# CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2020

HAVERFORD
COLLEGE

# Admin

- Lab 2 due **THURSDAY night**
  - Grace period til Friday at 5
  - Continuous features optional
  - Information gain optional (can do cond. entropy instead)

- Office hours today **4:30-6pm**

- **Reading for Friday**
  - Duame 7.6 (2+ pages)
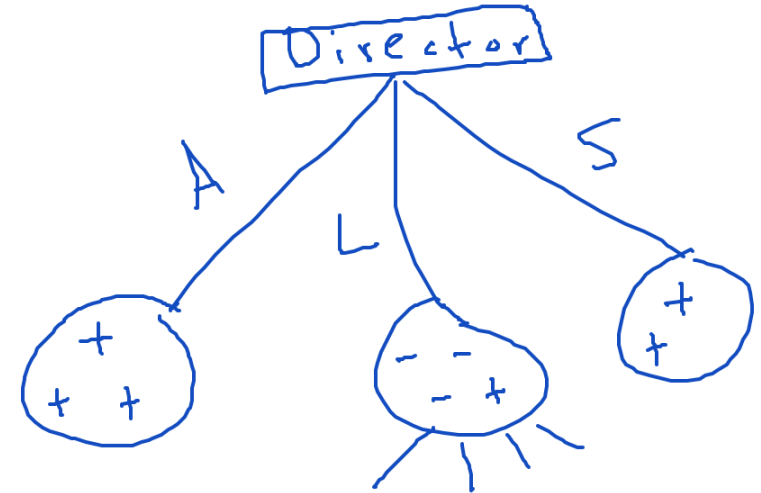  - ISL 59-63 (4+ pages)

# Outline for September 22

- Finish Decision Trees (recap continuous features)

- Learning problem so far + terminology

- Bias-Variance tradeoff

- Linear regression

# Linear Regression Goals

- Regression as a way to study *expected loss* and the *bias-variance tradeoff*

- Review matrix algebra and expected values

- As an introduction to optimization (specifically *stochastic gradient descent*)

# Outline for September 22

- Finish Decision Trees (recap continuous features)


- Learning problem so far + terminology


- Bias-Variance tradeoff


- Linear regression

| Movie | Type | Length | Director | Famous actors | Liked? |
|-------|------|--------|----------|---------------|--------|
| m1 | Comedy | Short | Adamson | No | Yes |
| m2 | Animated | Short | Lasseter | No | No |
| m3 | Drama | Medium | Adamson | No | Yes |
| m4 | Animated | Long | Lasseter | Yes | No |
| m5 | Comedy | Long | Lasseter | Yes | No |
| m6 | Drama | Medium | Singer | Yes | Yes |
| m7 | Animated | Short | Singer | No | Yes |
| m8 | Comedy | Long | Adamson | Yes | Yes |
| m9 | Drama | Medium | Lasseter | No | Yes |

← Y

Handout 3

root



$P(Li = yes) = \dfrac{2}{3}$

$H(Li) = -\left(\dfrac{2}{3}\log\dfrac{2}{3} + \dfrac{1}{3}\log\dfrac{1}{3}\right) = 0.92$

$H(Li \mid T) = 0.61$
$H(Li \mid Le) = 0.61$
$H(Li \mid D) = 0.36$  **MIN ENTROPY**
$H(Li \mid F) = 0.85$

$Gain(Li, T) = 0.92 - 0.61 = 0.31$
$Gain(Li, Le) = 0.92 - 0.61 = 0.31$
$Gain(Li, D) = 0.92 - 0.36 = 0.56$  **MAX INFO GAIN**
$Gain(Li, F) = 0.92 - 0.85 = 0.07$

| Movie | Type | Length | Director | Famous actors | Liked? |
|-------|------|--------|----------|---------------|--------|
| m1 | Comedy | Short | Adamson | No | Yes |
| m2 | Animated | Short | Lasseter | No | No |
| m3 | Drama | Medium | Adamson | No | Yes |
| m4 | Animated | Long | Lasseter | Yes | No |
| m5 | Comedy | Long | Lasseter | Yes | No |
| m6 | Drama | Medium | Singer | Yes | Yes |
| m7 | Animated | Short | Singer | No | Yes |
| m8 | Comedy | Long | Adamson | Yes | Yes |
| m9 | Drama | Medium | Lasseter | No | Yes |

