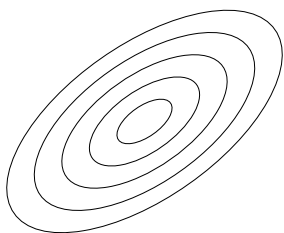
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



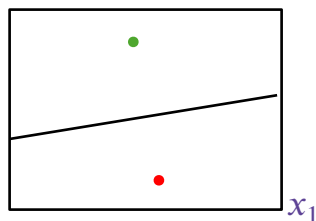
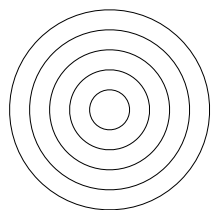
Mahalanobis norm  $(u - v)^T M (u - v)$



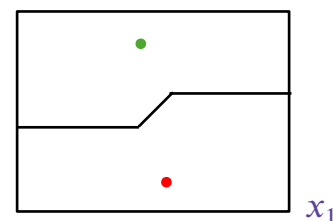
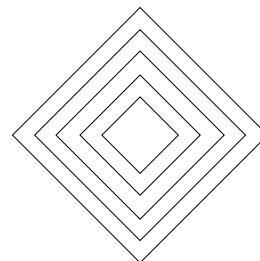
# Notable distance metrics & level sets



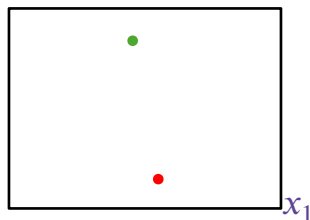
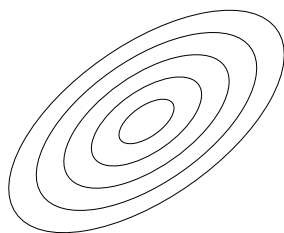
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



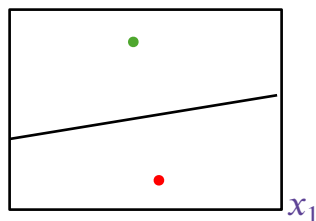
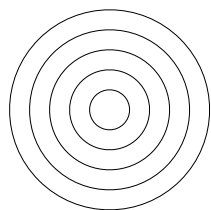
Mahalanobis norm  $(u - v)^T M (u - v)$



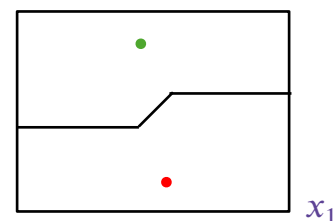
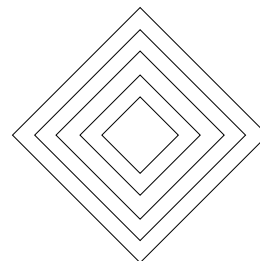
# Notable distance metrics & level sets



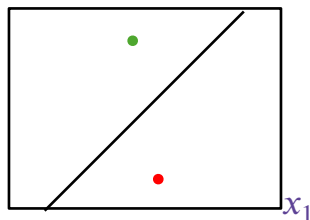
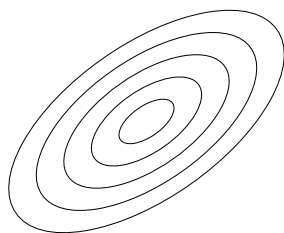
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



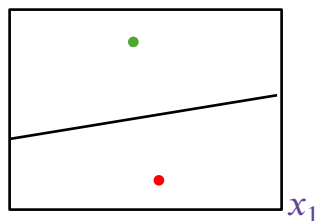
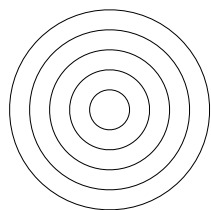
Mahalanobis norm  $(u - v)^T M (u - v)$



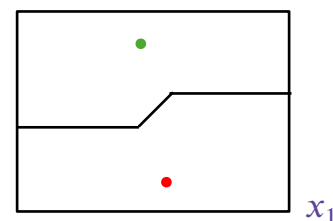
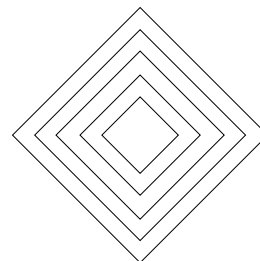
# Notable distance metrics & level sets



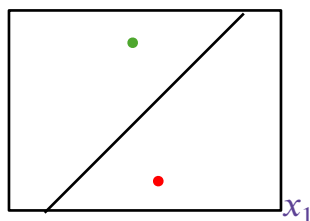
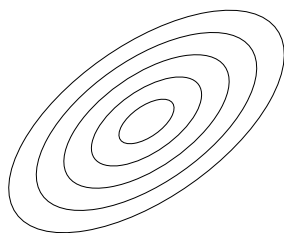
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



Mahalanobis norm  $(u - v)^T M (u - v)$

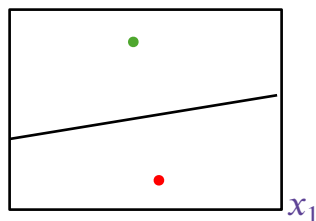
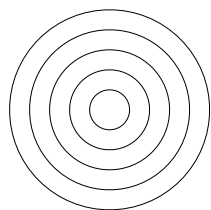


# Weight different dimensions differently

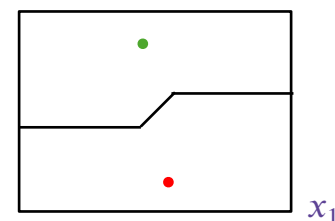
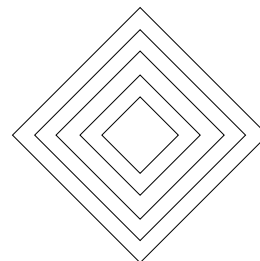
# Notable distance metrics & level sets



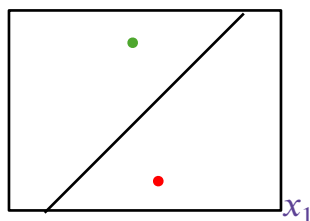
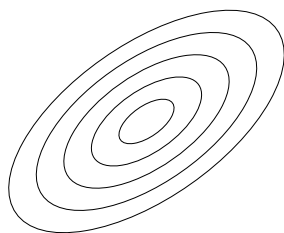
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



Mahalanobis norm  $(u - v)^T M (u - v)$



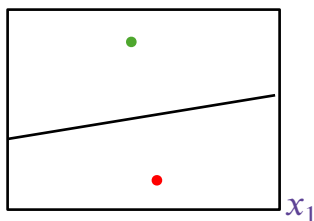
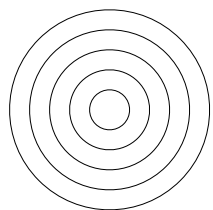
$\ell_\infty$  norm (max)

# Weight different dimensions differently

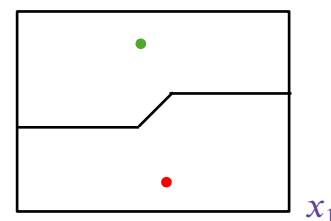
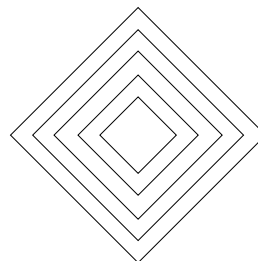
# Notable distance metrics & levels



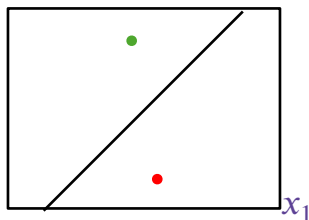
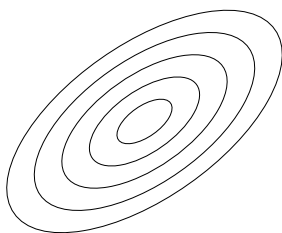
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



Mahalanobis norm  $(u - v)^T M (u - v)$



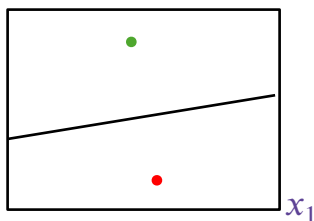
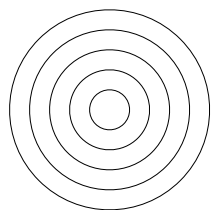
$\ell_\infty$  norm (max)  $\|u - v\|_\infty = \max_i |u_i - v_i|$

# Weight different dimensions differently

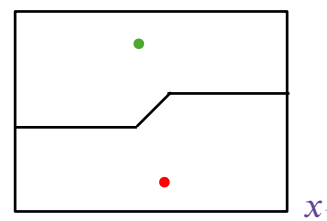
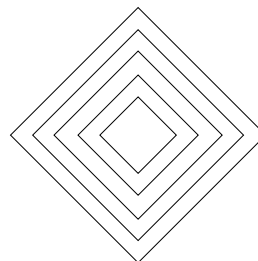
# Notable distance metrics & level sets



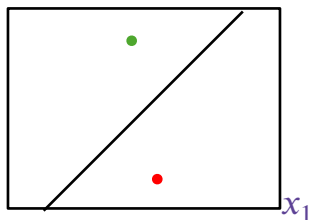
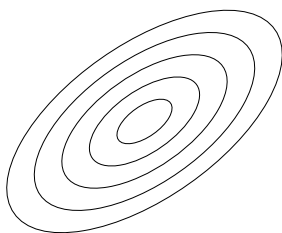
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



