

# Ch 2.1: What is Statistical Learning?

Lecture 2 - CMSE 381

Prof. Elizabeth Munch

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

Weds, Aug 28, 2024

# Announcements

## Last time:

- Discussed where to find everything
  - ▶ Course webpage
  - ▶ Slack
  - ▶ D2L
- Check out the syllabus!

Lec #	Date			Reading	HW
1	Mon	8/26	Intro / First day stuff / Python Review Pt 1	1	
2	Wed	8/28	What is statistical learning?	2.1	
	Fri	8/30	Class Cancelled (Dr Munch out of town)		
	Mon	9/2	No class - Labor day		
3	Wed	9/4	Assessing Model Accuracy	2.2.1, 2.2.2	
4	Fri	9/6	Linear Regression	3.1	HW #1 Due

## Announcements:

- Get on slack!
  - ▶ +1 point on the first homework if you post a gif in the thread
- First homework due Sun Sep 8
- First office hours next week

# Covered in this class

- Input/output variables
  - Prediction vs inference
  - Reduceable vs irreducible error
  - Overfitting
  - Classification vs regression
  - Supervised vs Unsupervised learning
- Please note: no jupyter notebook for today's class, slides only

# An example data set: Advertising

1		TV	Radio	Newspaper	Sales
2	1	230.1	37.8	69.2	22.1
3	2	44.5	39.3	45.1	10.4
4	3	17.2	45.9	69.3	9.3
5	4	151.5	41.3	58.5	18.5
6	5	180.8	10.8	58.4	12.9
7	6	8.7	48.9	75	7.2
8	7	57.5	32.8	23.5	11.8
9	8	120.2	19.6	11.6	13.2
10	9	8.6	2.1	1	4.8
11	10	199.8	2.6	21.2	10.6
12	11	66.1	5.8	24.2	8.6

- Sales of a product in 200 markets, along with amount spent on three different types of advertising
- Goal:
- Input variables:
- Output variable:

Data available at [msu-cmse-courses.github.io/CMSE381-F24/DataSets/DataSets.html](https://msu-cmse-courses.github.io/CMSE381-F24/DataSets/DataSets.html)

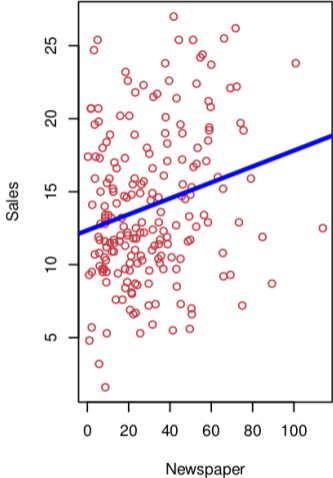
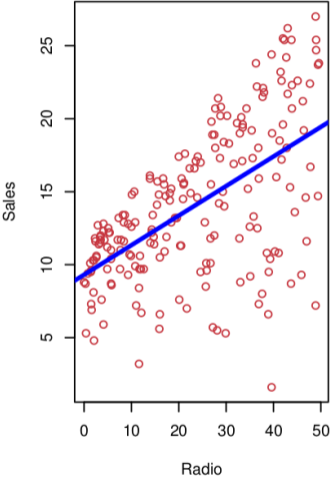
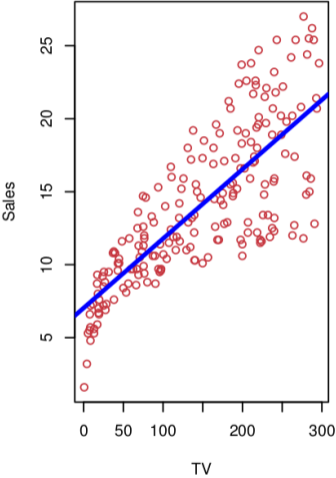
# Notation and Big Assumption

