

CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2020



HVERFORD
COLLEGE

- **Lab 5 due TODAY!**
 - Grace period until Wed night
 - No office hours today (were yesterday)
- Lab 6 posted today
- Welcome prospective students!

Outline for October 20

- Recap multi-class logistic regression
- Introduction to ensemble methods
- Bagging
- Random forests
- AdaBoost

Outline for October 20

- Recap multi-class logistic regression
- Introduction to ensemble methods
- Bagging
- Random forests
- AdaBoost

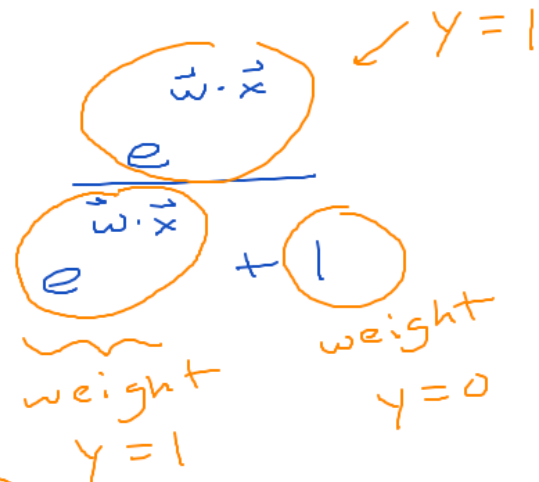
Multi-class Logistic Regression

K = num classes (political parties, blood groups, etc)

2 classes

$$h_{\vec{w}}(\vec{x}) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}} = \underbrace{\frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}}_{\text{prob } 1}$$

$1 - h(\vec{x}) = \text{prob } 0$



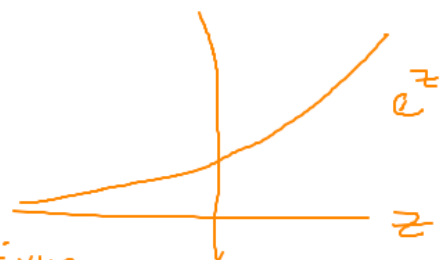
Generalize

$k = 1, 2, 3, \dots, K$

$$\vec{h}_w(\vec{x}) = \begin{bmatrix} p(y=1|\vec{x}) \\ p(y=2|\vec{x}) \\ \vdots \\ p(y=K|\vec{x}) \end{bmatrix}$$

$$= \frac{1}{\sum_{k=1}^K e^{\vec{w}^{(k)} \cdot \vec{x}}} \begin{bmatrix} e^{\vec{w}^{(1)} \cdot \vec{x}} \\ e^{\vec{w}^{(2)} \cdot \vec{x}} \\ \vdots \\ e^{\vec{w}^{(K)} \cdot \vec{x}} \end{bmatrix}$$

positive or negative
called logits



Probability distribution

$$W = \begin{bmatrix} | & | & \dots & | \\ \vec{w}^{(1)} & \vec{w}^{(2)} & \dots & \vec{w}^{(K)} \\ | & | & \dots & | \end{bmatrix} \quad (p+1) \times K$$

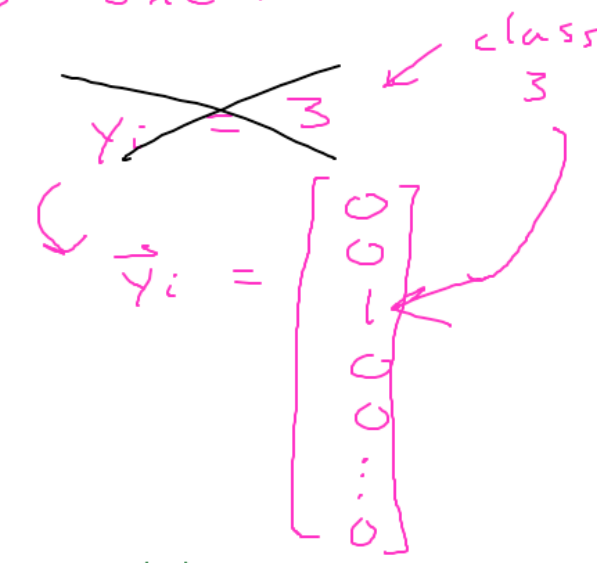
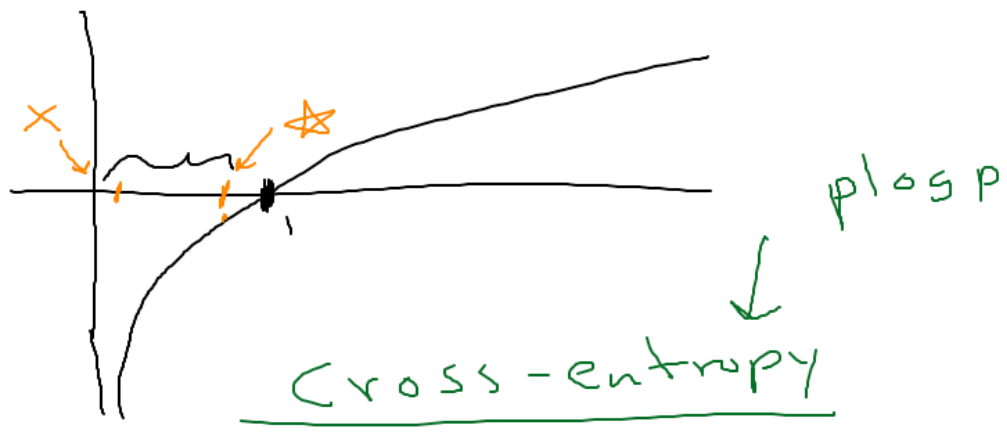
weight vector for each class

