

# Machine Learning

---

## Outline

- I. ML and applications
- II. Supervised learning
- III. Decision trees

Com S 474/574 *Introduction to Machine Learning* (Spring 2022)

\* A large portion of the material is drawn from Dr. Jin Tian's notes.

\*\* Figures are from either the [textbook site](#) or Dr. Jin Tian's notes.

# AI and Machine Learning

---

- ◆ AI is the enterprise of design and analysis of intelligent agents.
- ◆ Intelligent behavior requires **knowledge** (e.g., model of the environment).
- ◆ Explicit specifications of the knowledge needed for specific tasks are hard, and often infeasible.
- ◆ How to acquire knowledge?

# AI and Machine Learning

---

- ◆ AI is the enterprise of design and analysis of intelligent agents.
- ◆ Intelligent behavior requires **knowledge** (e.g., model of the environment).
- ◆ Explicit specifications of the knowledge needed for specific tasks are hard, and often infeasible.
- ◆ How to acquire knowledge?

**Machine learning (ML)** refers to the process in which a computer

- observes some data,
- builds a **model** based on the data, and
- uses the model as both a hypothesis about the world and a piece of problem solving software.

# Learning Agents

---

- ♣ Learning modifies the agent's decision mechanisms to **improve** performance.
  - Which **component** is to be improved.
  - Which **prior knowledge** the agent has, which influences the model.
  - What **data** and **feedback** on that data is available.
- ♣ Environment changes over time – learning needs to adapt to changes.
- ♣ Learning is essential for unknown environments.

# Applications of ML

---

- Agriculture
- Anatomy
- Adaptive websites
- Affective computing
- Banking
- Bioinformatics
- Brain–machine interfaces
- Cheminformatics
- Citizen science
- Computer networks
- Computer vision
- Credit-card fraud detection
- Data quality
- DNA sequence classification
- Economics
- Financial market analysis<sup>[75]</sup>
- General game playing
- Handwriting recognition
- Information retrieval
- Insurance
- Internet fraud detection
- Linguistics
- Machine learning control
- Machine perception
- Machine translation
- Marketing
- Medical diagnosis
- Natural language processing
- Natural language understanding
- Online advertising
- Optimization
- Recommender systems
- Robot locomotion
- Search engines
- Sentiment analysis
- Sequence mining
- Software engineering
- Speech recognition
- Structural health monitoring
- Syntactic pattern recognition
- Telecommunication
- Theorem proving
- Time series forecasting
- User behavior analytics

# Applications of ML

---

- Agriculture
- Anatomy
- Adaptive websites
- Affective computing
- Banking
- Bioinformatics
- Brain–machine interfaces
- Cheminformatics
- Citizen science
- Natural language processing
- Natural language understanding
- Online advertising
- Optimization
- Recommender systems
- Robot locomotion
- Search engines
- Sentiment analysis
- Computer networks
- Computer vision
- Credit-card fraud detection
- Data quality
- DNA sequence classification
- Economics
- Financial market analysis<sup>[75]</sup>
- General game playing
- Handwriting recognition
- Sequence mining
- Software engineering
- Speech recognition
- Structural health monitoring
- Syntactic pattern recognition
- Telecommunication
- Theorem proving
- Time series forecasting
- User behavior analytics
- Information retrieval
- Insurance
- Internet fraud detection
- Linguistics
- Machine learning control
- Machine perception
- Machine translation
- Marketing
- Medical diagnosis

# Data Mining

---

- ◆ Huge amounts of data are available from science, medicine, economics, geography, environment, sports, ...
- ◆ Data is a potentially valuable resource.
- ◆ Raw data are useless – need techniques to automatically extract information from it.
  - Data: recorded facts
  - Information: patterns underlying the data
- ◆ Machine learning techniques automatically find patterns in data.

# The Game-Weather Problem

---

Weather condition for playing a certain game:

Outlook	Temperature	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	Normal	False	Yes
...	...	...	...	...