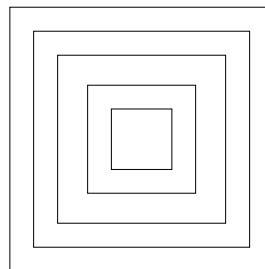
$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



Mahalanobis norm  $(u - v)^T M (u - v)$



$\ell_\infty$  norm (max)  $\|u - v\|_\infty = \max_i |u_i - v_i|$

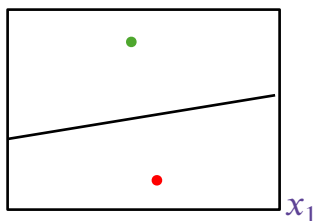
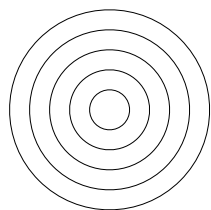


# Weight different dimensions differently

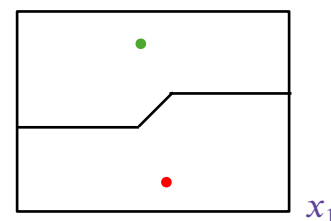
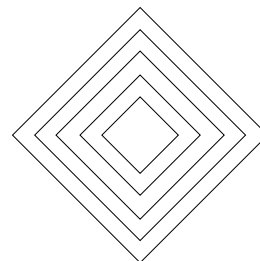
# Notable distance metrics & level sets



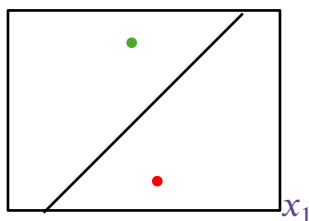
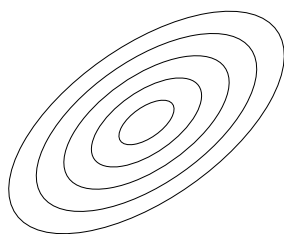
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



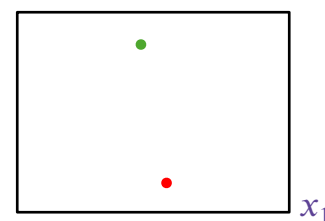
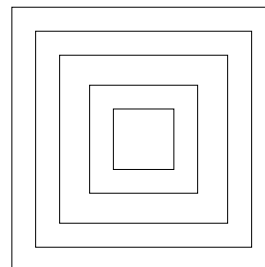
$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



Mahalanobis norm  $(u - v)^T M (u - v)$



$\ell_\infty$  norm (max)  $\|u - v\|_\infty = \max_i |u_i - v_i|$

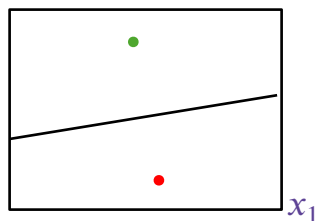
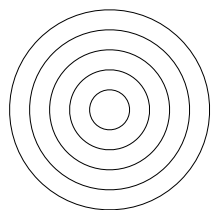


# Weight different dimensions differently

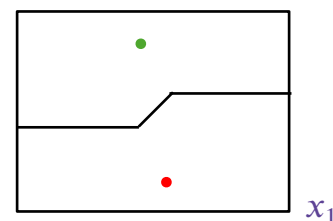
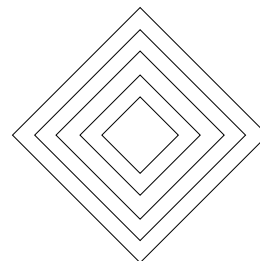
# Notable distance metrics & level sets



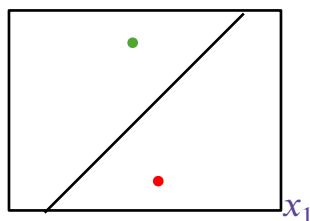
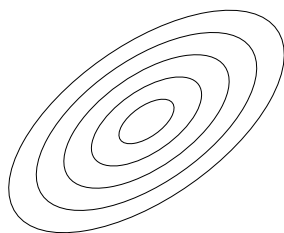
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



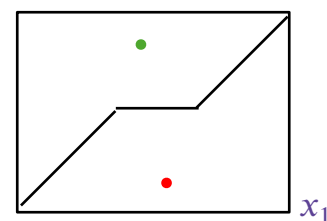
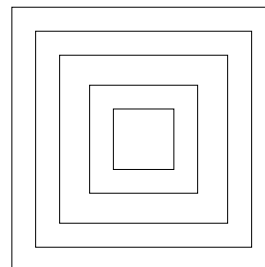
$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



Mahalanobis norm  $(u - v)^T M (u - v)$



$\ell_\infty$  norm (max)  $\|u - v\|_\infty = \max_i |u_i - v_i|$

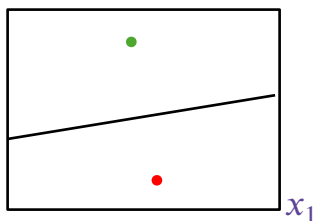
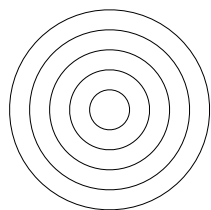


# Weight different dimensions differently

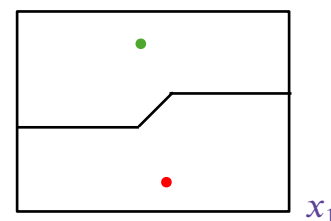
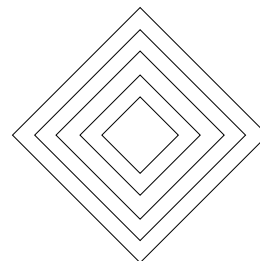
# Notable distance metrics & level sets



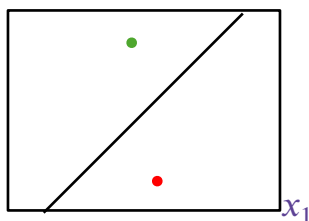
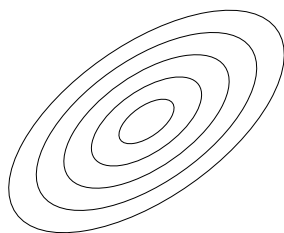
$\ell_2$  norm (Euclidean)  $d(u, v) = \|u - v\|_2^2$



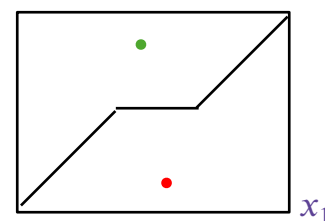
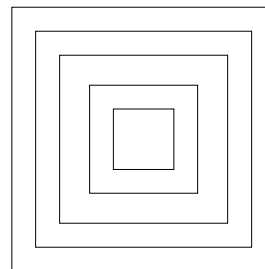
$\ell_1$  norm (Manhattan, taxicab)  $\|u - v\|_1 = \sum_{i=1}^d |u_i - v_i|$



Mahalanobis norm  $(u - v)^T M (u - v)$



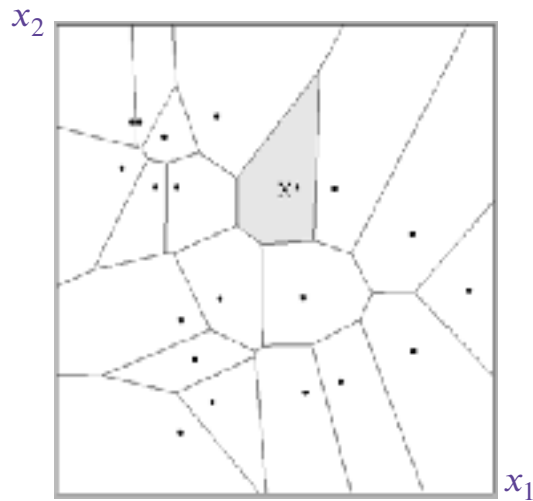
$\ell_\infty$  norm (max)  $\|u - v\|_\infty = \max_i |u_i - v_i|$



# Weight different dimensions differently # Max distance between vectors

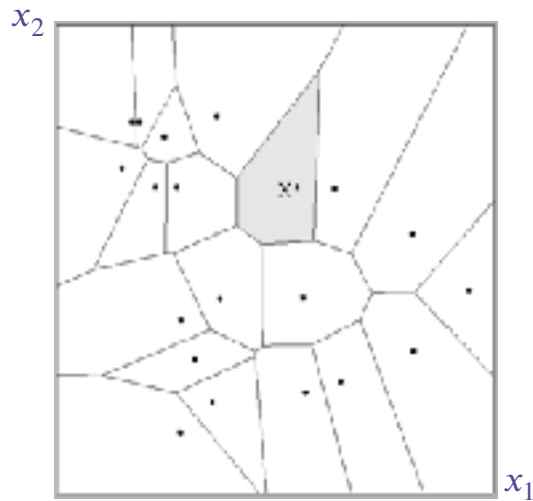
# Example: distance metrics with $k = 1$ NN