Multi-class Logistic Regression: cost function

still neg log likelihood $\{h^y (1-h)^{1-y}\}$

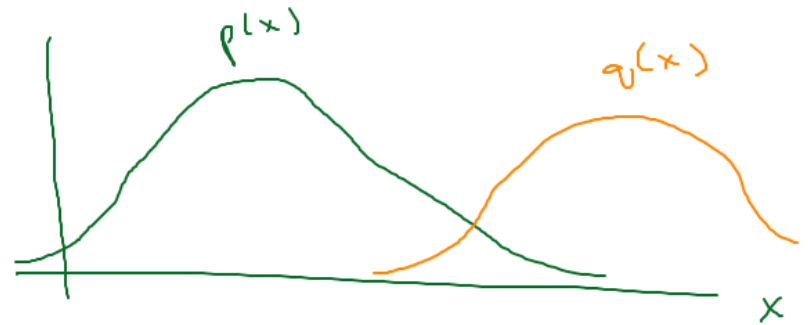
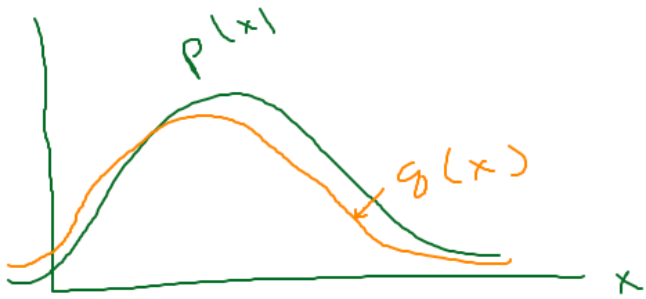
$$J(W) = - \sum_{i=1}^n \sum_{k=1}^K \underbrace{y_{ik}}_{\text{indicator}} \log \underbrace{p(y_i=k|\vec{x}_i)}_{\text{"one-hot"}}$$

$\frac{e}{\text{normalizer}}$



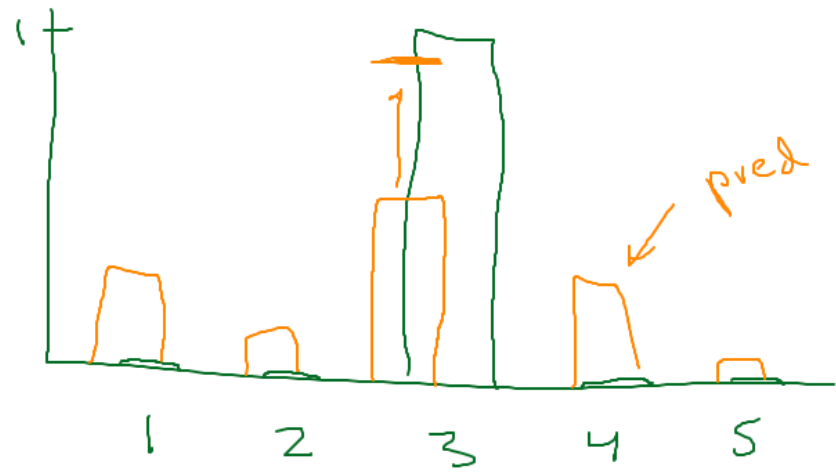
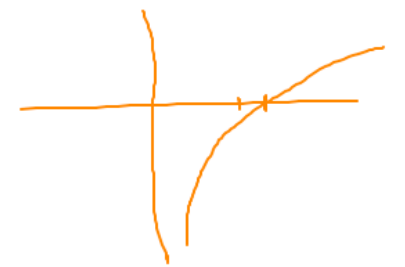
$$H(p, q) = - \sum_x p(x) \log q(x)$$