$$0.36$$

$$-\left(\frac{1}{4}\log\frac{1}{4} + \frac{3}{4}\log\frac{3}{4}\right)$$

weighted avg

$$H(Li \mid D) = \frac{3}{9} H(Li \mid D=A) + \frac{4}{9} \cdot H(Li \mid D=L) + \frac{2}{9} H(Li \mid D=S)$$

$$H(Li \mid D=A) = -P(Y \mid D=A)\log P(Y \mid D=A)$$
$$-P(N \mid D=A)\log P(N \mid D=A)$$

und

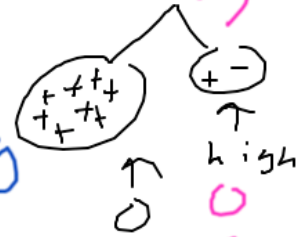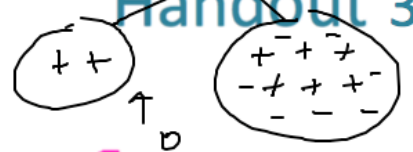$$P(Y \mid D=A) = \frac{P(Y \text{ and } D=A)}{P(D=A)} = \frac{3}{3} = 1$$

| Movie | Type | Length | Director | Famous actors | Liked? | Y |
|-------|----------|--------|----------|---------------|--------|---|
| m1 | Comedy | Short | Adamson | No | Yes | |
| m2 | Animated | Short | Lasseter | No | No | |
| m3 | Drama | Medium | Adamson | No | Yes | |
| m4 | Animated | Long | Lasseter | Yes | No | |
| m5 | Comedy | Long | Lasseter | Yes | No | |
| m6 | Drama | Medium | Singer | Yes | Yes | |
| m7 | Animated | Short | Singer | No | Yes | |
| m8 | Comedy | Long | Adamson | Yes | Yes | |
| m9 | Drama | Medium | Lasseter | No | Yes | |

$$0.36$$

$$-\left(\frac{1}{4}\log\frac{1}{4} + \frac{3}{4}\log\frac{3}{4}\right)$$

weighted avg

$$H(Li \mid D) = \frac{3}{9} H(Li \mid D=A) + \frac{4}{9} \cdot H(Li \mid D=L) + \frac{2}{9} H(Li \mid D=S)$$

output Y

$$H(Li \mid D=A) = -P(Yes \mid D=A)\log P(Yes \mid D=A)$$
$$-P(No \mid D=A)\log P(No \mid D=A)$$

$\{Yes, No\}$

$$P(Y \mid D=A) = \frac{P(Y \text{ and } D=A)}{P(D=A)} = \frac{3}{3} = 1$$

# Outline for September 22

- **Finish Decision Trees (recap continuous features)**

  main

  $dtree = DecisionTree(train)$

- Learning problem so far + terminology

  "Sun"
  "rain"

  children
  key: string
  value: DTree

  newpart = Partition( )
  child = DTree (newpart, d+1)
  self.children[v] = child

- Bias-Variance tradeoff

- Linear regression

# Continuous Features

(do this for the TRAIN only!)

| X | Y |
|---|---|
| 10 | Y |
| 7 | Y |
| 8 | N |
| 3 | Y |
| 7 | N |
| 12 | Y |
| 2 | Y |

1) Sort examples based on given feature

2   3   7   7   8   10   12
Y   Y   Y   N   N   Y   Y

2) Different label with same feature value, collapse to "None"

2   3   7   8   10   12
Y   Y   None   N   Y   Y

3) Whenever label changes, make a feature (use avg)

$x \le 5$

$x \le 7.5$

$x \le 9$

$x \le 5$

F
F
F

T
F
F
T

# Continuous Features (pair exercise)

(do this for the TRAIN only!)

3 new cols

| temp | Y |
|------|---|
| 80 | Y |
| 48 | Y |
| 60 | N |
| 48 | Y |
| 40 | N |
| 48 | Y |
| 90 | Y |

x ≤ 44

F, F, F, F, T, F, F

1) Sort examples based on feature "temp"

40    48    48    48    60    80    90
N     Y     Y     N     X     Y     Y
                                N

2) Different label with same feature value, collapse to "None"

40    48    60    80    90
N     None  X     Y     X Y
      Y     N

3) Whenever label changes, make a feature (use avg)

x ≤ 44      x ≤ 54      x ≤ 70

# Any other Lab 2 questions?

# Outline for September 22

- Finish Decision Trees (recap continuous features)

- Learning problem so far + terminology

- Bias-Variance tradeoff

- Linear regression

# Loss Functions

- E.g., zero-one loss

  - Simple accuracy - is prediction right?