$$d(x, x') = (x_1 - x'_1)^2 + (x_2 - x'_2)^2$$

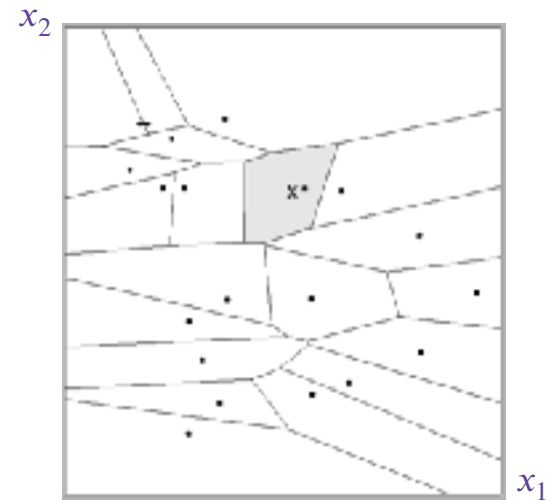


# Example: distance metrics with $k = 1$ NN

$$d(x, x') = (x_1 - x'_1)^2 + (x_2 - x'_2)^2$$



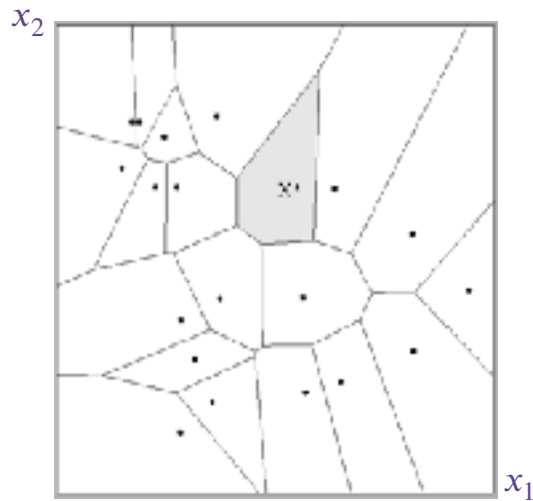
$$d(x, x') = (x_1 - x'_1)^2 + 9(x_2 - x'_2)^2$$



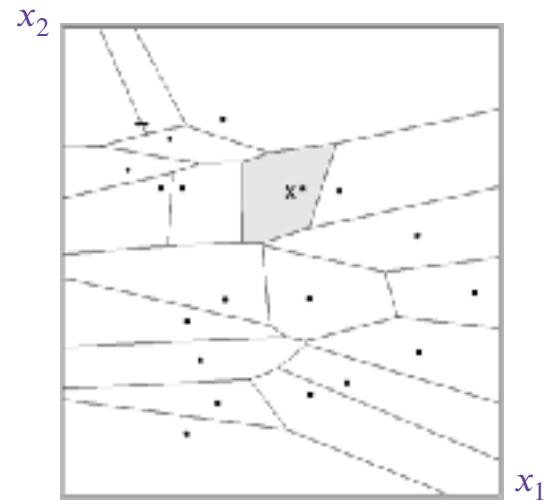
# Example: distance metrics with $k = 1$ NN

# L2

$$d(x, x') = (x_1 - x'_1)^2 + (x_2 - x'_2)^2$$



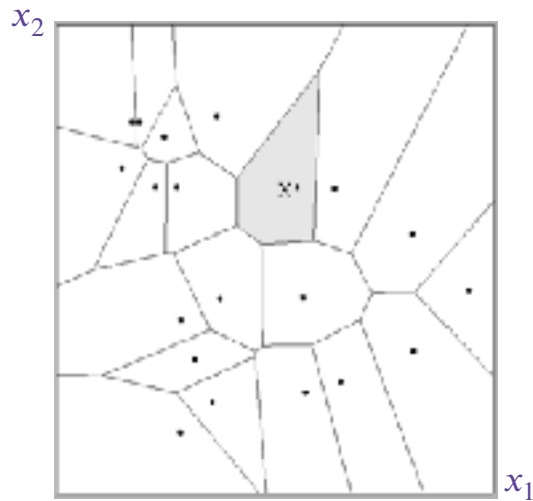
$$d(x, x') = (x_1 - x'_1)^2 + 9(x_2 - x'_2)^2$$



# Example: distance metrics with $k = 1$ NN

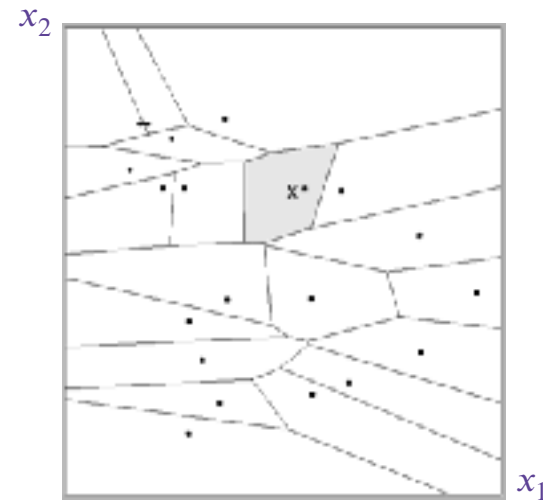
# L2

$$d(x, x') = (x_1 - x'_1)^2 + (x_2 - x'_2)^2$$



# Mahalanobis

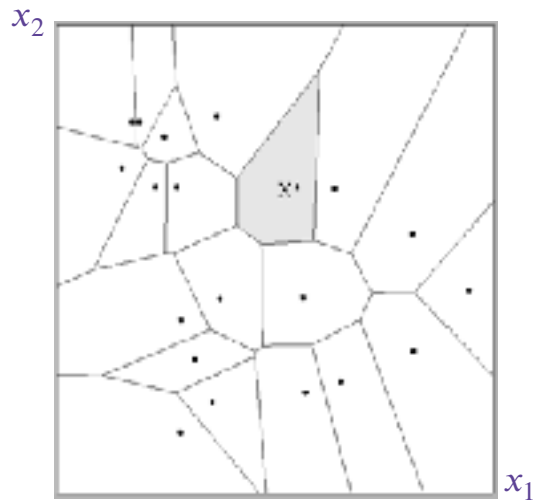
$$d(x, x') = (x_1 - x'_1)^2 + 9(x_2 - x'_2)^2$$



# Example: distance metrics with $k = 1$ NN

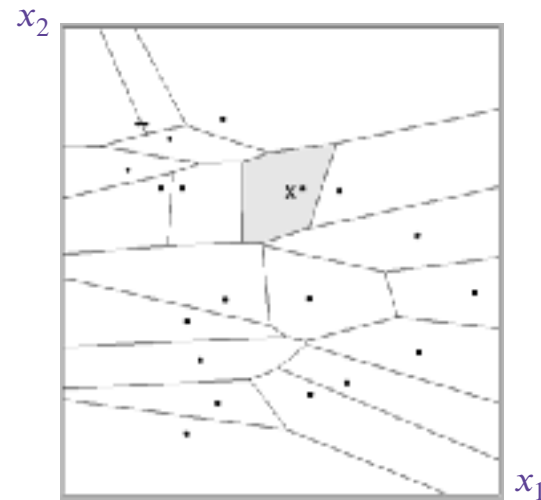
# L2

$$d(x, x') = (x_1 - x'_1)^2 + (x_2 - x'_2)^2$$



# Mahalanobis

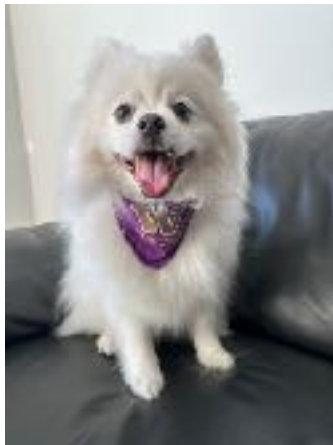
$$d(x, x') = (x_1 - x'_1)^2 + 9(x_2 - x'_2)^2$$



# Know the geometry of your feature space

# Learned distance metrics

Training data

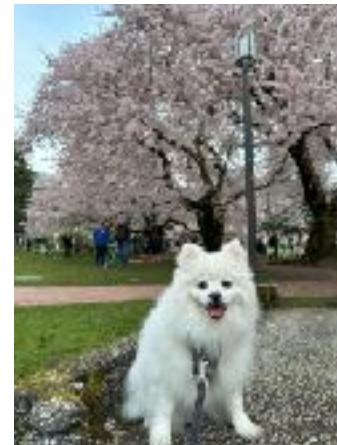


Dog



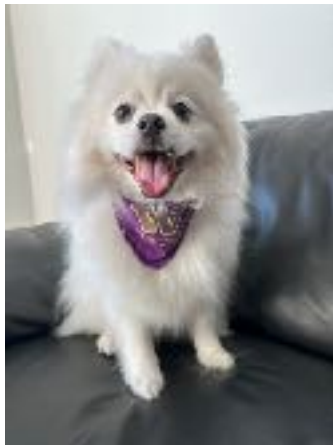
Cat

Test data



# Learned distance metrics

Training data

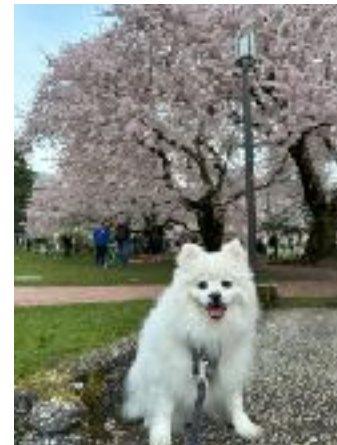


Dog



Cat

Test data



# Use a neural net to learn the distance function

## kNN demo

- <http://vision.stanford.edu/teaching/cs231n-demos/knn/>

# 1-NN classification: Theoretical guarantees

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

# 1-NN classification: Theoretical guarantees

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P$$

# Assume: we have enough data and **our true function is smooth**

# 1-NN classification: Theoretical guarantees

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d$$

# Assume: we have enough data and **our true function is smooth**

# 1-NN classification: Theoretical guarantees

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

# Assume: we have enough data and **our true function is smooth**

# 1-NN classification: Theoretical guarantees

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

# Assume: we have enough data and **our true function is smooth**

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

# 1-NN classification: Theoretical guarantees

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

# Assume: we have enough data and **our true function is smooth**

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

As  $n \rightarrow \infty$ ,

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

As  $n \rightarrow \infty$ ,

$$P(y | x_{\text{NN}}) \rightarrow P(y | x)$$

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

As  $n \rightarrow \infty$ ,

$$P(y | x_{\text{NN}}) \rightarrow P(y | x) \quad \# \text{ Why?}$$

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

As  $n \rightarrow \infty$ ,

$P(y | x_{\text{NN}}) \rightarrow P(y | x)$  # Why?