Cross Entropy



Log Regression :

$p(x) \rightarrow$ true
 $q(x) \rightarrow$ pred

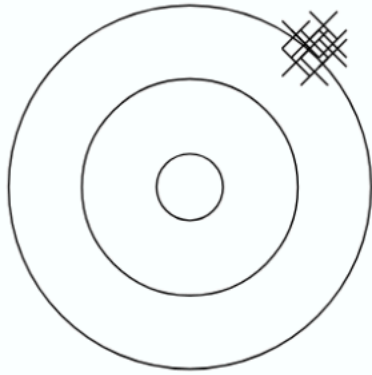


$K = S \approx 1$
 high
 $\sum_{k=1}^K \text{true} \cdot \log(\text{pred})$
 small
 1 for class 3

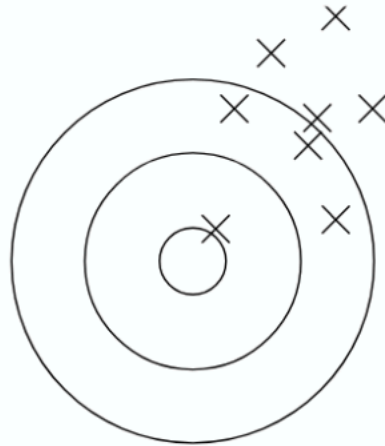
Outline for October 20

- Recap multi-class logistic regression
- Introduction to ensemble methods
- Bagging
- Random forests
- AdaBoost

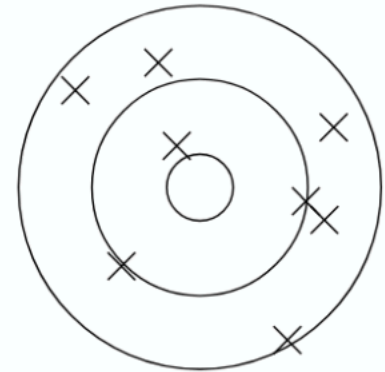
Quiz: recap bias and variance



A



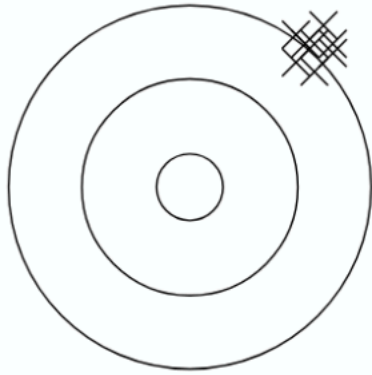
B



C

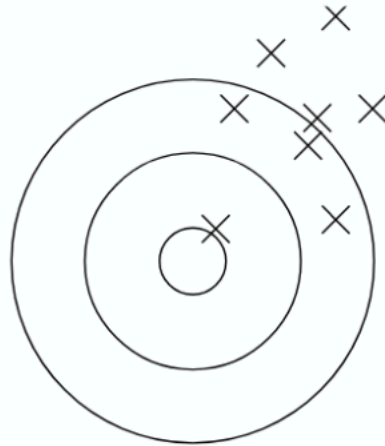
Label each picture with variance (high or low) and bias (high or low)

Quiz: recap bias and variance

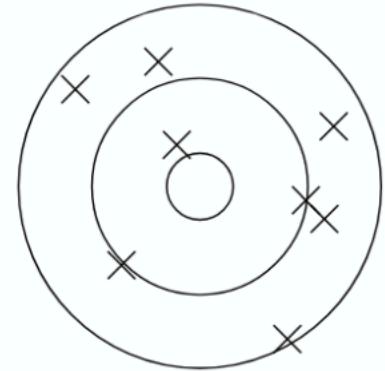


A

Variance: low
Bias: high



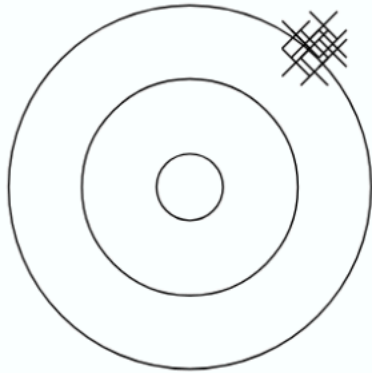
B



C

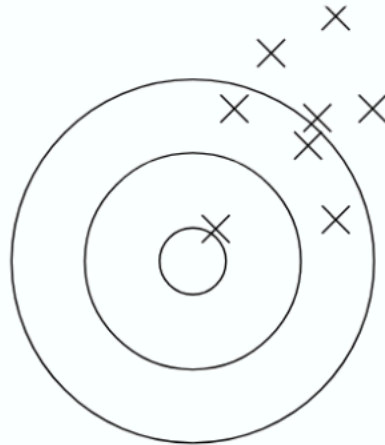
Label each picture with variance (high or low) and bias (high or low)