# The Game-Weather Problem

---

Weather condition for playing a certain game:

Outlook	Temperature	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	Normal	False	Yes
...	...	...	...	...

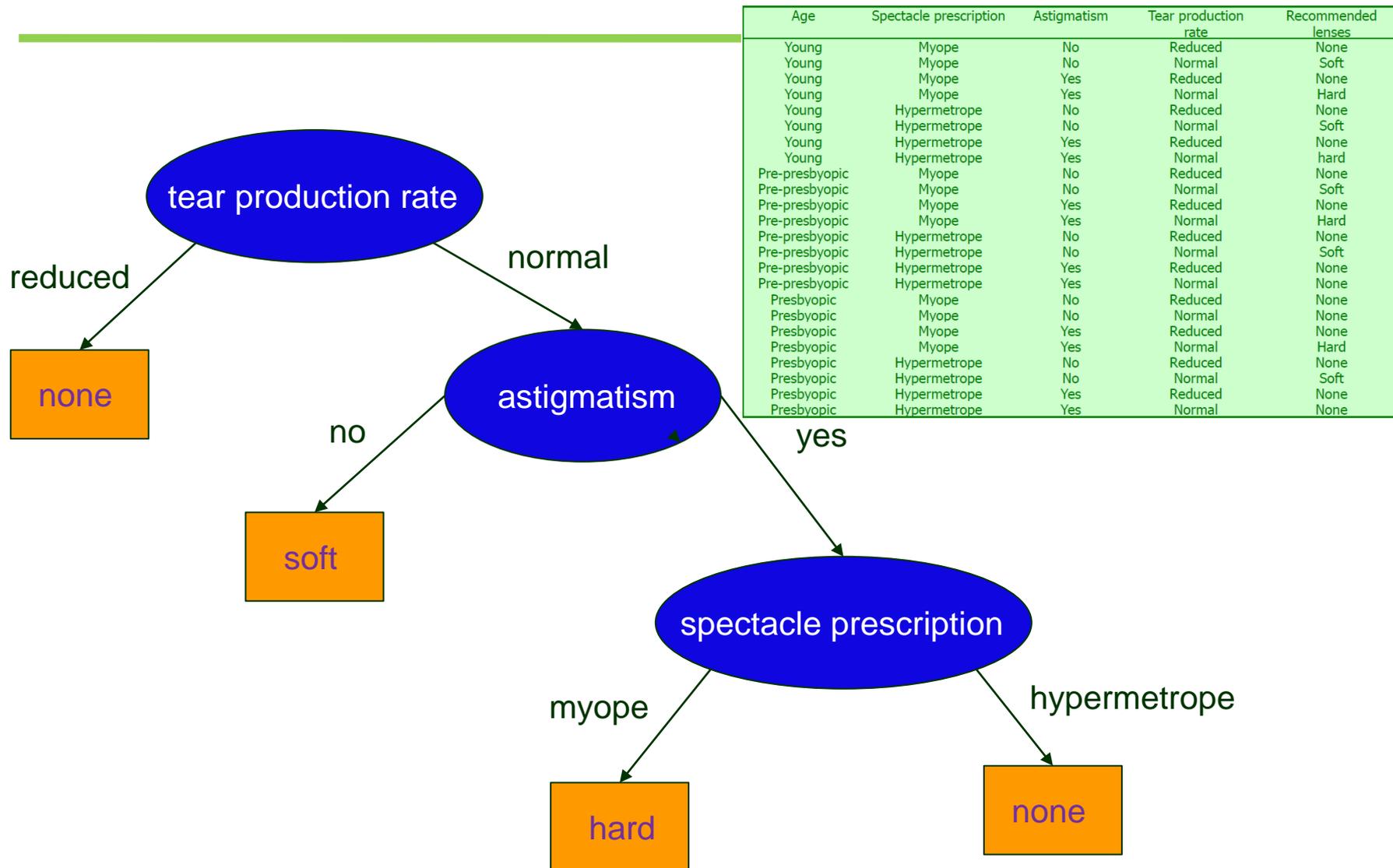
Learned classification rules:

```
If outlook = sunny and humidity = high then play = no
If outlook = rainy and windy = true then play = no
If outlook = overcast then play = yes
If humidity = normal then play = yes
If none of the above then play = yes
```

