

CS 260: Foundations of Data Science

Prof. Thao Nguyen

Fall 2025



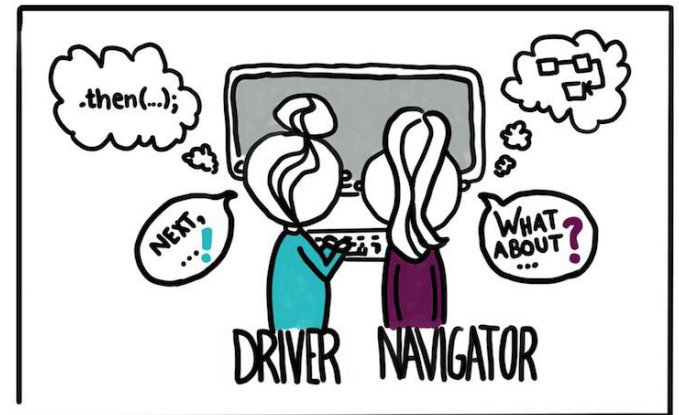
HAVERFORD
COLLEGE

Admin

- **Sit somewhere new!**
- **Lab 1** is due tonight at midnight
- **Lab 2** posted (due next Tuesday)
 - This lab will be done in pairs, please **find a partner**

Pair Programming

- One person is **driver** (at the keyboard)
- One person is **navigator**
- Switch every 30 min!
- Always be working on the assignment together
- No “divide and conquer”
- Make sure to push frequently and pull if there have been any changes



Outline for today

- Data representation and featurization
- Introduction to modeling
- Why are models useful?
- Linear models

Outline for today

- Data representation and featurization
- Introduction to modeling
- Why are models useful?
- Linear models

Tennis Data

Day	Outlook	Temperature	Humidity	Wind	PlayTennis (y)
x_1	Sunny	Hot	High	Weak	No
x_2	Sunny	Hot	High	Strong	No
x_3	Overcast	Hot	High	Weak	Yes
x_4	Rain	Mild	High	Weak	Yes
x_5	Rain	Cool	Normal	Weak	Yes
x_6	Rain	Cool	Normal	Strong	No
x_7	Overcast	Cool	Normal	Strong	Yes
x_8	Sunny	Mild	High	Weak	No
x_9	Sunny	Cool	Normal	Weak	Yes
x_{10}	Rain	Mild	Normal	Weak	Yes

Data from Machine Learning by Tom Mitchell (Table 3.2)

- Input or **features**: outlook, temp, humidity, wind
- Output or “**label**”: play tennis (yes or no)

Sea Ice data (Lab 2)

Year **Sea Ice Extent***

1996	7.88
1997	6.74
1998	6.56
1999	6.24
2000	6.32
2001	6.75
2002	5.96
2003	6.15
2004	6.05
2005	5.57
2006	5.92
2007	4.3
2008	4.63

- Input or **feature**: year
- Output or **“label”**: sea ice extent

*Arctic sea ice extent (1,000,000 km²)

Data Representation Notation

Data Representation

$X =$ matrix

name year Id classes

$\left[\begin{array}{c} \text{---} \end{array} \right]$ \vec{x}_i^T

n examples

$n \times p$

p features

Usually : want to model y as some function of x

label/output

$\vec{y} =$

n

$n \times 1$

i.e. major

Feature Terminology

- *Features*: feature names
 - shape
 - sea ice extent
- *Feature values*: what values are possible
 - {circle, square, triangle}
 - all non-negative values
- *Feature vector*: values for a particular example/data point
 - $\mathbf{x} = [x_1, x_2, x_3, \dots, x_p]$

Featurization: make numerical

- Real-valued features get copied directly. *Duame, Chap 3*
- Binary features become 0 (for false) or 1 (for true).
- Categorical features with V possible values get mapped to V -many binary indicator features.

Q: what about features that might already be on a spectrum
(e.g. sunny, rain, overcast)?

• Regression : $y \in \mathbb{R}$

• Binary classification :

• Multi-class classification

Humidity $\in \{\text{normal, high}\}$
 \Downarrow \Downarrow
 0 1

Outlook $\in \{\text{sunny, overcast, rain}\}$
 \Downarrow \Downarrow \Downarrow
 0 1 2

(continuous)

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

$y \in \{1, 2, \dots, K\}$ (image recognition)

Shape $\in \{\bigcirc, \triangle, \square\}$
 \Downarrow \Downarrow \Downarrow
 0 1 2

	is \bigcirc ?	is \triangle ?	is \square ?
$n=3$ \square	0	0	1
\triangle	0	1	0
\triangle	0	1	0