Quiz: recap bias and variance



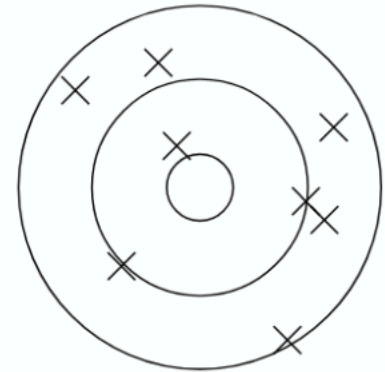
A

Variance: low
Bias: high



B

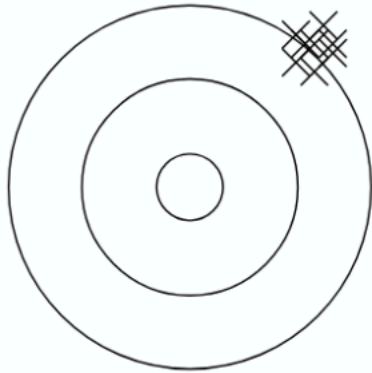
Variance: high
Bias: high



C

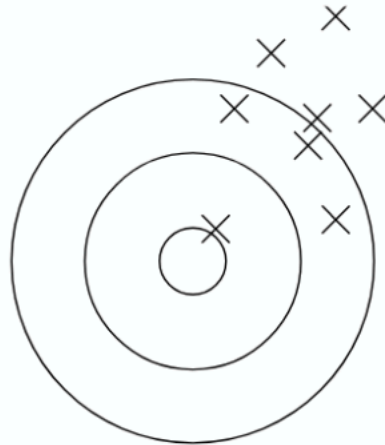
Label each picture with variance (high or low) and bias (high or low)

Quiz: recap bias and variance



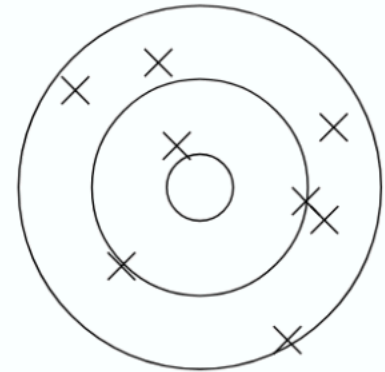
A

Variance: low
Bias: high



B

Variance: high
Bias: high

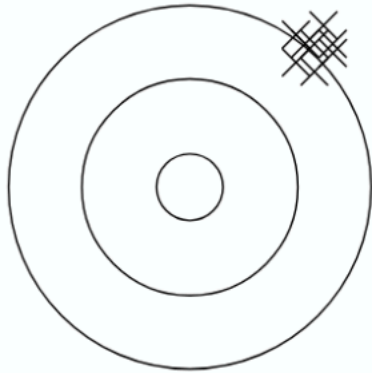


C

Variance: high
Bias: low

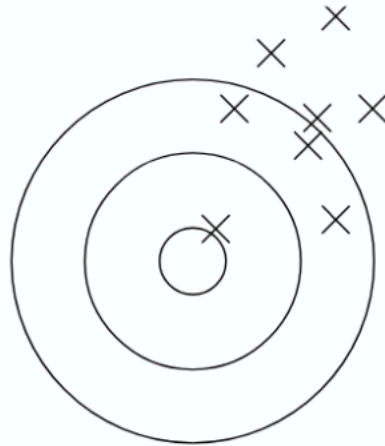
Label each picture with variance (high or low) and bias (high or low)

Quiz: recap bias and variance



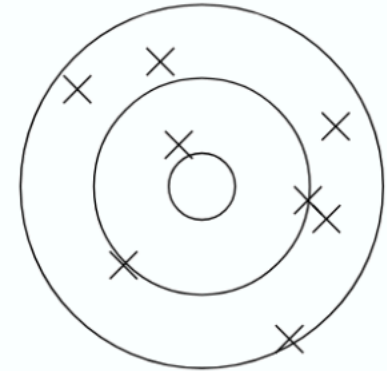
A

Variance: low
Bias: high



B

Variance: high
Bias: high



C

Variance: high
Bias: low

This is the type of classifier we want to average!

Label each picture with variance (high or low) and bias (high or low)

Ensemble Idea

- Average the results from several models with high variance and low bias
 - Important that models be diverse (don't want them to be wrong in the same ways)
- If n observations each have variance s^2 , then the mean of the observations has variance s^2/n (reduce variance by averaging!)