By assumptions, distance between  $x$  and  $x_{\text{NN}} \rightarrow 0$  as  $n \rightarrow \infty$

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

As  $n \rightarrow \infty$ ,

$P(y | x_{\text{NN}}) \rightarrow P(y | x)$  # Why?

Error:  $2P(y = 1 | x)P(y = 0 | x)$

By assumptions, distance between  $x$  and  $x_{\text{NN}} \rightarrow 0$  as  $n \rightarrow \infty$

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

As  $n \rightarrow \infty$ ,

$P(y | x_{\text{NN}}) \rightarrow P(y | x)$  # Why?

Error:  $2P(y = 1 | x)P(y = 0 | x)$

By assumptions, distance between  $x$  and  $x_{\text{NN}} \rightarrow 0$  as  $n \rightarrow \infty$

# Define Bayes error:

$$p^* = \min\{P(y = 1 | x), P(y = 0 | x)\}$$

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

As  $n \rightarrow \infty$ ,

$P(y | x_{\text{NN}}) \rightarrow P(y | x)$  # Why?

Error:  $2P(y = 1 | x)P(y = 0 | x)$

$$= 2p^*(1 - p^*)$$

By assumptions, distance between  $x$  and  $x_{\text{NN}} \rightarrow 0$  as  $n \rightarrow \infty$

# Define Bayes error:

$$p^* = \min\{P(y = 1 | x), P(y = 0 | x)\}$$

# 1-NN classification: Theoretical guarantees

# Assume: we have enough data and **our true function is smooth**

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d \quad y \in \{0, 1\}$$

Given test point  $x$ , let  $x_{\text{NN}}$  be the nearest neighbour in  $D$

Error if  $y_{\text{NN}} \neq y$

Case 1:  $y_{\text{NN}} = 1, y = 0$  w.p.  $P(y = 1 | x_{\text{NN}})P(y = 0 | x)$

Case 2:  $y_{\text{NN}} = 0, y = 1$  w.p.  $P(y = 0 | x_{\text{NN}})P(y = 1 | x)$

As  $n \rightarrow \infty$ ,

$P(y | x_{\text{NN}}) \rightarrow P(y | x)$  # Why?

Error:  $2P(y = 1 | x)P(y = 0 | x)$

$$= 2p^*(1 - p^*) \leq 2p^*$$

By assumptions, distance between  $x$  and  $x_{\text{NN}} \rightarrow 0$  as  $n \rightarrow \infty$

# Define Bayes error:

$$p^* = \min\{P(y = 1 | x), P(y = 0 | x)\}$$

# 1-NN classification: Theoretical guarantees

As  $n \rightarrow \infty$ ,  $\text{Error} = 2p^*(1 - p^*) \leq 2p^*$

**Theorem**[Cover, Hart, 1967] If  $P_X$  is supported everywhere in  $\mathbb{R}^d$  and  $P(Y = 1|X = x)$  is smooth everywhere, then as  $n \rightarrow \infty$  the 1-NN classification rule has error at most twice the Bayes error rate.

# 1-NN classification: Theoretical guarantees

As  $n \rightarrow \infty$ ,  $\text{Error} = 2p^*(1 - p^*) \leq 2p^*$

**Theorem**[Cover, Hart, 1967] If  $P_X$  is supported everywhere in  $\mathbb{R}^d$  and  $P(Y = 1|X = x)$  is smooth everywhere, then as  $n \rightarrow \infty$  the 1-NN classification rule has error at most twice the Bayes error rate.

# Bayes error is the best you can do given measurement error, so with infinite data nearest neighbors is a really good thing to do!

# 1-NN classification: Theoretical guarantees

As  $n \rightarrow \infty$ ,  $\text{Error} = 2p^*(1 - p^*) \leq 2p^*$

**Theorem**[Cover, Hart, 1967] If  $P_X$  is supported everywhere in  $\mathbb{R}^d$  and  $P(Y = 1|X = x)$  is smooth everywhere, then as  $n \rightarrow \infty$  the 1-NN classification rule has error at most twice the Bayes error rate.

# Bayes error is the best you can do given measurement error, so with infinite data nearest neighbors is a really good thing to do  
# Can fit arbitrarily complicated functions

# 1-NN classification: Theoretical guarantees

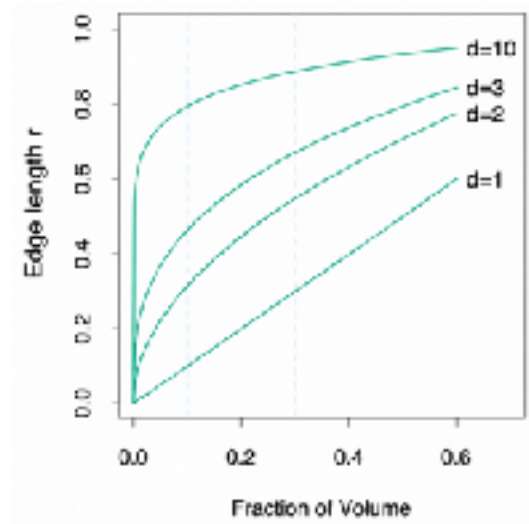
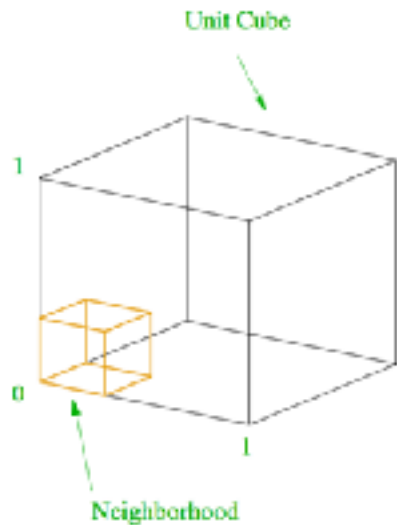
As  $n \rightarrow \infty$ ,  $\text{Error} = 2p^*(1 - p^*) \leq 2p^*$

**Theorem**[Cover, Hart, 1967] If  $P_X$  is supported everywhere in  $\mathbb{R}^d$  and  $P(Y = 1|X = x)$  is smooth everywhere, then as  $n \rightarrow \infty$  the 1-NN classification rule has error at most twice the Bayes error rate.

# Bayes error is the best you can do given measurement error, so with infinite data nearest neighbors  $\rightarrow$  fitting to  $\rightarrow$  ~~to~~ can fit arbitrarily complicated functions



# Curse of dimensionality, example 1



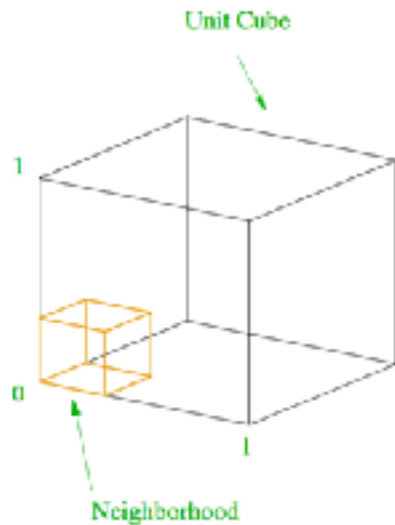
side length  $r$

$X$  is uniformly distributed over  $[0,1]^d$ . What is  $P(X \in [0,r]^d)$ ?

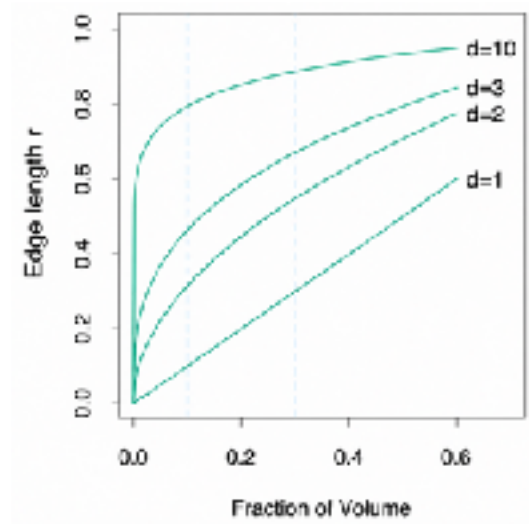
How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?