# Contact Lense Data

Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	No	Reduced	None
Young	Myope	No	Normal	Soft
Young	Myope	Yes	Reduced	None
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	No	Reduced	None
Young	Hypermetrope	No	Normal	Soft
Young	Hypermetrope	Yes	Reduced	None
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	No	Reduced	None
Pre-presbyopic	Myope	No	Normal	Soft
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	No	Reduced	None
Pre-presbyopic	Hypermetrope	No	Normal	Soft
Pre-presbyopic	Hypermetrope	Yes	Reduced	None
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	No	Reduced	None
Presbyopic	Myope	No	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	No	Reduced	None
Presbyopic	Hypermetrope	No	Normal	Soft
Presbyopic	Hypermetrope	Yes	Reduced	None
Presbyopic	Hypermetrope	Yes	Normal	None

# Learned Decision Tree



# Predicting CPU Performance

---

209 different computer configurations

	Cycle time (ns)	Main memory (Kb)		Cache (Kb)	Channels		Performance
	MYCT	MMIN	MMAX	CACH	CHMIN	CHMAX	PRP
1	125	256	6000	256	16	128	198
2	29	8000	32000	32	8	32	269
...							
208	480	512	8000	32	0	0	67
209	480	1000	4000	0	0	0	45

# Predicting CPU Performance

209 different computer configurations

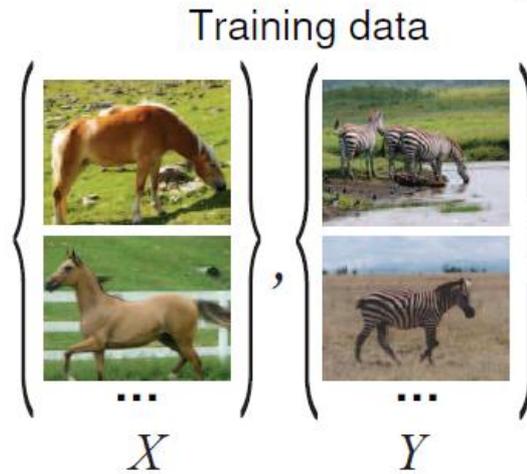
	Cycle time (ns)	Main memory (Kb)		Cache (Kb)	Channels		Performance
	MYCT	MMIN	MMAX	CACH	CHMIN	CHMAX	PRP
1	125	256	6000	256	16	128	198
2	29	8000	32000	32	8	32	269
...							
208	480	512	8000	32	0	0	67
209	480	1000	4000	0	0	0	45

Function obtained through linear regression (fitting):

$$\text{PRP} = -55.9 + 0.0489 \text{ MYCT} + 0.0153 \text{ MMIN} + 0.0056 \text{ MMAX} \\ + 0.6410 \text{ CACH} - 0.2700 \text{ CHMIN} + 1.480 \text{ CHMAX}$$

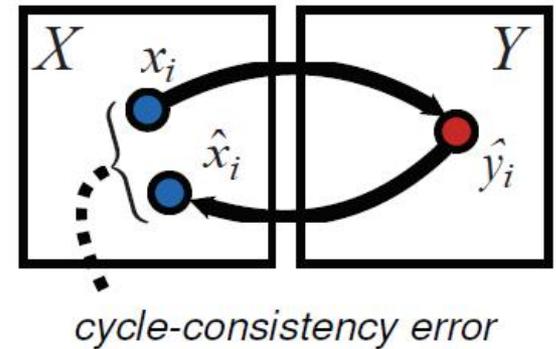
# Image Translation

Translate a horse into a zebra  
(find corresponding pairs):