} binary

Featurization: make numerical

Handout 3

Handout 3

Q1: $n=10, p=4$

Day	Outlook	Temperature	Humidity	Wind	PlayTennis (y)
x_1	Sunny	Hot	High	Weak	No
x_2	Sunny	Hot	High	Strong	No
x_3	Overcast	Hot	High	Weak	Yes
x_4	Rain	Mild	High	Weak	Yes
x_5	Rain	Cool	Normal	Weak	Yes
x_6	Rain	Cool	Normal	Strong	No
x_7	Overcast	Cool	Normal	Strong	Yes
x_8	Sunny	Mild	High	Weak	No
x_9	Sunny	Cool	Normal	Weak	Yes
x_{10}	Rain	Mild	Normal	Weak	Yes

Q2

Sunny: {0,1}
 Overcast: {0,1}
 Rain: {0,1}
 Temperature: {0, 1, 2} (Cool, Mild, Hot)
 Humidity: {0,1} (Normal, High)
 Wind {0,1} (Weak, Strong)

Data from Machine Learning by Tom Mitchell (Table 3.2)

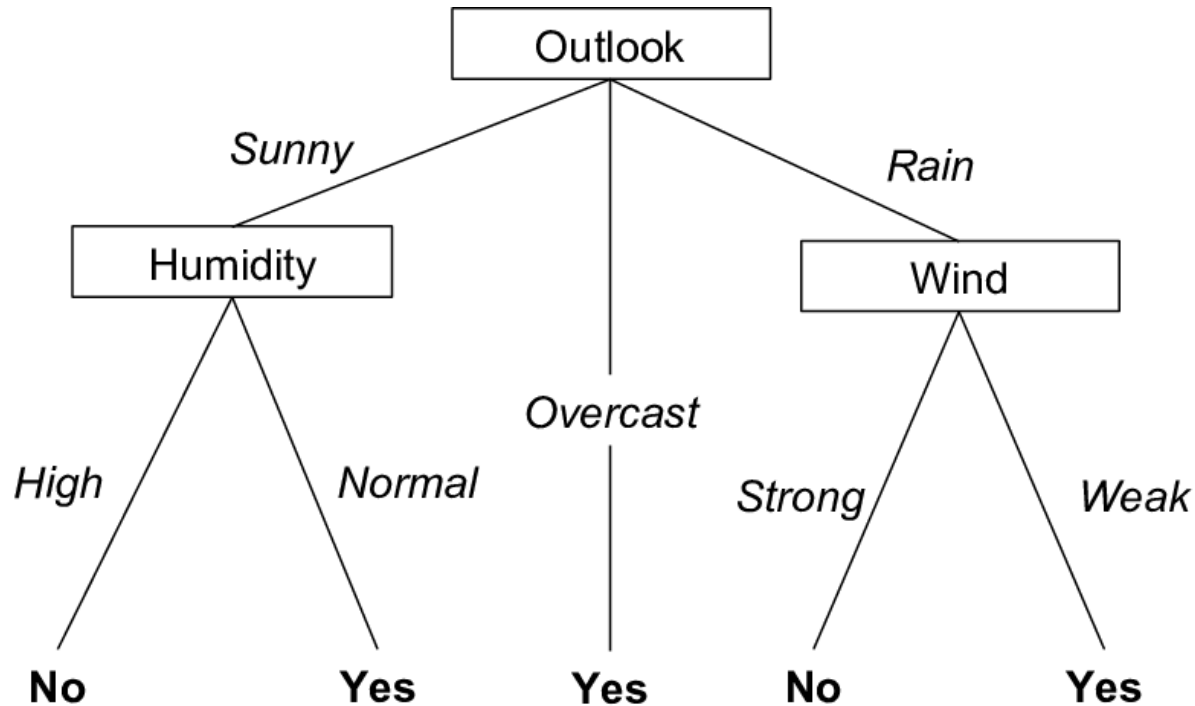
Q3

	Sunny	Overcast	Rain	Temp	Humidity	Wind
x_1	1	0	0	2	1	0

Outline for today

- Data representation and featurization
- Introduction to modeling
- Why are models useful?
- Linear models