Learning Theory

Let H be the hypothesis space

Three sources of **limitations** for traditional classifiers:

- ❖ Statistical - H is too large relative to size of data
 - ❖ Many hypotheses can fit the data by chance
- ❖ Computational - H is too large to completely search for “best” model
- ❖ Representational - H is not expressive enough

Learning Theory

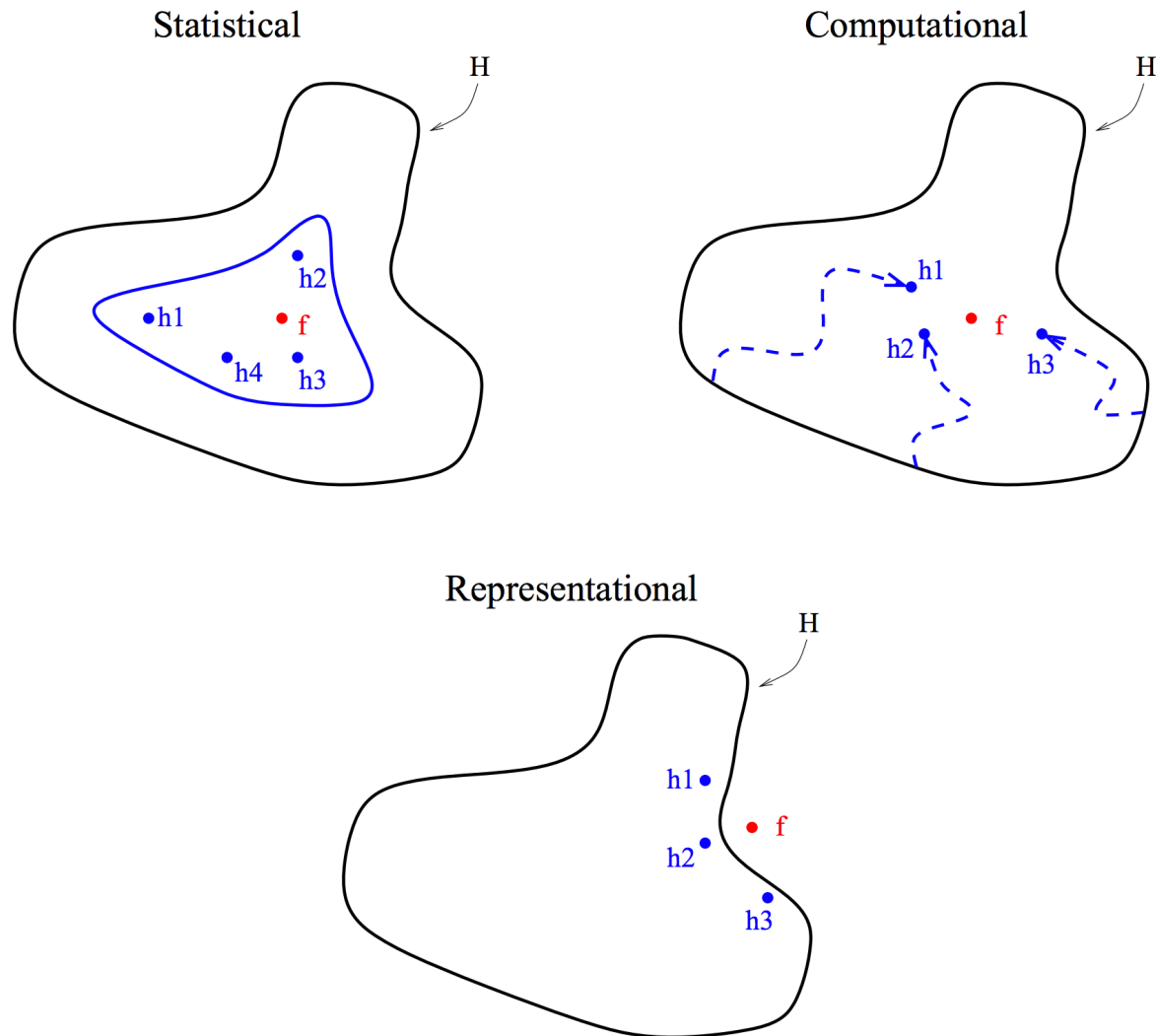
- ❖ Statistical: Average of unstable models (high variance) has more stability
- ❖ Computational: searching from multiple starting points is better approximation than one starting point
- ❖ Representational: sum of many models can represent more hypotheses than an individual model

Learning Theory

- ❖ Statistical: Average of unstable models (high variance) has more stability
- ❖ Computational: searching from multiple starting points is better approximation than one starting point
- ❖ Representational: sum of many models can represent more hypotheses than an individual model

Ensembles can address all 3!