Training

Objective



Test

Input



Result



# Machine Learning Models

---

## ♣ Supervised learning

- ◆ Naïve Bayes classifier
- ◆ Nearest neighbor methods
- ◆ Linear models
- ◆ Decision trees
- ◆ Neural networks
- ◆ Support vector machines
- ◆ Ensemble learning

## ♣ Probabilistic graphical models

- ◆ Bayesian networks
- ◆ Markov random fields

## ♣ Unsupervised learning

- ◆ Clustering: mixture models, K-means, hierarchical clustering
- ◆ Principal component analysis
- ◆ Independent component analysis

## ♣ Sequential data

- ◆ HMMs
- ◆ Recurrent neural networks

## ♣ Markov decision process

## ♣ Reinforcement Learning

# Machine Learning Models

---

## ♣ Supervised learning

- ◆ Naïve Bayes classifier
- ◆ Nearest neighbor methods
- ◆ Linear models
- ◆ Decision trees
- ◆ Neural networks
- ◆ Support vector machines
- ◆ Ensemble learning

## ♣ Probabilistic graphical models

- ◆ Bayesian networks
- ◆ Markov random fields

## ♣ Unsupervised learning

- ◆ Clustering: mixture models, K-means, hierarchical clustering
- ◆ Principal component analysis
- ◆ Independent component analysis

## ♣ Sequential data

- ◆ HMMs
- ◆ Recurrent neural networks

## ♣ Markov decision process

## ♣ Reinforcement Learning

# Supervised Learning

---

The agent observes input-output pairs and learns a function that maps from input to output.

# Supervised Learning

---

The agent observes input-output pairs and learns a function that maps from input to output.

**Problem** Given a *training set* of  $N$  input-output pairs:

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$$

where each pair was generated by an unknown function  $y = f(x)$ , discover a function  $h$  to approximate  $f$ .

# Supervised Learning

---

The agent observes input-output pairs and learns a function that maps from input to output.

**Problem** Given a *training set* of  $N$  input-output pairs:

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$$

where each pair was generated by an unknown function  $y = f(x)$ , discover a function  $h$  to approximate  $f$ .

|  
*hypothesis* (or *model*) drawn from a  
*hypothesis space*  $\mathcal{H}$  of possible functions

# Supervised Learning

---

The agent observes input-output pairs and learns a function that maps from input to output.

**Problem** Given a *training set* of  $N$  input-output pairs:

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$$

where each pair was generated by an unknown function  $y = f(x)$ , discover a function  $h$  to approximate  $f$ .

|  
*hypothesis* (or *model*) drawn from a  
*hypothesis space*  $\mathcal{H}$  of possible functions

chosen according to some prior  
knowledge about data generation

# Supervised Learning

---

The agent observes input-output pairs and learns a function that maps from input to output.

**Problem** Given a *training set* of  $N$  input-output pairs:

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$$

where each pair was generated by an unknown function  $y = f(x)$ , discover a function  $h$  to approximate  $f$ .

|  
*hypothesis* (or *model*) drawn from a  
*hypothesis space*  $\mathcal{H}$  of possible functions

chosen according to some prior  
knowledge about data generation

$y_1, y_2, \dots, y_n$ : *ground truth* to be predicted by our model

# Best-Fit Function

---

- ♣ With noise in data, we cannot expect an exact match with the ground truth, namely,  $h(x_i) = y_i$ , for  $1 \leq i \leq N$ .

# Best-Fit Function

---

- ♣ With noise in data, we cannot expect an exact match with the ground truth, namely,  $h(x_i) = y_i$ , for  $1 \leq i \leq N$ .
- ♣ Instead, we look for a best-fit function  $h$  for which each  $h(x_i)$  is close  $y_i$ .