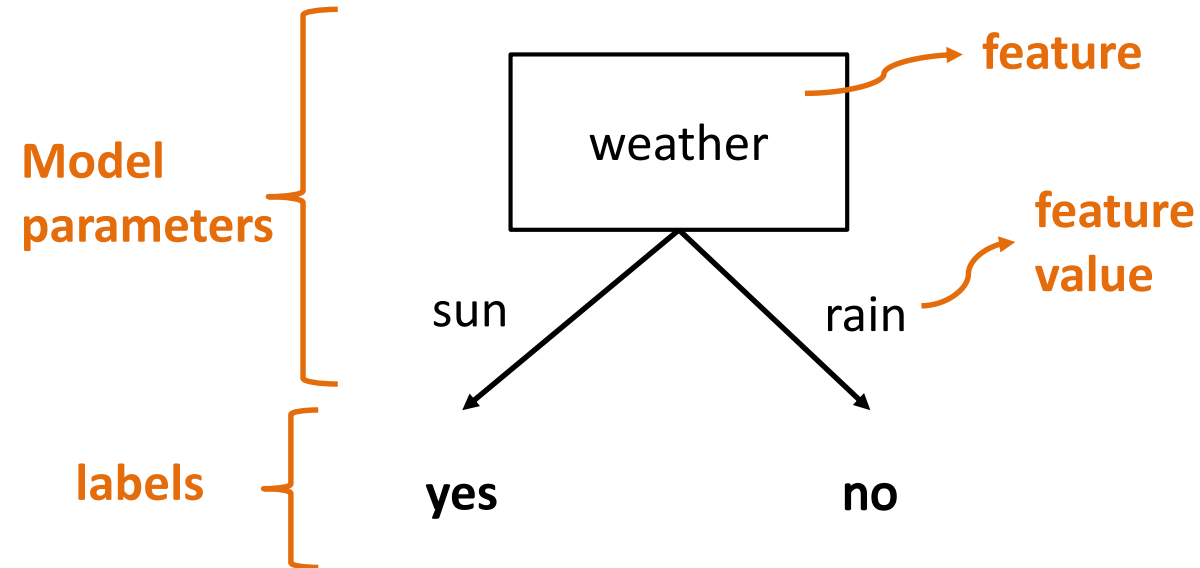
Example of a model



- Each internal node: one feature
- Each branch from node: selects one value of the feature
- Each leaf node: predict y

Model Examples

1) Decision Tree



Data

weather	tennis
sun	yes
rain	no
rain	no
sun	yes
sun	no

[1, 2]

[2, 0]