  - For binary or multi-class prediction

$$l(y, \hat{y}) = \begin{cases} 0 & \text{if } y = \hat{y} \\ 1 & \text{otherwise} \end{cases}$$

- E.g., squared loss

  - For regression

$$l(y, \hat{y}) = (y - \hat{y})^2$$

- Absolute loss (also for regression)

$$\ell(y, \hat{y}) = |y - \hat{y}|$$

# Formalizing the learning problem

❖ Given:

   ❖ Loss function, $\ell$

   ❖ A sample of data $D$ from an unknown distribution of all data $\mathcal{D}$

   ❖ A hypothesis space $H = \{h | h : X \to Y\}$

❖ Do:

   ❖ Find a function $f(X) \to y$ that

   ❖ minimize error over $\mathcal{D}$ with respect to $\ell$

# Why might learning fail?

# Inductive Bias
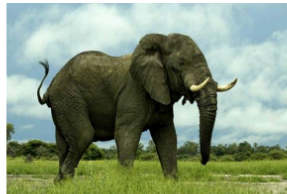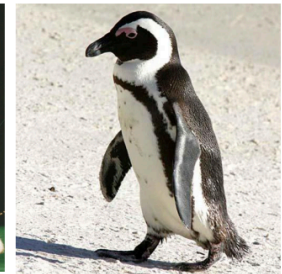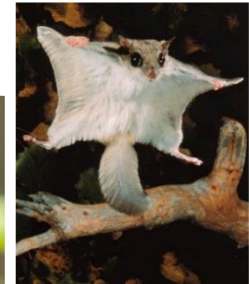
Training Data

class A

class B

Testing Data

# Inductive Bias
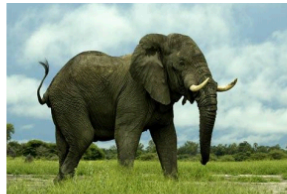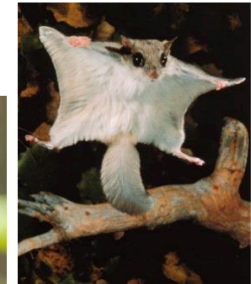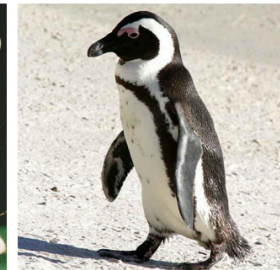
Training Data

class A
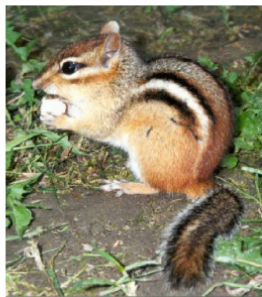
class B

Testing Data

A

A

B

B

A: "fly"
B: "no fly"

# Inductive Bias
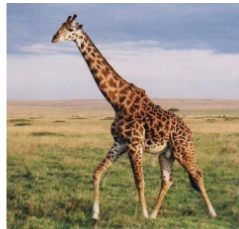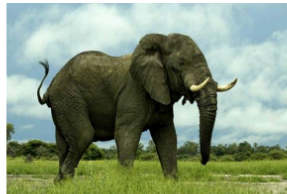
Training Data

class A

class B

Testing Data

A

B

B

A

A: "bird"
B: "mammal"

# Why might learning fail?

- Noise in the training data
  - Typos in a restaurant review


- Available features are insufficient
  - x-ray does not capture the medical issue


- "Correct" prediction is up to interpretation
  - Parental controls on web content


- Learning algorithm cannot cope with the data

# Hyperparameters

- Difficult to define precisely, but typically a parameter that controls other parameters

- What is one hyperparameter in decision trees?

<p style="text-align:center; color:blue;">Max depth!</p>

- We can't choose hyperparameters via test data (breaks cardinal rule!)

- But we can use *validation data*

# General approach to training

1. Split your data into 70% training data, 10% development data and 20% test data.  (validation data)

2. For each possible setting of your hyperparameters:

    (a) Train a model using that setting of hyperparameters on the training data.

    (b) Compute this model's error rate on the development data.