# Curse of dimensionality, example 1

# If  $d = 1$



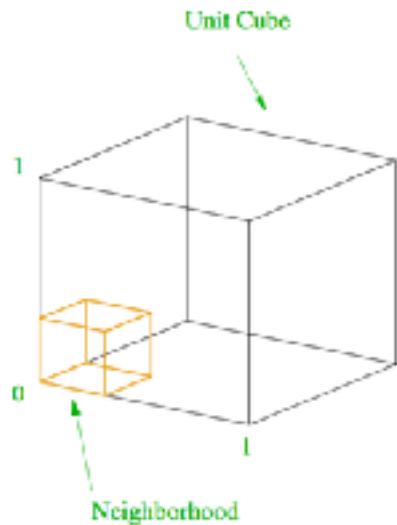
side length  $r$



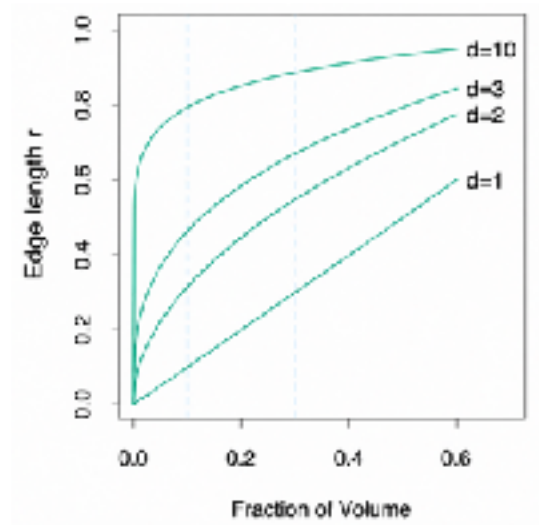
$X$  is uniformly distributed over  $[0,1]^d$ . What is  $P(X \in [0,r]^d)$ ?

How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?

# Curse of dimensionality, example 1



side length  $r$



# If  $d = 1$

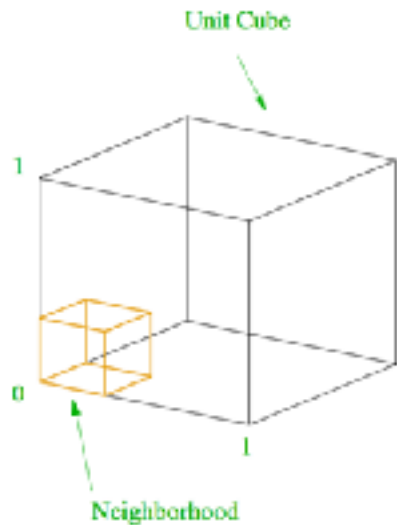


$X$  is uniformly distributed over  $[0,1]^d$ . What is  $P(X \in [0,r]^d)$

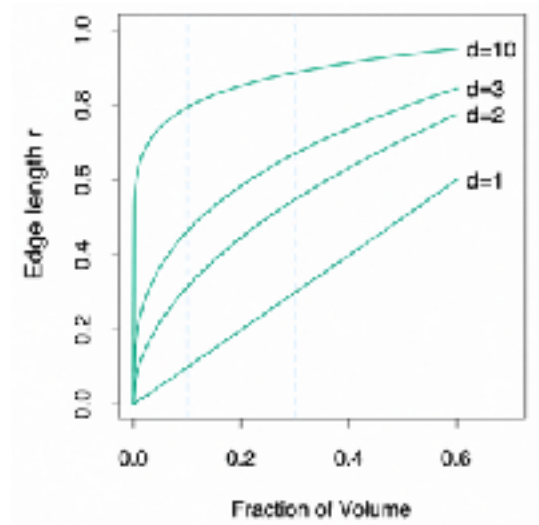
?

How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?

# Curse of dimensionality, example 1



side length  $r$



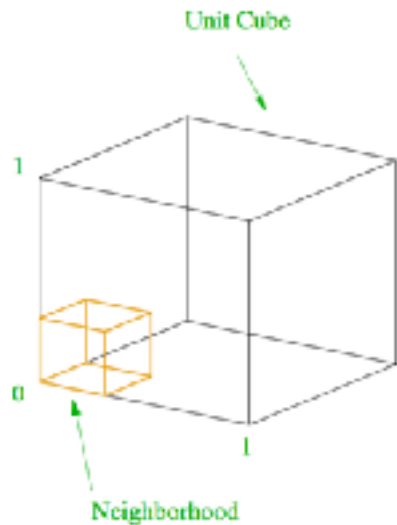
# If  $d = 10$   $P = 0.3$



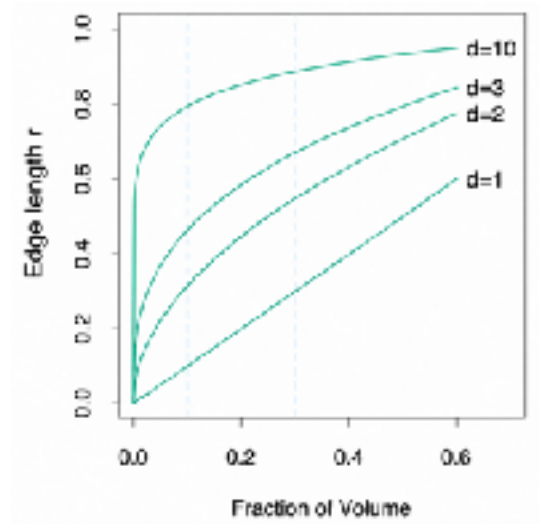
$X$  is uniformly distributed over  $[0,1]^d$ . What is  $P(X \in [0,r]^d)$ ?

How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?

# Curse of dimensionality, example 1



side length  $r$



# If  $d = 1$   $P = 0.3$



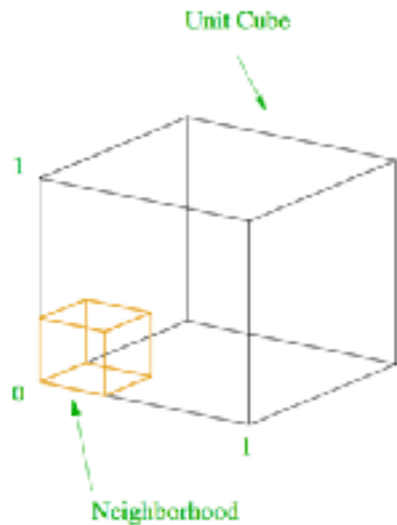
# If  $d = 2$

$X$  is uniformly distributed over  $[0,1]^d$ . What is  $P(X \in [0,r]^d)$

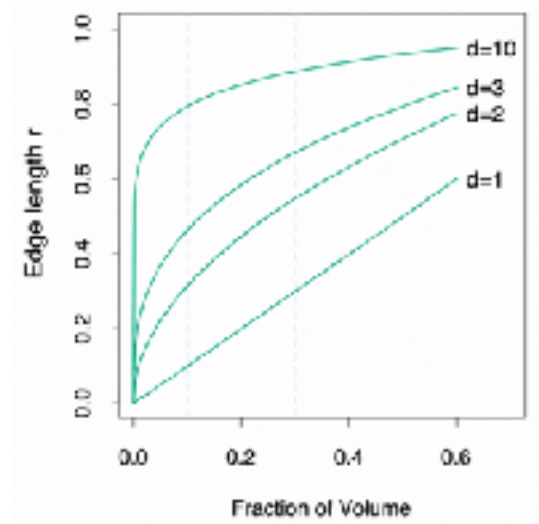
?

How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?

# Curse of dimensionality, example 1



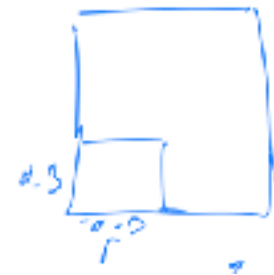
side length  $r$



# If  $d = 1$   $P = 0.3$



# If  $d = 2$

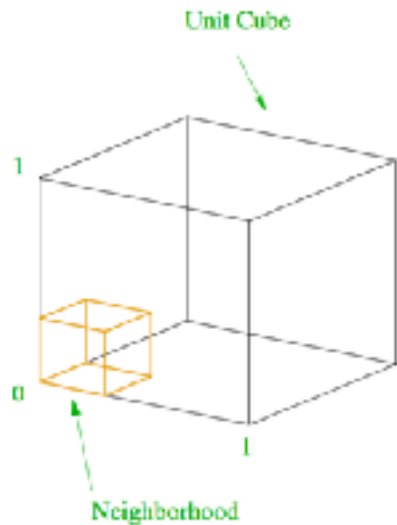


$X$  is uniformly distributed over  $[0,1]^d$ . What is  $P(X \in [0,r]^d)$

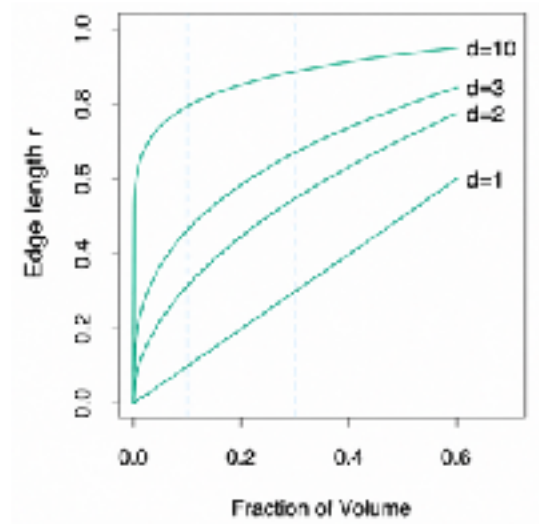
?

How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?