→ “no” and “yes” label counts

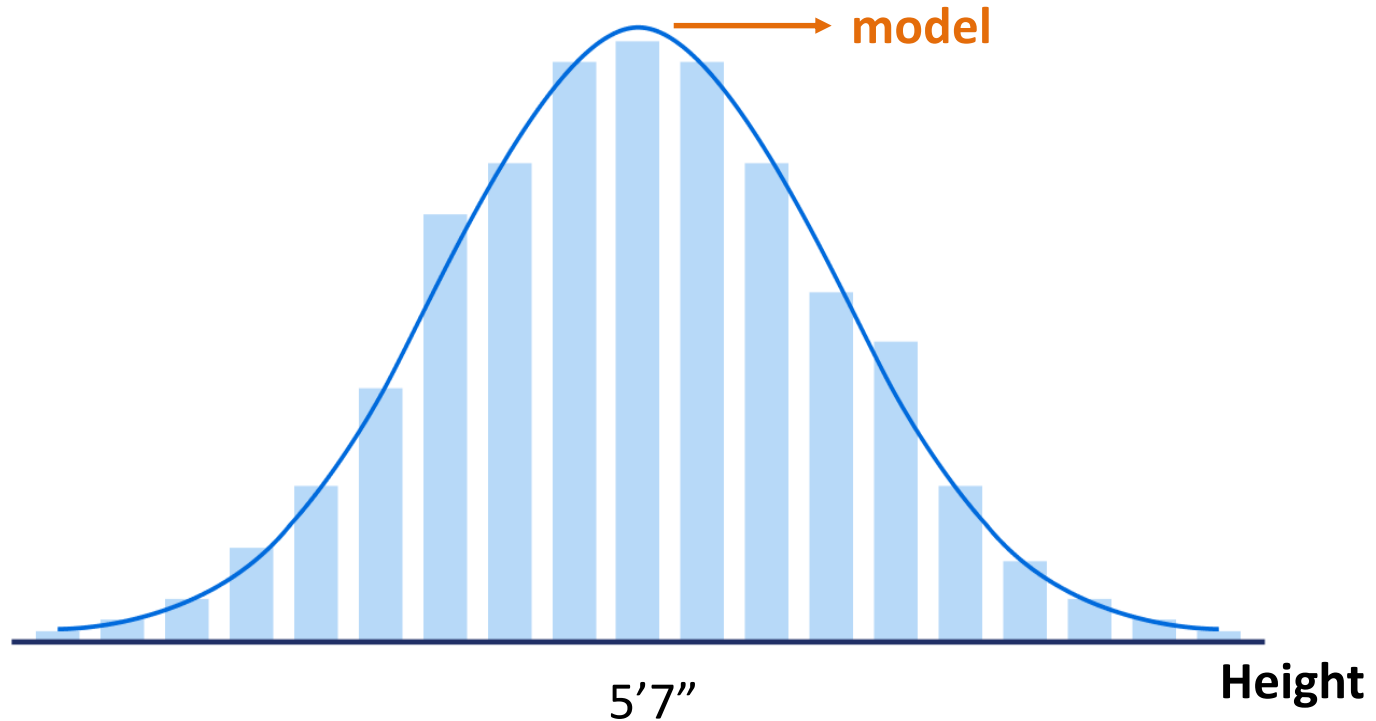
2/3 correct

2/2 correct

=> 80% accuracy

Model Examples

2) Normal distribution

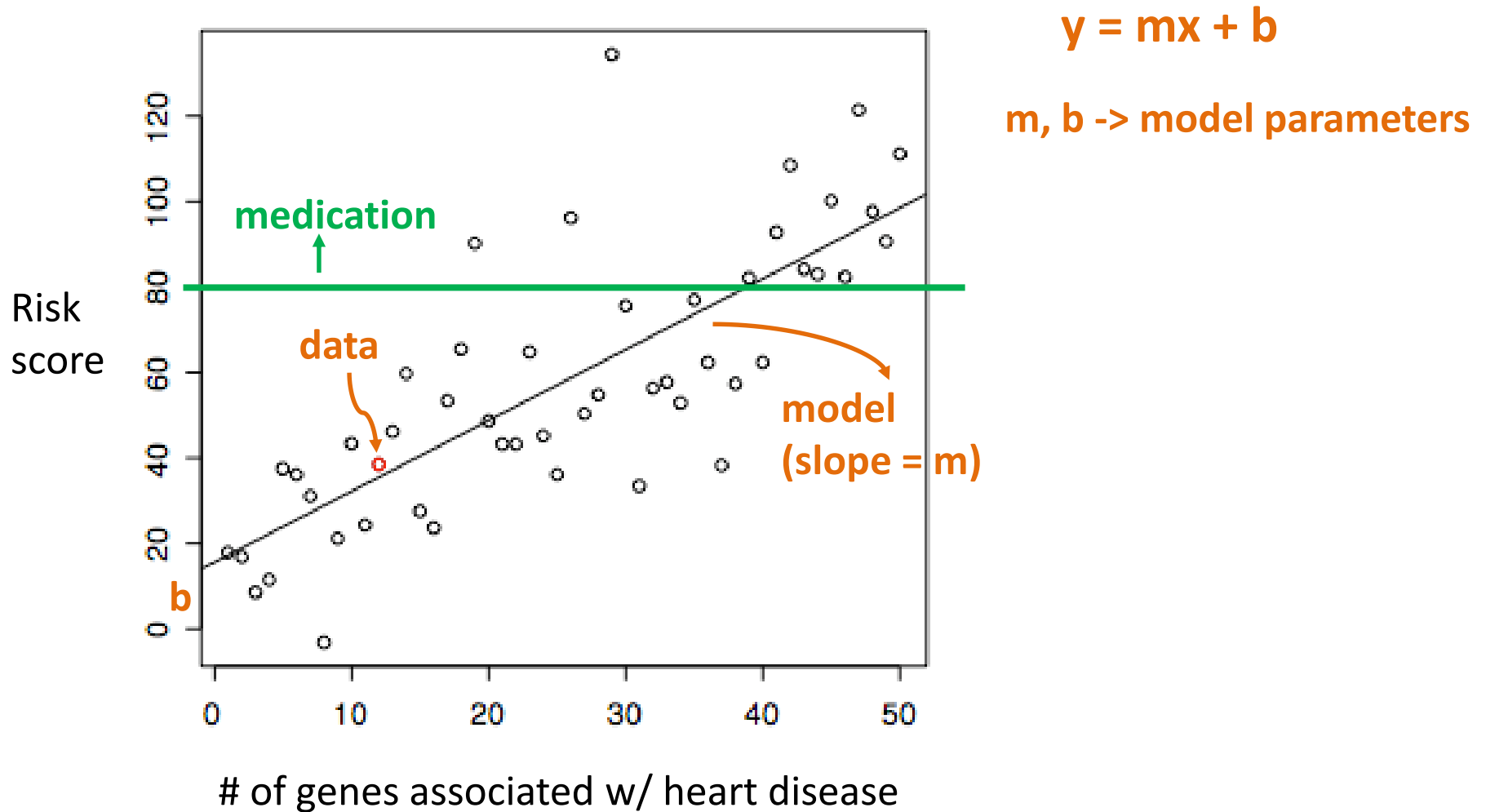


mean: 5'7"
variance: 2"