Input variables:  $X_1, X_2, \dots, X_p$

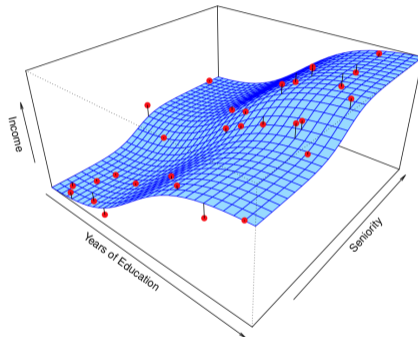
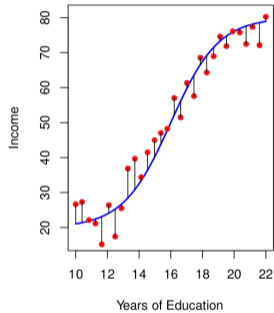
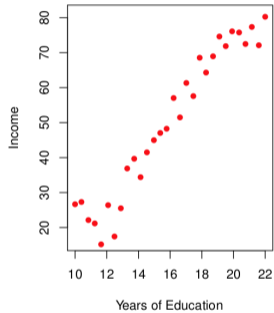
Output variable:  $Y$

$$Y = f(X) + \varepsilon$$

# Advertising Example



# More examples



# Section 1

## Prediction vs Inference



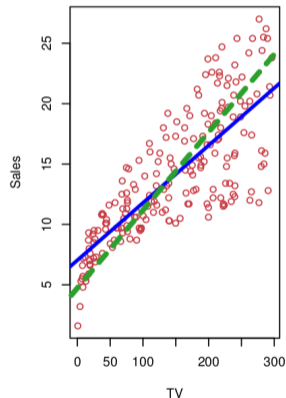
# Prediction

Given a value  $X$ , try to provide an estimate for  $f(X)$ .

Build a model:

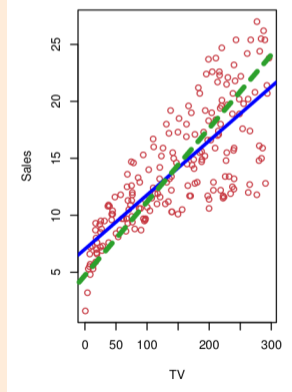
$$\hat{Y} = \hat{f}(X)$$

Example: If we spend \$250 on TV advertising, what do we predict we will we make in sales?



## Group question:

1		TV	Radio	Newspaper	Sales
2	1	230.1	37.8	69.2	22.1
3	2	44.5	39.3	45.1	10.4
4	3	17.2	45.9	69.3	9.3



The blue solid line is  $f$ . The green dashed line is  $\hat{f}$ .

- What is the predicted sales for the first three data points using the green dashed line  $\hat{f}$  shown in the graph?
- Using the dashed green line as the predicted model  $\hat{f}$ , what is the error in each of the three predictions?

# Reduceable vs irreducible error

All models are wrong, some are useful.

**Reduceable Error**

$$Y - \hat{Y}$$

**Irreducible Error**

## More on error

- Given estimate  $\hat{f}$  (fixed)
- Set of predictors  $X$  (fixed)
- Prediction  $\hat{Y} = \hat{f}(X)$

$$E(Y - \hat{Y})^2 =$$

Want  $f$ , but not for prediction  
(or possibly combined with  
prediction)

- Which predictors are associated with the response?
- What is the relationship between the response and each predictor?
- Can the relationship between  $Y$  and each predictor be adequately summarized using a linear equation? Is it more complicated?

Determine whether each scenario is prediction, inference, or both.

<b>Application</b>	<b>Prediction</b>	<b>Inference</b>
Predict effectiveness of vaccine		
Determine the address written on the image of an envelope.		
Identify risk factors for getting long covid.		
Transcribe an audio file of a person talking.		
Predict stock prices.		

## Section 2

How to estimate  $f$ ?