Learning Theory

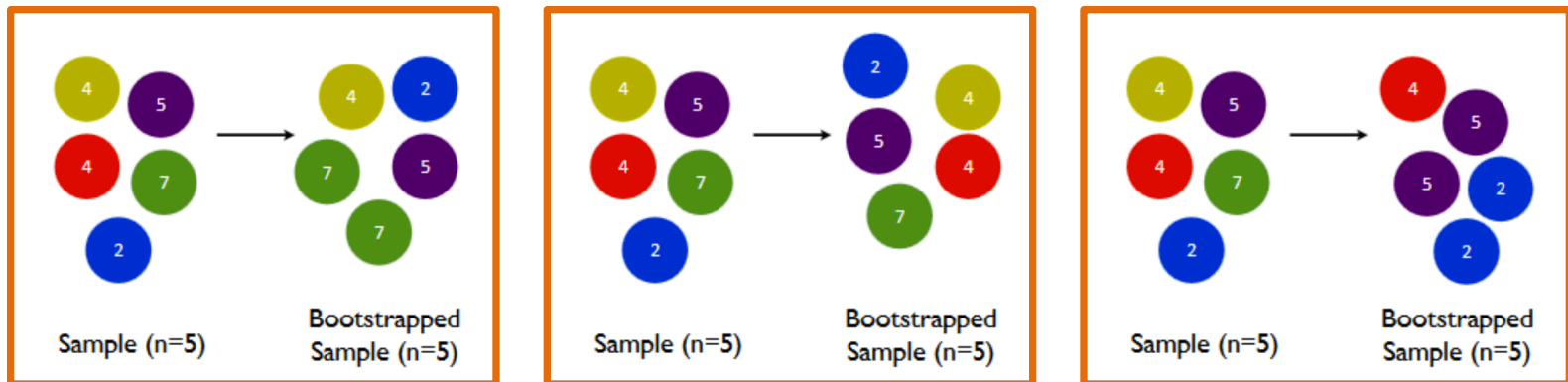


Outline for October 20

- Recap multi-class logistic regression
- Introduction to ensemble methods
- **Bagging**
- Random forests
- AdaBoost