} **Model
parameters**

Model Examples

3) Linear models

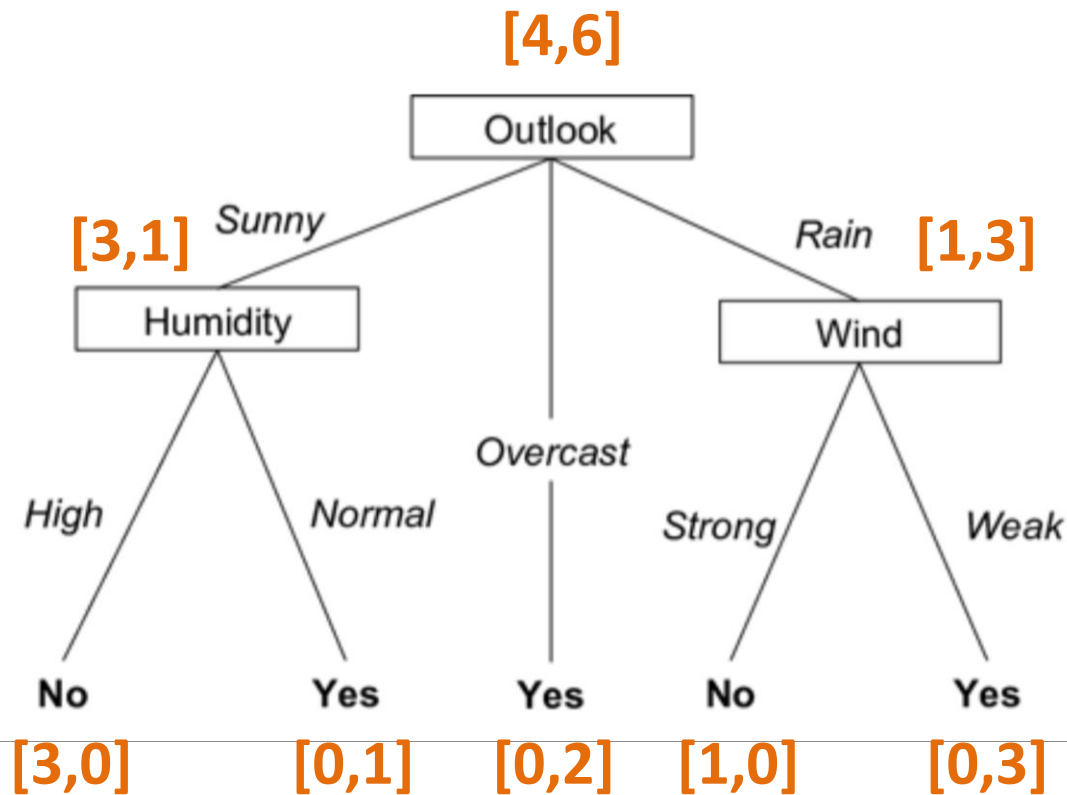


Handout 3

Handout 3

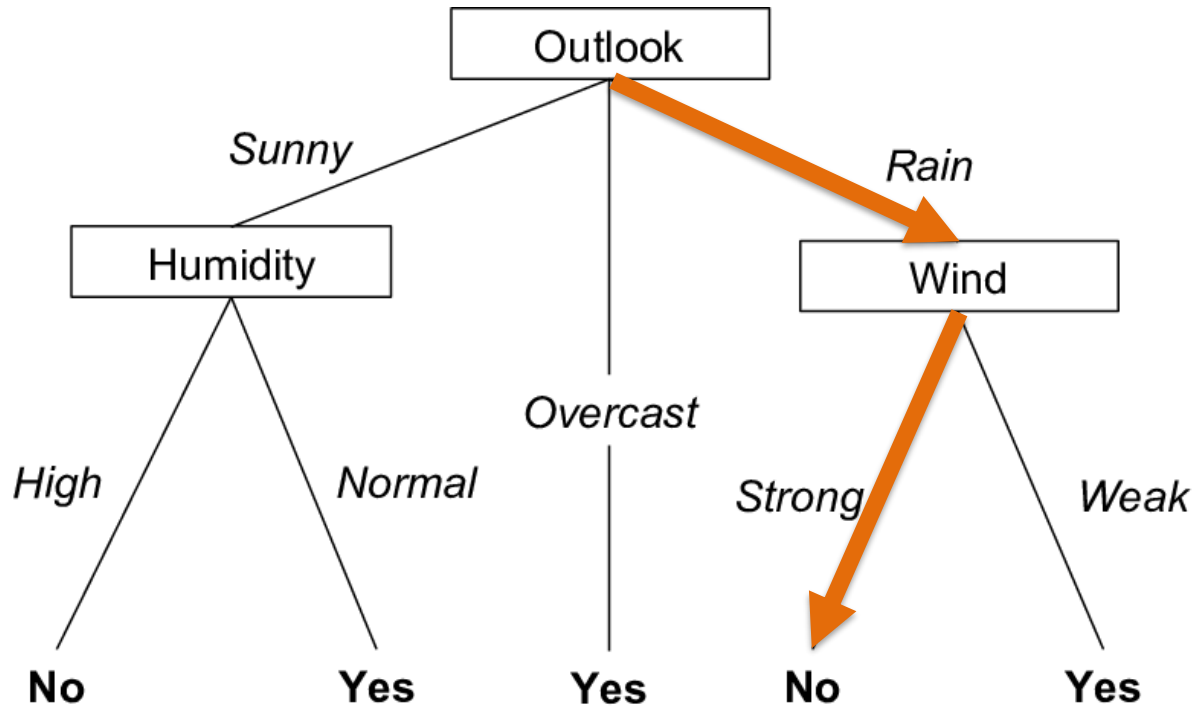
Q4

In the model below, the children of each node divide the data into partitions. Label each node (both internal nodes and leaves) with the counts of “No” and “Yes” labels based on the partition. For example, the counts for the node labeled *Outlook* would be [4,6]. Does this model perfectly classify all examples?



Handout 3

Q5



(test example) $x =$

Outlook	Temp	Humidity	Wind
Rain	Hot	High	Strong

$y_{pred} = \text{No}$

Outline for today

- Data representation and featurization
- Introduction to modeling
- **Why are models useful?**
- Linear models

Why are models useful?

- Understand/explain/interpret the phenomena
- Predict outcomes for future examples

What are the most important features?

X

Color	Shape	Size
red	square	big
blue	square	big
red	circle	small
blue	square	small
red	circle	big

Y

Likes toy?
+
+
-
-
+

What are the most important features?

X

Y

Color	Shape	Size
red	square	big
blue	square	big
red	circle	big
blue	square	big
red	circle	big

Likes toy?
+
+
-
-
+

Outline for today

- Data representation and featurization
- Introduction to modeling
- Why are models useful?
- Linear models

Linear Models

* features: \vec{x} ($p=1$
call it x)

* label: $y \in \mathbb{R}$
output

Goals

- ① describe linear dependence.
- ② predict output given new data

model

$$h_{\vec{w}}(x) = w_0 + w_1 x = \hat{y}$$

"b" "m"

$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$

prediction

how good is our model?