# Input: Training data

- $n$  data points observed
- $x_{ij}$  is the  $j$ th predictor for observation  $i$
- $y_i$  is the response variable for the  $i$ th observation
- Training data:
  - ▶  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
  - ▶  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$

		TV	Radio	Newspaper	Sales
1					
2	1	230.1	37.8	69.2	22.1
3	2	44.5	39.3	45.1	10.4
4	3	17.2	45.9	69.3	9.3
5	4	151.5	41.3	58.5	18.5
6	5	180.8	10.8	58.4	12.9
7	6	8.7	48.9	75	7.2
8	7	57.5	32.8	23.5	11.8
9	8	120.2	19.6	11.6	13.2
10	9	8.6	2.1	1	4.8
11	10	199.8	2.6	21.2	10.6
12	11	66.1	5.8	24.2	8.6



# Parametric methods

**Step 1:** Select a model

Example:

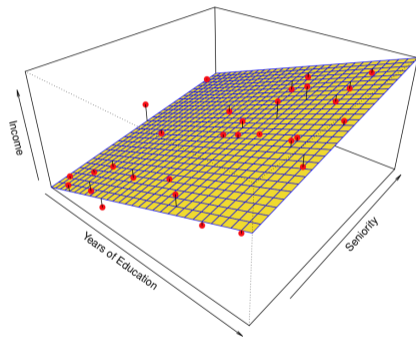
$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

**Step 2:** Train the model

Example:

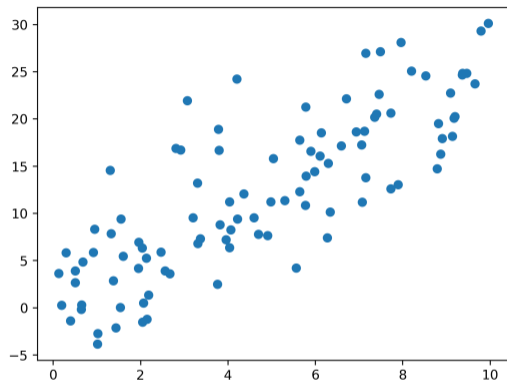
Find  $\beta_i$ 's so that

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$



# How do you decide on the coefficients?

$$Y \approx \beta_0 + \beta_1 X_1$$

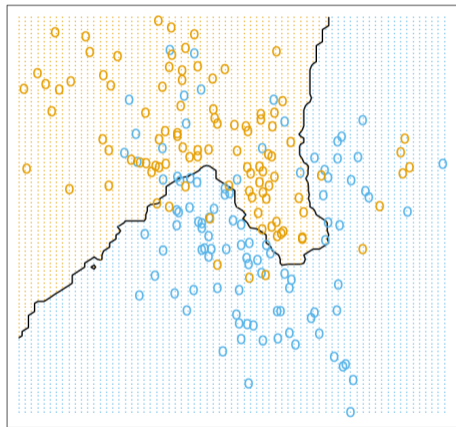


Desmos toy: <https://www.desmos.com/calculator/skvt8c7317>

## Example Non-parametric method: Nearest Neighbors

$N_k(x)$  = Set of  $k$  nearest neighbors of  $x$

$$\hat{f}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$$



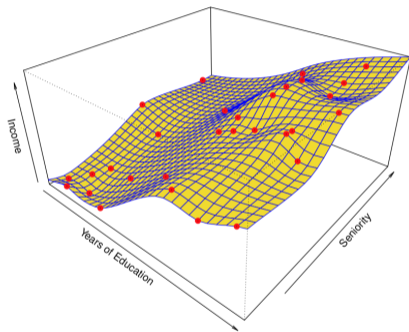
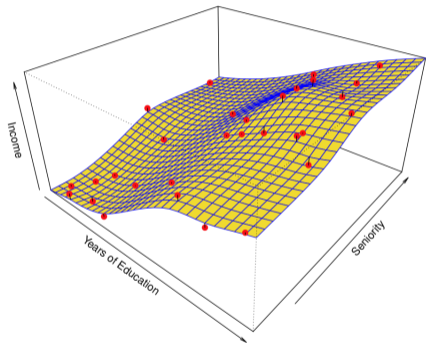
$k = 15$