3. From the above collection of models, choose the one that achieved the lowest error rate on development data.

4. Evaluate that model on the test data to estimate future test performance.

# Outline for September 22

- Finish Decision Trees (recap continuous features)

- Learning problem so far + terminology

- Bias-Variance tradeoff
  *whiteboard*
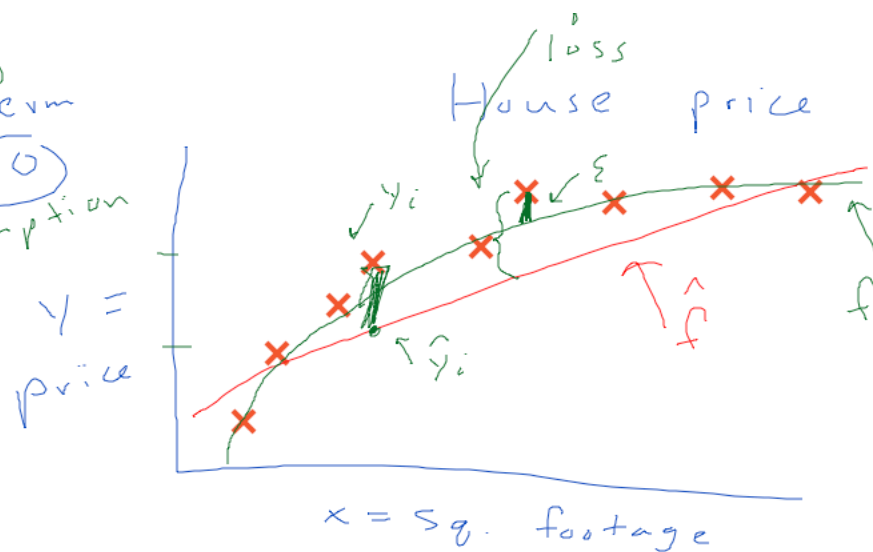
- Linear regression

# Regression setup and Expected Value

Regression setup

model: true $y = f(x) + \varepsilon$    (noise) error term    mean 0  assumption

$\hat{f}$ = estimate of $f$

$\hat{y} = \hat{f}(x)$    prediction

loss  House price

$y = $ price    $x = sq.$ footage

GOAL    $\ell(y, \hat{y}) = (y - \hat{y})^2$

$\underline{\underline{E[(y - \hat{y})^2]}}$    expected loss

$E(D) = \frac{1}{10}(1 + 2 \cdots 5) + \frac{1}{2} \cdot 6$
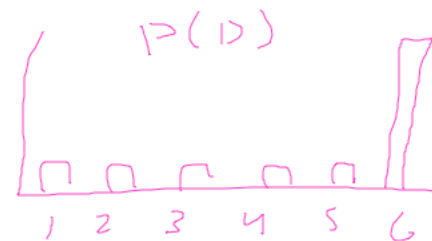
$= 4.5$

expected value = weighted avg

$E[X] = \sum p(x = v) \cdot v$
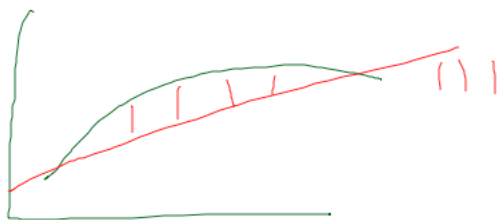
$P(D = 6) = \frac{1}{2}$

$P(D \neq 6) = \frac{1}{10}$

$P(D)$

1  2  3  4  5  6

# Compute Expected Loss

$$E\left[(y-\hat{y})^2\right] = E\left[\underbrace{(\underbrace{y-f}_{\varepsilon} + \underbrace{f-\hat{f}}_{0})^2}_{}\right] \qquad \hat{f} = \hat{y}$$

a, b labels above $(y-f)$ and $(f-\hat{f})$

$\hat{f} \overset{\uparrow}{=} \hat{y}$
$(x)$

$$= \underbrace{Var(\varepsilon)}_{\substack{noise \\ irreducible \\ error}} \overset{a^2}{} + \underbrace{E\left[(f-\hat{f})^2\right]}_{\substack{reducible \\ error}} \overset{b^2}{}$$