residuals

$$y_i - \hat{y}_i$$

truth prediction

one example

Overall

want to minimize

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$(w_0 + w_1 x_i)$

RSS or SSE



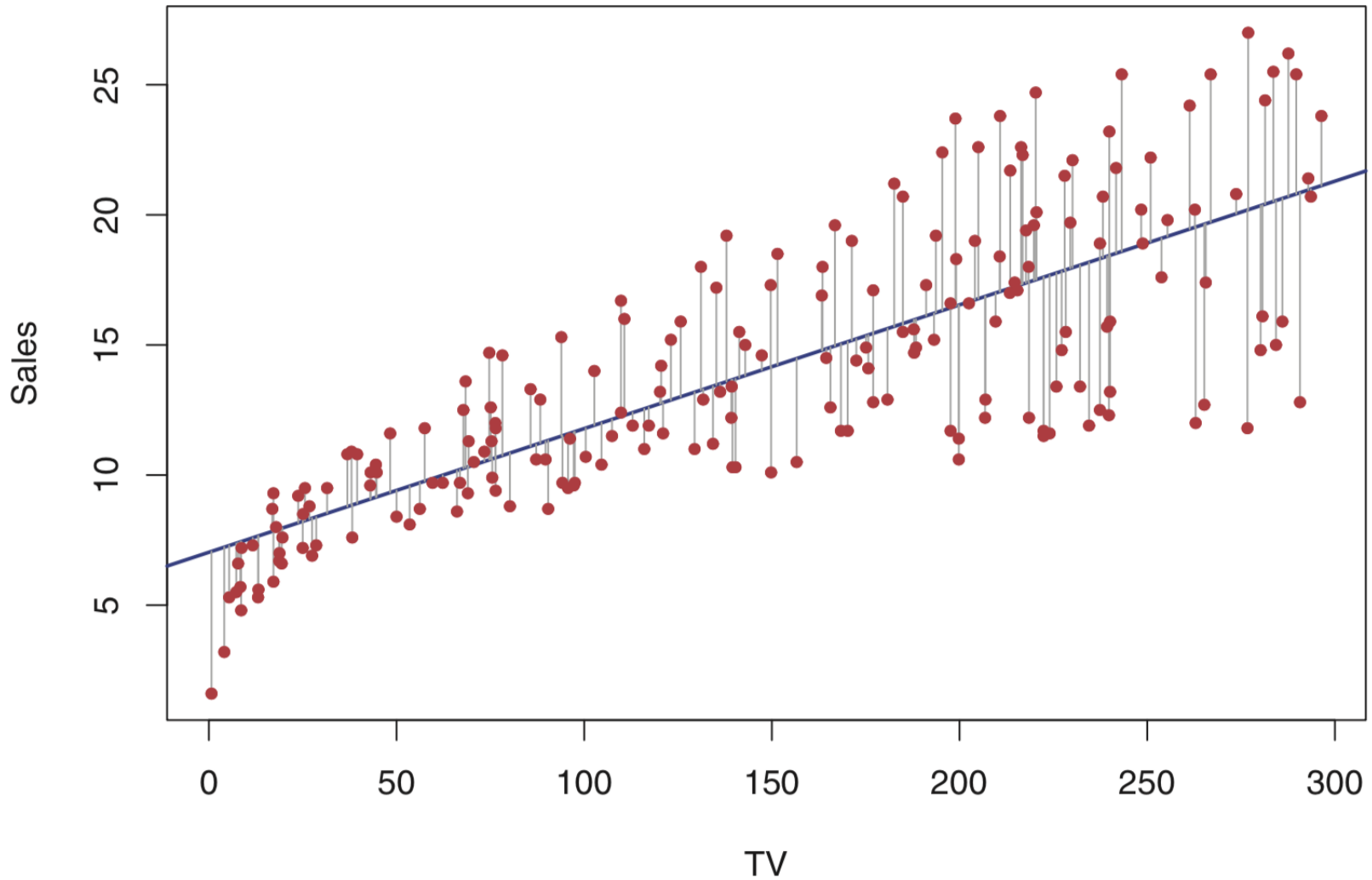
RSS : residual sum of squares