# Curse of dimensionality, example 1



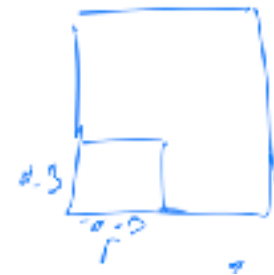
side length  $r$



# If  $d = 1$   $P = 0.3$



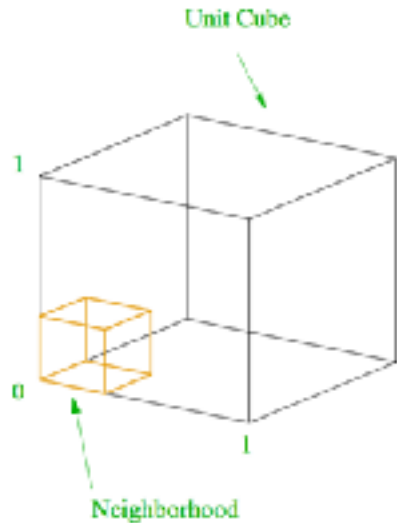
# If  $d = 2$   $P = 0.3^2 = 0.09$



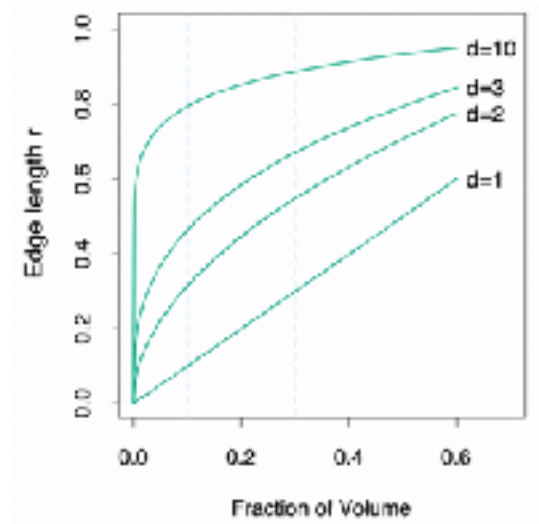
$X$  is uniformly distributed over  $[0,1]^d$ . What is  $P(X \in [0,r]^d)$ ?

How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?

# Curse of dimensionality, example 1



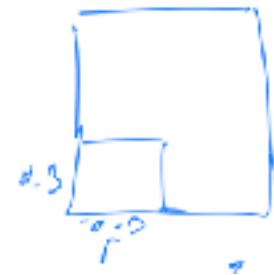
side length  $r$



# If  $d = 1$   $P = 0.3$



# If  $d = 2$   $P = 0.3^2 = 0.09$



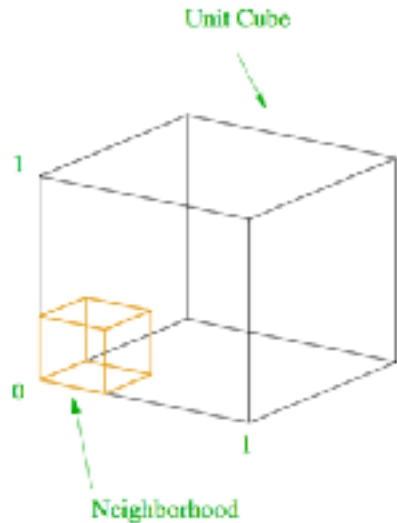
$X$  is uniformly distributed over  $[0,1]^d$ . What is  $P(X \in [0,r]^d)$

?

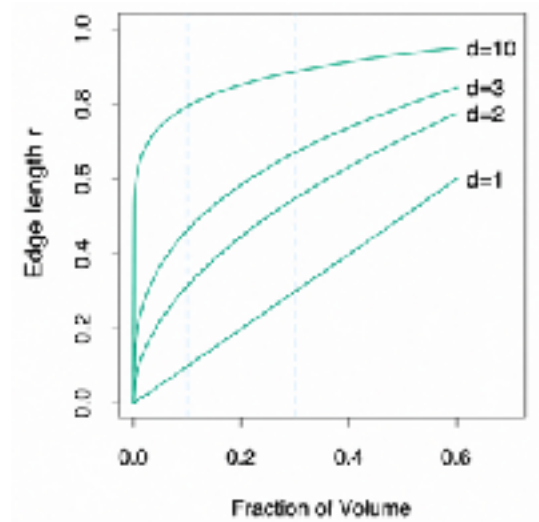
How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?

$$= \frac{1}{r^d}$$

# Curse of dimensionality, example 1



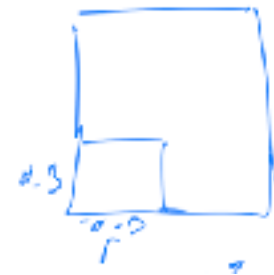
side length  $r$



# If  $d = 1$   $P = 0.3$



# If  $d = 2$   $P = 0.3^2 = 0.09$



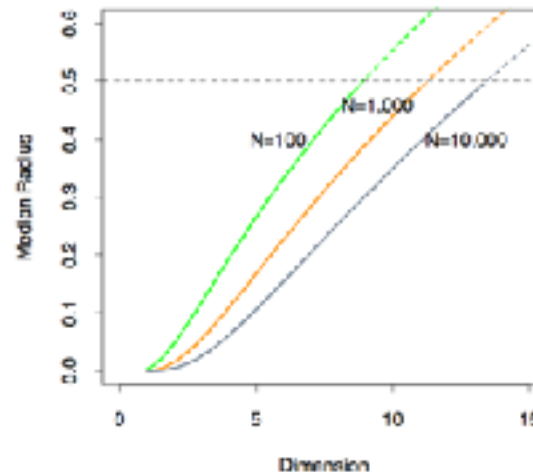
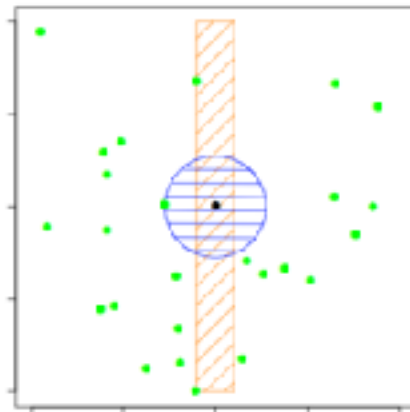
$X$  is uniformly distributed over  $[0, 1]^d$ . What is  $P(X \in [0, r]^d)$

$$= \frac{1}{r^d}$$

?  
How many samples do we need so that a nearest neighbor is within a cube of length  $r$ ?  
# Increases exponentially with dimension  $p$ !

# Curse of dimensionality, example 2

$\{X_i\}_{i=1}^n$  are uniformly distributed over  $[-.5, .5]^p$ .

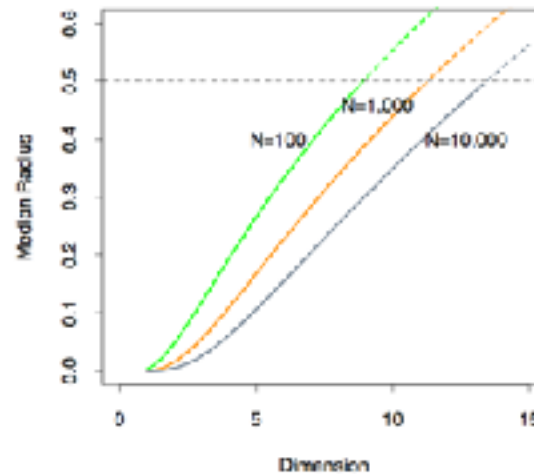
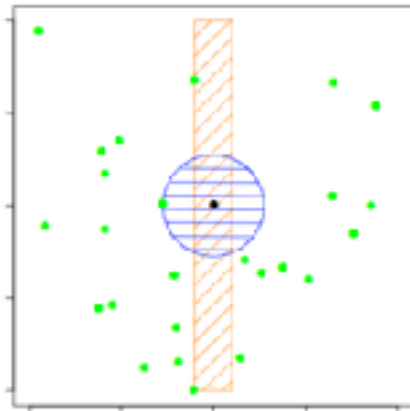


What is the median distance from a point at origin to its 1NN?

How many samples do we need so that a median Euclidean distance is within  $r$ ?

# Curse of dimensionality, example 2

$\{X_i\}_{i=1}^n$  are uniformly distributed over  $[-.5, .5]^p$ .



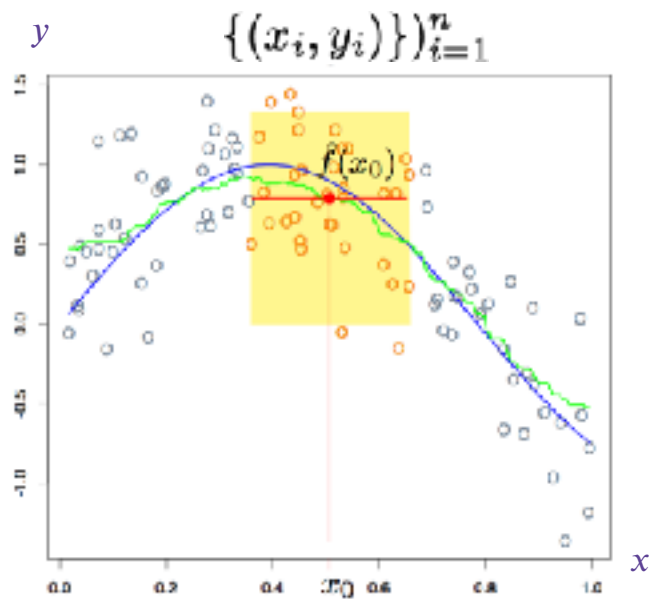
# High dimensional space is hard to cover

What is the median distance from a point at origin to its 1NN?