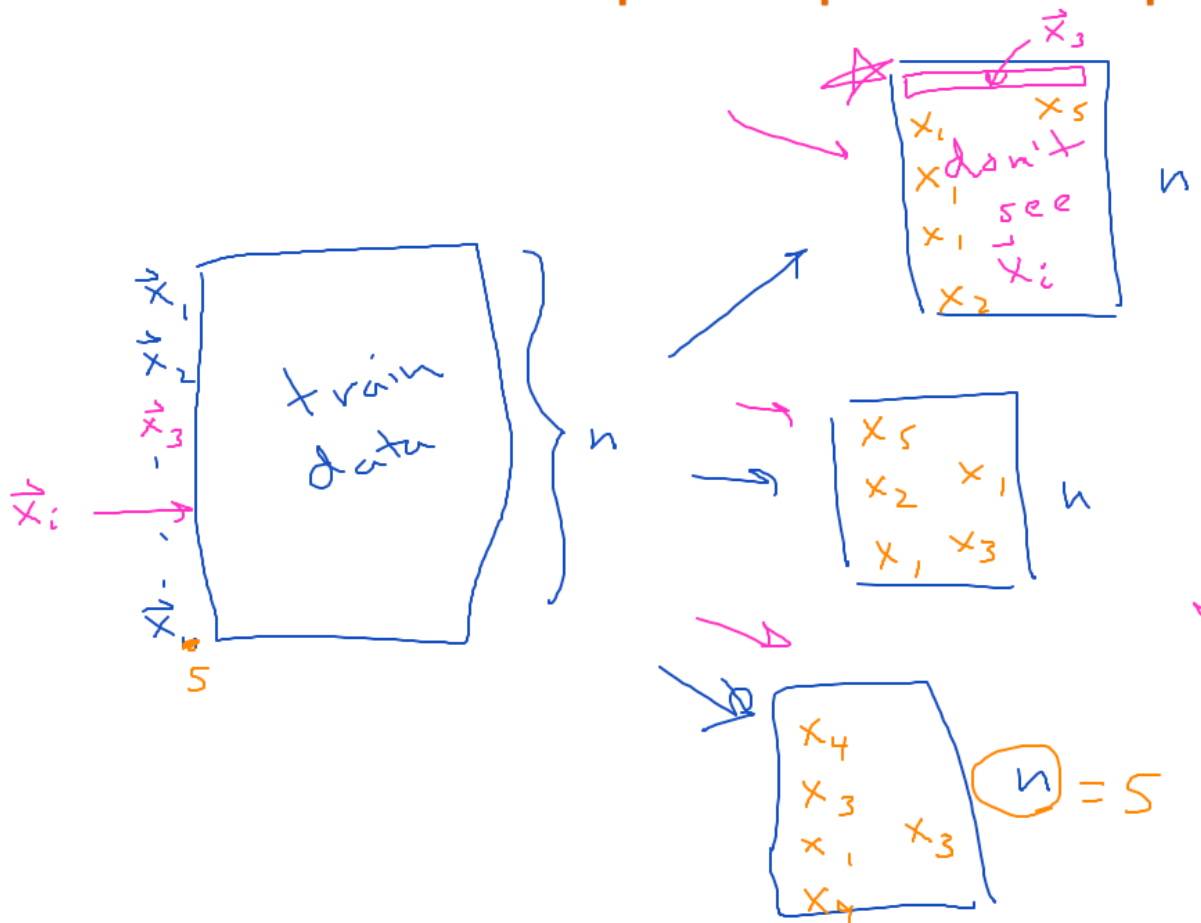
Bagging Algorithm

- ❖ Bagging = Bootstrap Aggregation [Brieman, 1996]
- ❖ *Bootstrap* (randomly sample with replacement) original data to create many different training sets
- ❖ Run base learning algorithm on each new data set independently



Desmond Ong, Stanford

Bootstrap: sample with replacement



why?

prob didn't choose a datapoint

$$\left(\frac{n-1}{n}\right)^n = \left(1 - \frac{1}{n}\right)^n$$

$$\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n \rightarrow e^x$$

$n \rightarrow \infty$
 $x_2 + x_5$
 not present

$$e^{-1} \approx 0.37$$

Notation for ensembles

- T = # models/classifiers (index t)
- x = test example (could be a vector)
- $X^{(t)}$ = bootstrap training dataset t
- $h^{(t)}(x)$ = hypothesis about x from model t
 $\in \{0, 1\}$ for binary problems
- r = probability of error of individual model
- R = # votes for wrong class

Bagging (Bootstrap Aggregation)

Train:

$T = 15 \rightarrow$

votes	Y
8	$\Rightarrow 0$
7	$\Rightarrow 1$

for t in range(T):

$$h(x) = 0$$

* create bootstrap sample $X^{(t)}$ of size n from training data

* train on $X^{(t)}$ to get model $h^{(t)}$

Test:

for each test example x

$$h(x) = \underset{y \in \{0,1\}}{\operatorname{argmax}} \sum_{t=1}^T \mathbb{1}(h^{(t)}(x) = y)$$

Bagging: motivating example

$T = 3$

test

(x, y)

$r =$ prob one model is wrong

$R =$ # models that are wrong

$$R = \sum_{t=1}^T \mathbb{1}(h^{(t)}(x) = \tilde{y})$$

wrong class

$$P(R=k) = \binom{T}{k} r^k (1-r)^{T-k}$$

wrong right

incorrect overall

$$P(R > \frac{T}{2}) = \sum_{k=\frac{T+1}{2}}^T \binom{T}{k} r^k (1-r)^{T-k} \quad \binom{4}{2} = \frac{4 \cdot 3}{2}$$

more than half models

wrong

$$\binom{T}{k} = \frac{T!}{k!(T-k)!}$$

$h^{(1)}$	$h^{(2)}$	$h^{(3)}$
0	0	0
0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

$r = \frac{1}{4}$

T choose k

all possibilities

Bagging: motivating example

$T=3$

test

(x, y)

$r =$ prob one model is wrong

$R =$ # models that are wrong

$$R = \sum_{t=1}^T \mathbb{I}(h^{(t)}(x) = \tilde{y}) \leftarrow \text{wrong class}$$

$h^{(1)}$	$h^{(2)}$	$h^{(3)}$
0	0	0
0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

$r = \frac{1}{4}$

$$P(R=k) = \binom{T}{k} r^k (1-r)^{T-k}$$

r^k wrong
 $(1-r)^{T-k}$ right

incorrect overall

$$P(R > \frac{T}{2}) = \sum_{k=\frac{T+1}{2}}^T \binom{T}{k} r^k (1-r)^{T-k}$$

more than half models wrong

Example

$T=3$
 $r = \frac{1}{4}$
all wrong

$$P(R > 1.5) = 3 \left(\frac{1}{4}\right)^2 \left(\frac{3}{4}\right) + \left(\frac{1}{4}\right)^3$$

$= 0.16$

84% accuracy!