SSE : sum of squared errors

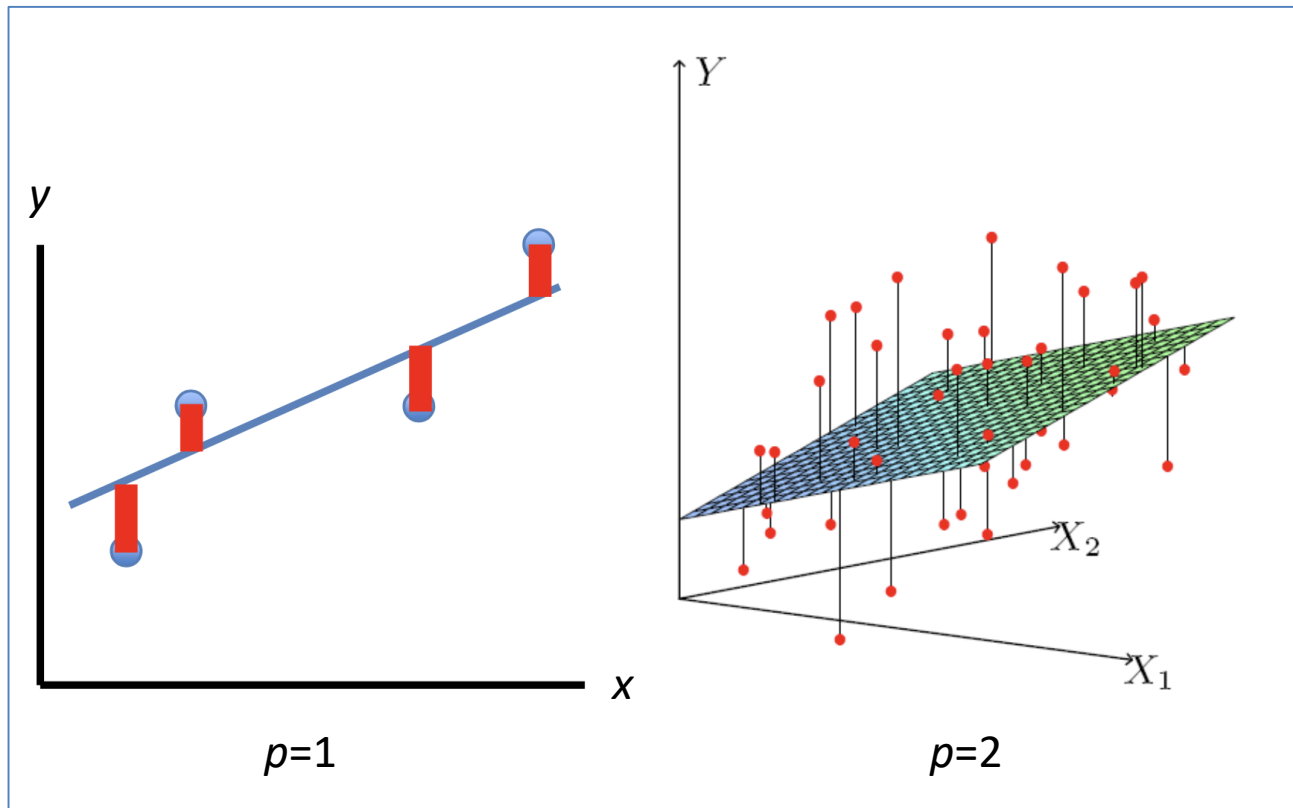
Goals of fitting a linear model

- 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) Can we (accurately) predict y given a new x ?
- 4) Is a linear model enough?