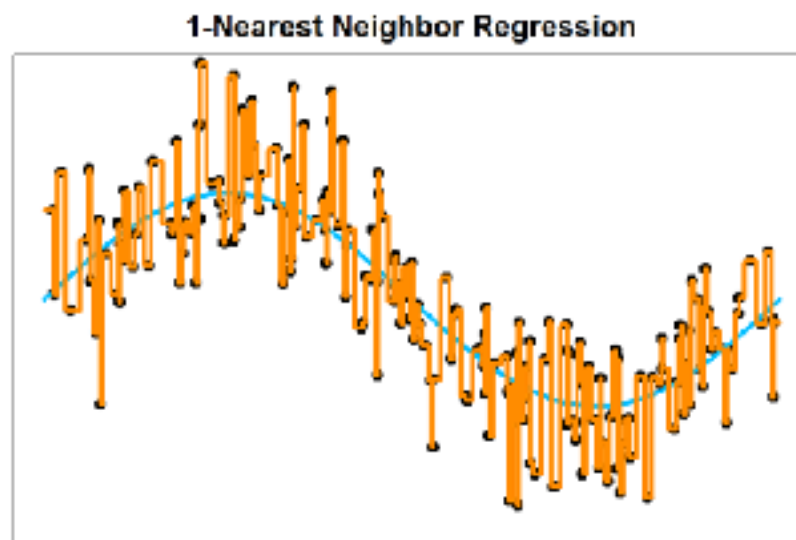
How many samples do we need so that a median Euclidean distance is within  $r$ ?

Slides are not updated beyond this point,  
download the latest ones from the course website

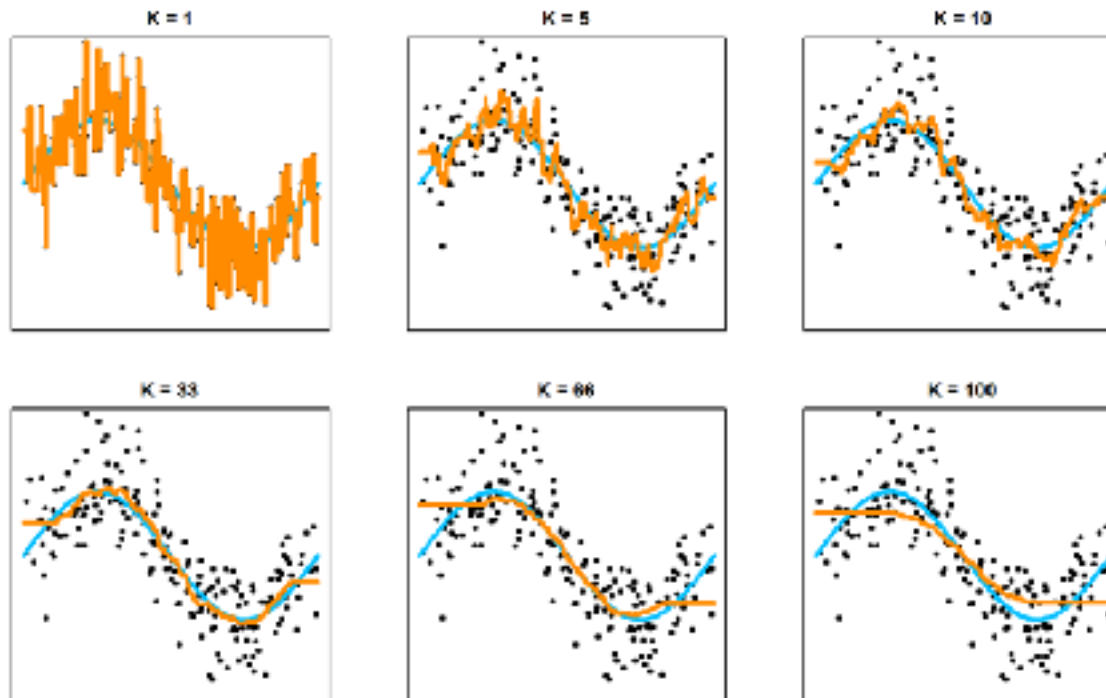
# Nearest neighbor regression



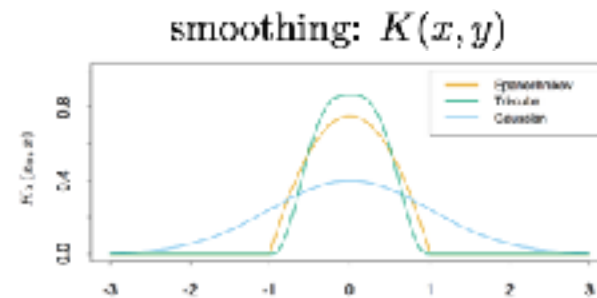
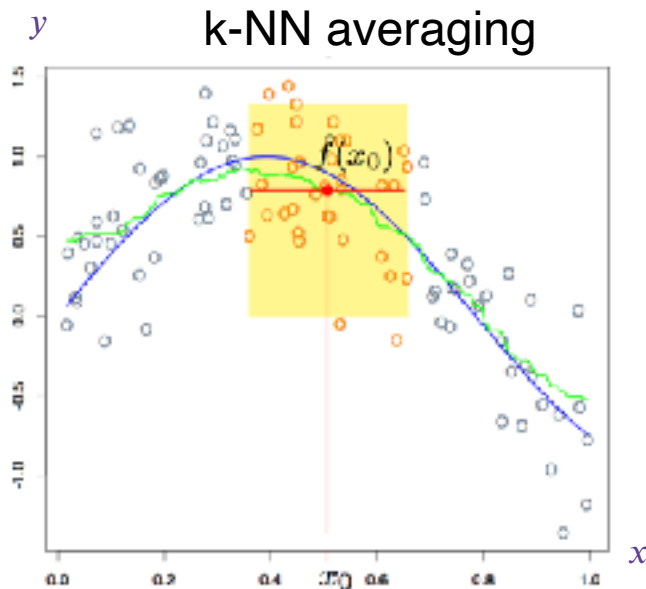
# Overfitting



# Bias vs variance



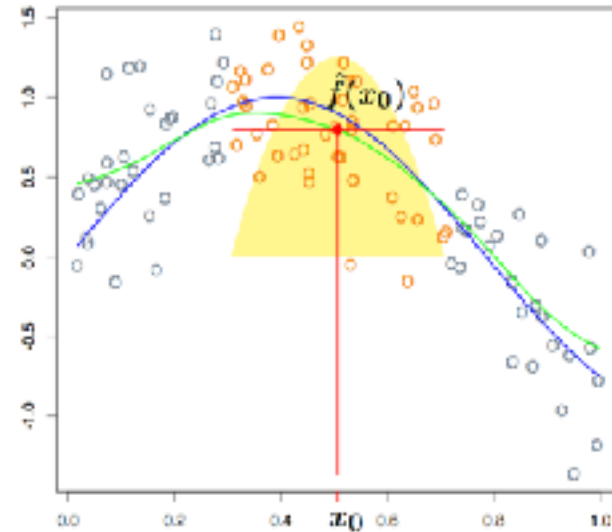
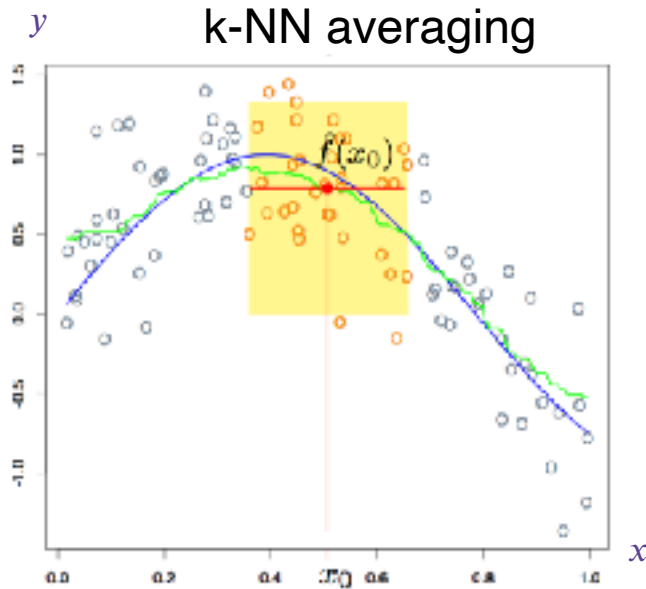
# Smoothed nearest neighbor regression



$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$

# Smoothed nearest neighbor regression

$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$

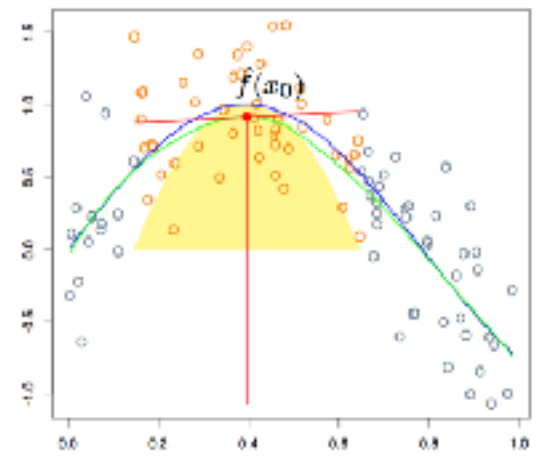
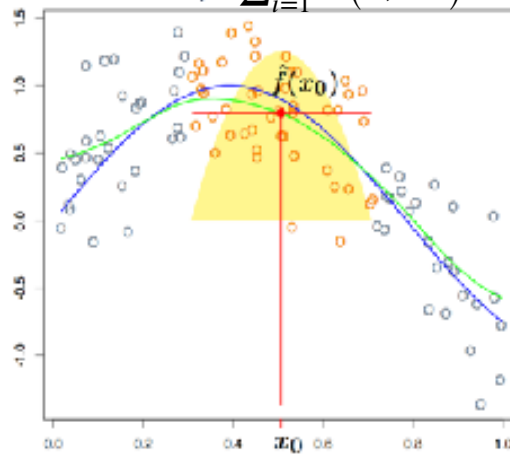
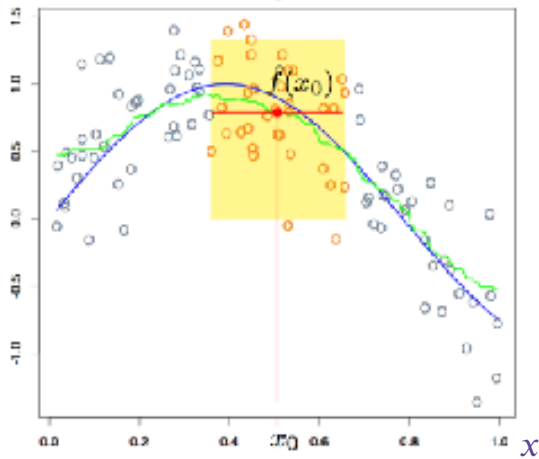


# Locally linear regression

$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$

y

k-NN averaging



Have we seen non-parametric methods before?