all possibilities

Outline for October 20

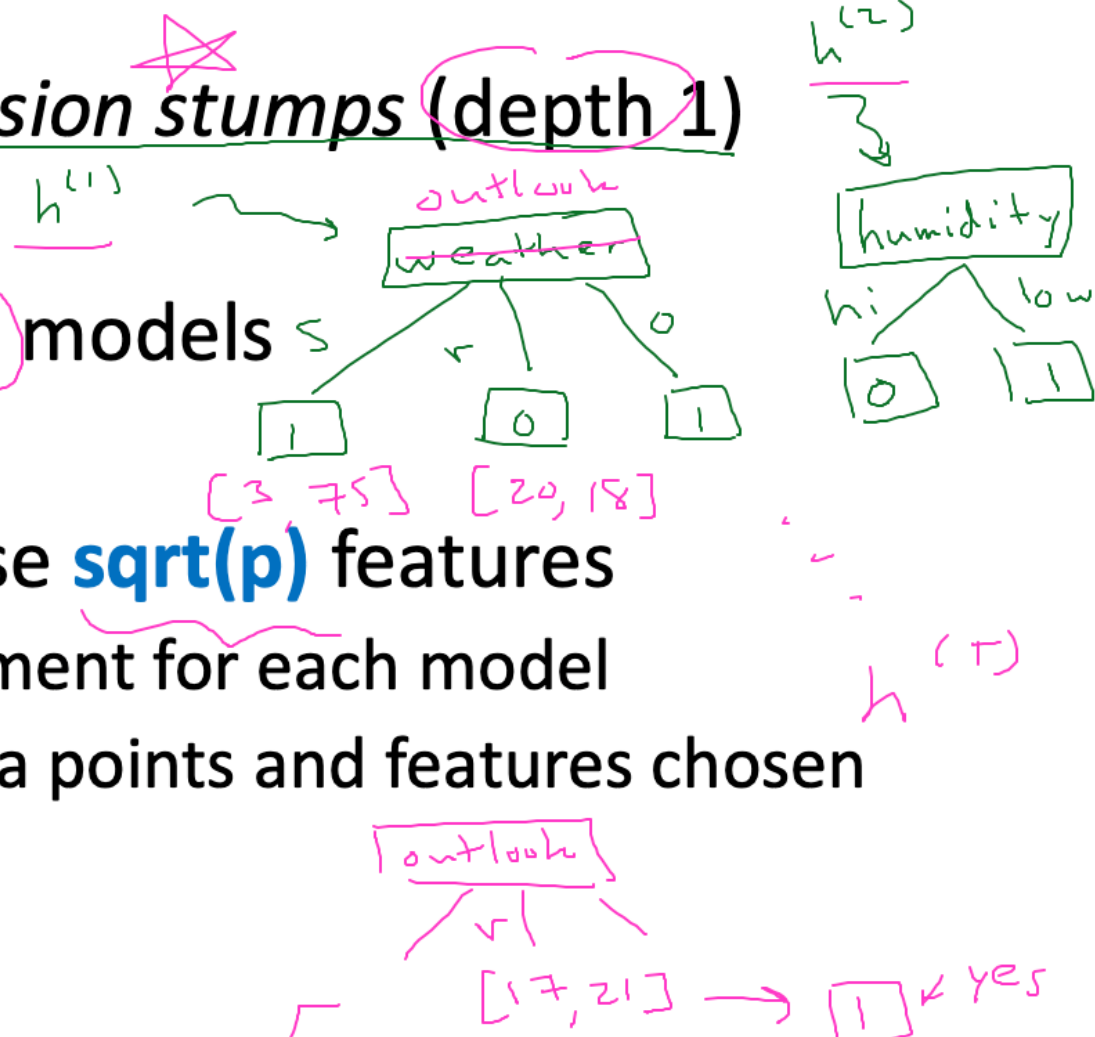
- Recap multi-class logistic regression
- Introduction to ensemble methods
- Bagging
- **Random forests**
- AdaBoost

Random Forests

- Idea: choose a different **subset of features** for every classifier t
- Typically use *decision stumps* (depth 1)
- Goal: decorrelate models
- In practice: choose **\sqrt{p}** features
 - Without replacement for each model
 - Every model: data points and features chosen independently

Random Forests

- Idea: choose a different subset of features for every classifier t
- Typically use decision stumps (depth 1)
- Goal: decorrelate models
 $1 = \text{yes}$
 $2 = \text{no}$
- In practice: choose \sqrt{p} features
 - Without replacement for each model
 - Every model: data points and features chosen independently



Outline for October 20

- Recap multi-class logistic regression
- Introduction to ensemble methods
- Bagging
- Random forests
- **AdaBoost**

NEXT TIME!