Example: predict sales from TV advertising budget



Linear model with 1 or 2 features



Linear Regression

- Output (y) is continuous, not a discrete label
- Learned model: *linear function* mapping input to output (a *weight* for each feature + *bias*)
- Goal: minimize the *RSS* (residual sum of squares) or *SSE* (sum of squared errors)

Maybe a linear model is not enough

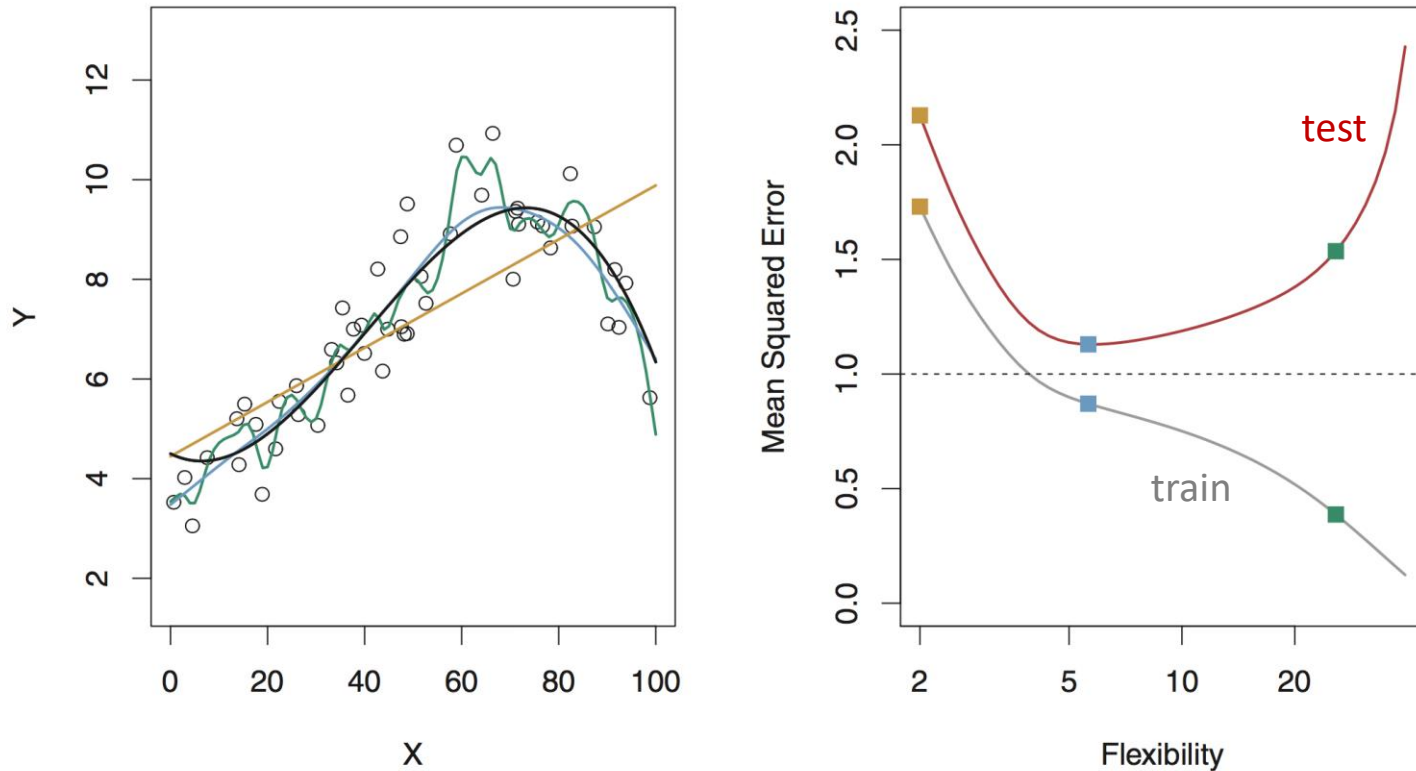


FIGURE 2.9. Left: Data simulated from f , shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line).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