# Have we seen non-parametric methods before?

- Kernel methods are non-parametric:

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

# Have we seen non-parametric methods before?

- Kernel methods are non-parametric:

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

# parameters goes up with # data

# Have we seen non-parametric methods before?

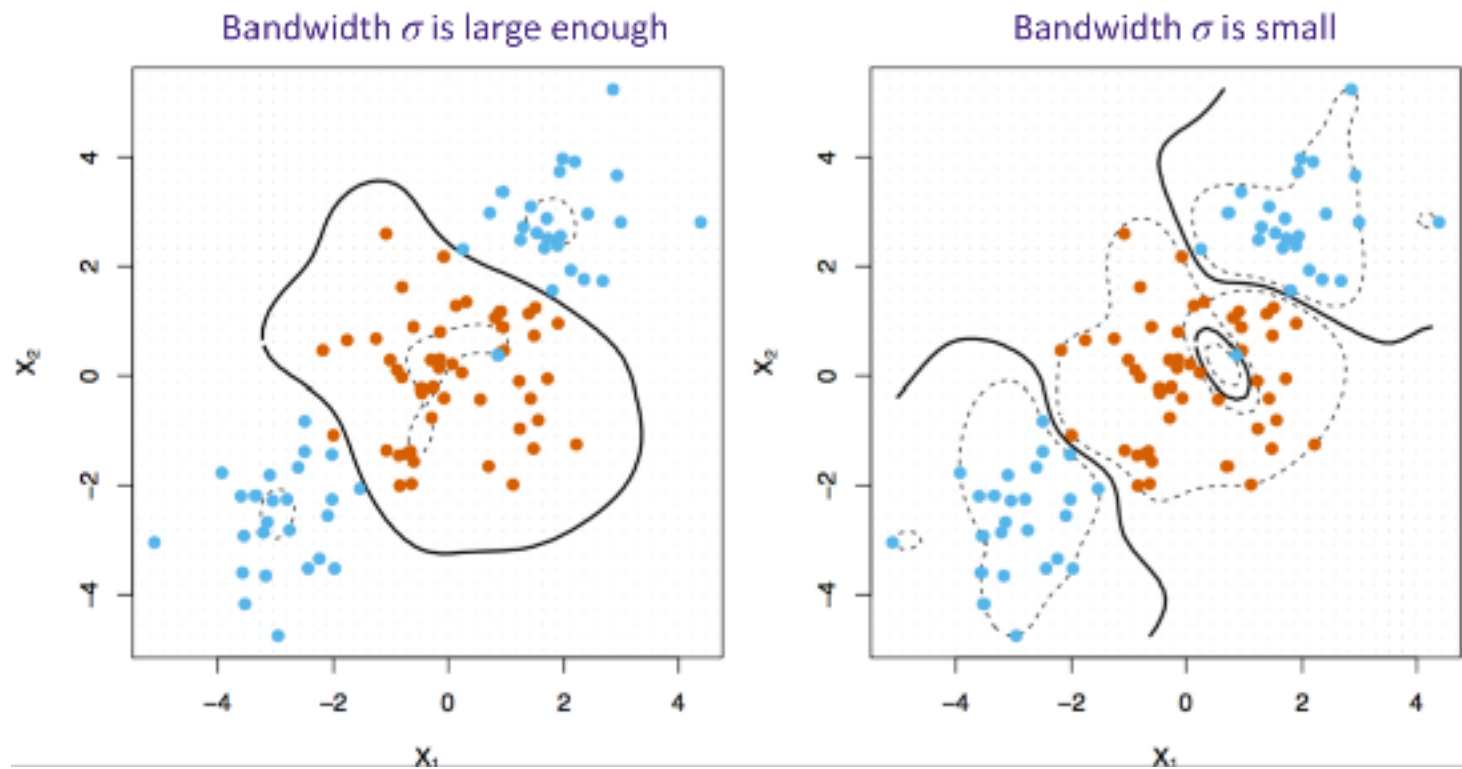
- Kernel methods are non-parametric:

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

# parameters goes up with # data

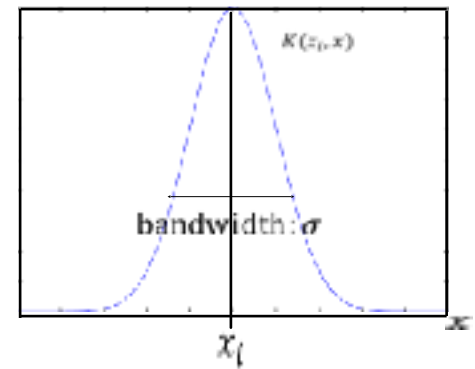
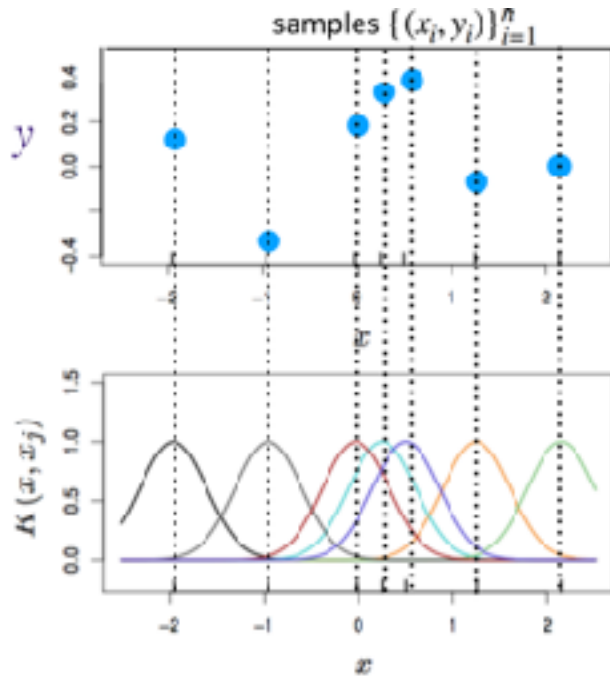
- Compare with ~~(smoothed) nearest neighbors:~~  
 $\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$

# The Radial Basis Function (RBF) kernel $\exp\left(-\frac{\|x^{(i)} - x\|_2^2}{2\sigma^2}\right)$



# The Radial Basis Function (RBF) kernel $\exp\left(-\frac{\|x^{(i)} - x\|_2^2}{2\sigma^2}\right)$

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$



# Kernel methods

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

# Kernel methods

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

- Can be seen as a soft, learned version of “nearest” neighbors

# Kernel methods

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

- Can be seen as a soft, learned version of “nearest” neighbors
- $K(x^{(i)}, x) = \phi(x^{(i)})^\top \phi(x)$  defines “similarity” between  $x^{(i)}$  and  $x$

# Kernel methods

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

- Can be seen as a soft, learned version of “nearest” neighbors
- $K(x^{(i)}, x) = \phi(x^{(i)})^\top \phi(x)$  defines “similarity” between  $x^{(i)}$  and  $x$
- How many parameters?

# Takeaways

# Takeaways

- k-NN is very simple to explain and implement

# Takeaways

- k-NN is very simple to explain and implement
- No training! But inference can still be computationally demanding.

# Takeaways

- k-NN is very simple to explain and implement
- No training! But inference can still be computationally demanding.
- You can use other forms of distance (not just Euclidean)

# Takeaways

- k-NN is very simple to explain and implement
- No training! But inference can still be computationally demanding.
- You can use other forms of distance (not just Euclidean)
- Smoothing and local linear regression can improve performance (at the cost of higher variance)

# Takeaways

- k-NN is very simple to explain and implement
- No training! But inference can still be computationally demanding.
- You can use other forms of distance (not just Euclidean)
- Smoothing and local linear regression can improve performance (at the cost of higher variance)
- With a lot of data, “local methods” have strong, simple theoretical guarantees

# Takeaways

- k-NN is very simple to explain and implement
- No training! But inference can still be computationally demanding.
- You can use other forms of distance (not just Euclidean)
- Smoothing and local linear regression can improve performance (at the cost of higher variance)
- With a lot of data, “local methods” have strong, simple theoretical guarantees
- Without a lot of data, neighborhoods aren’t “local” and methods suffer (curse of dimensionality)