# Parametric methods: Pros and Cons

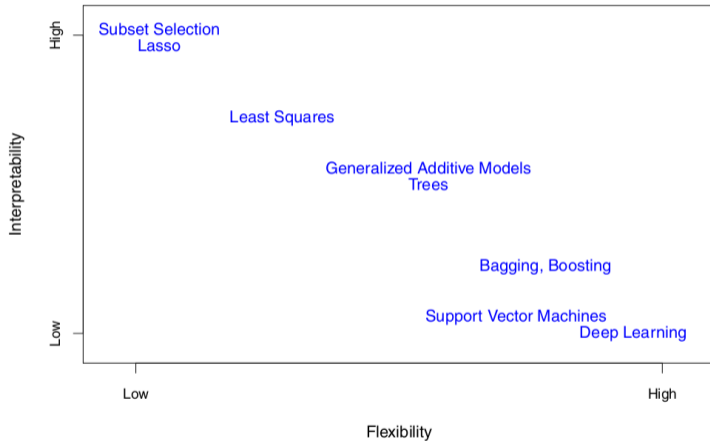
**Pros**

**Cons**

# Overfitting

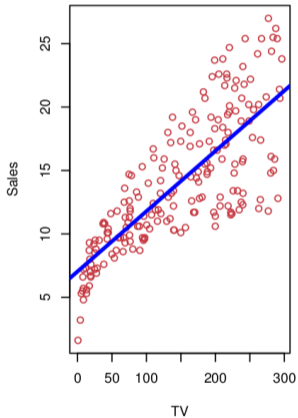


# Prediction Accuracy vs Model Interpretability



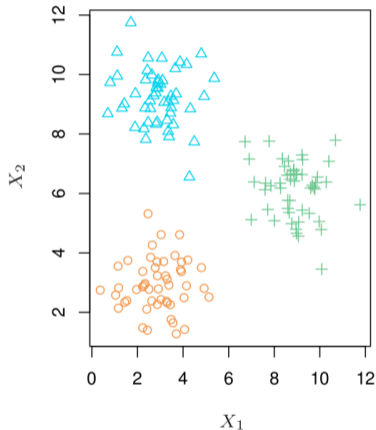
## Supervised learning:

Training data has response variable  $y$  for every input  $x$



## Unsupervised Learning:

Training data does not have response variable  $y$  for every input  $x$



# Regression vs Classification

## Types of variables:

- Quantitative
  
  
  
  
  
  
  
  
  
  
- Qualitative / Categorical



## Section 3

Group work

(a) We collect a set of data on the top 500 firms in the US. For each firm we record profit, number of employees, industry and the CEO salary. We are interested in understanding which factors affect CEO salary.

- Is this classification or regression?
- Do we want inference or prediction?
- What is  $n$ , the number of data points?
- What is  $p$ , the number of variables?

(b) We are considering launching a new product and wish to know whether it will be a success or a failure. We collect data on 20 similar products that were previously launched. For each product we have recorded whether it was a success or failure, price charged for the product, marketing budget, competition price, and ten other variables.

- Is this classification or regression?
- Do we want inference or prediction?
- What is  $n$ , the number of data points?
- What is  $p$ , the number of variables?



## Next time:

- Friday 8/30 and Monday 9/2
  - ▶ NO CLASS!!!!
- Wednesday 9/4
  - ▶ Bring Laptop!
  - ▶ First homework due Sun Sep 8
  - ▶ There will be a quiz next week

Lec #	Date		Reading	HW
1	Mon 8/26	Intro / First day stuff / Python Review Pt 1	1	
2	Wed 8/28	What is statistical learning?	2.1	
	Fri 8/30	Class Cancelled (Dr Munch out of town)		
	Mon 9/2	No class - Labor day		
3	Wed 9/4	Assessing Model Accuracy	2.2.1, 2.2.2	
4	Fri 9/6	Linear Regression	3.1	HW #1 Due Sun 9/8
5	Mon 9/9	More Linear Regression	3.1/3.2	

## Announcements:

- Get on slack!
  - ▶ +1 point on the first homework if you post a gif in the thread
- Office hours!
  - ▶ Dr. Munch: TBD
  - ▶ Christy Lu: TBD