$(a+b)^2$
$= a^2 + b^2 + 2ab$   0

$Var(x) = E\left[(x-\mu)^2\right]$

$$E\left[(f-\hat{f})^2\right] = E\left[\left(\underbrace{f - E[\hat{f}]}_{bias} + \underbrace{E[\hat{f}] - \hat{f}}_{0}\right)^2\right]$$
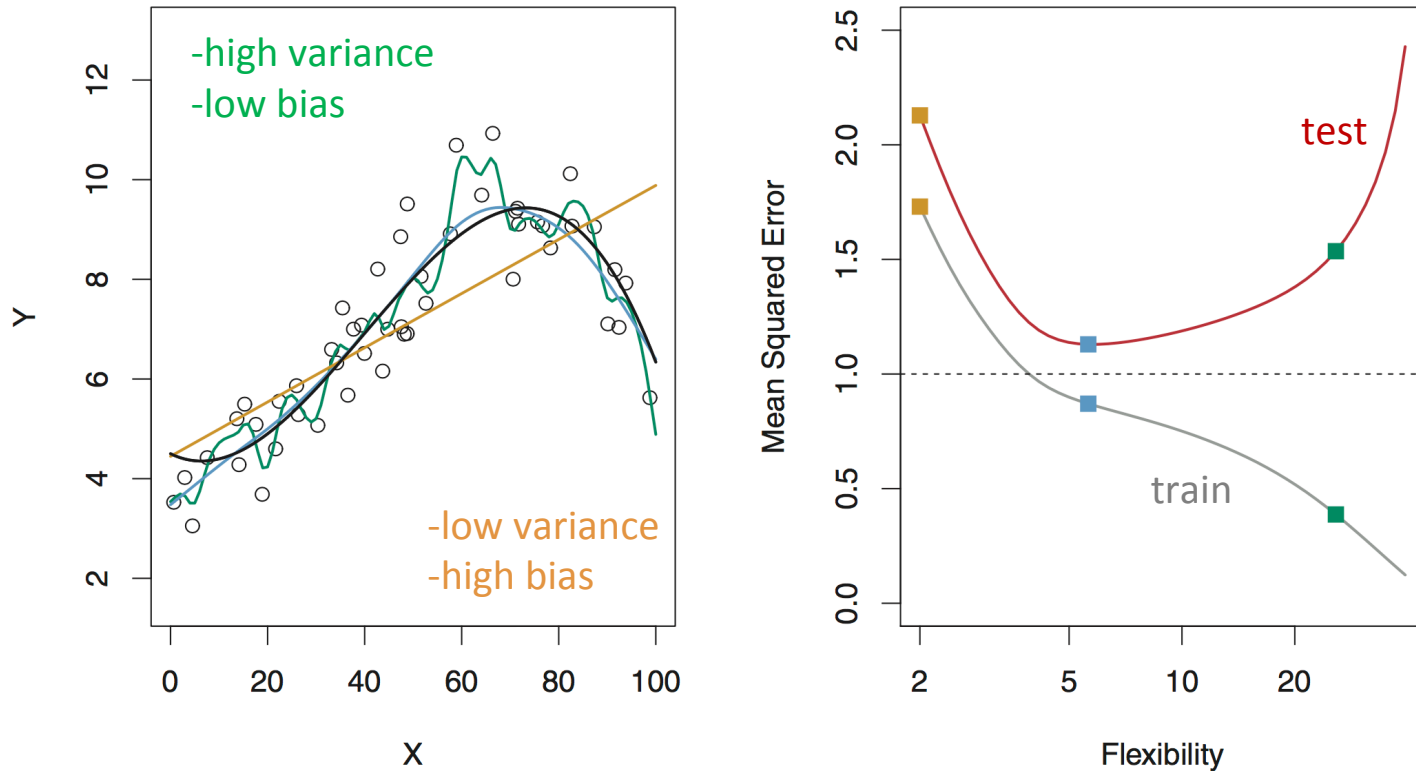
variance

bias

too simple

too complex

# Outline for September 22

- Finish Decision Trees (recap continuous features)

- Learning problem so far + terminology

- Bias-Variance tradeoff

- Linear regression

# Assessing Model Accuracy



**FIGURE 2.9.** Left: *Data simulated from f, shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves).* Right: *Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.*

# Goals of Inference

1) Which of the features/explanatory variables/ predictors (x) are associated with the response variable (y)?

2) What is the relationship between x and y?

3) Is a linear model enough?

4) Can we predict y given a new x?

# Regression Example