# Best-Fit Function

---

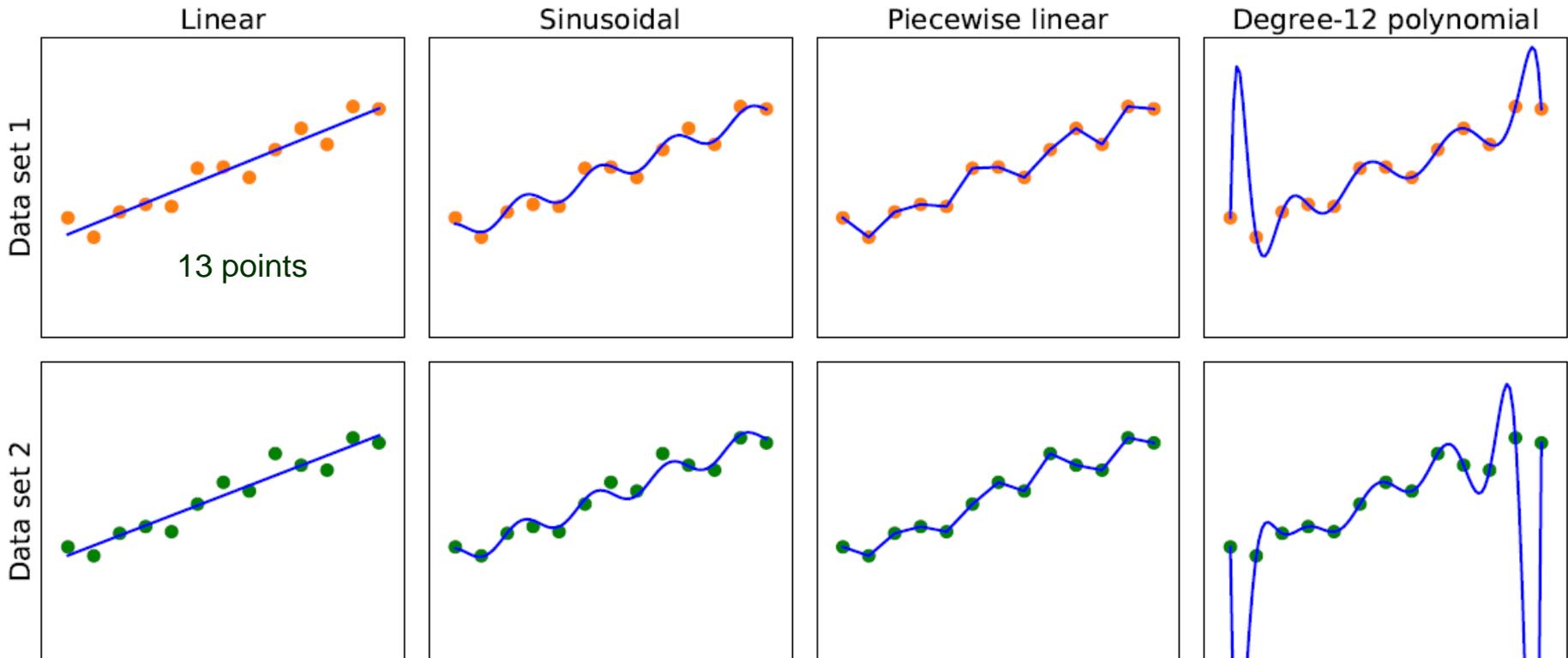
- ♣ With noise in data, we cannot expect an exact match with the ground truth, namely,  $h(x_i) = y_i$ , for  $1 \leq i \leq N$ .
- ♣ Instead, we look for a best-fit function  $h$  for which each  $h(x_i)$  is close  $y_i$ .
- ♣ The true measure of  $h$  is how it handles inputs it has not seen.

*Test set:* a second sample of  $(x_i, y_i)$  pairs

- ♣  $h$  generalizes well if it matches the test set with high accuracy.

# Fitting

(Least-squares fitting: <https://faculty.sites.iastate.edu/jia/files/inline-files/data-fit.pdf>)



$$h(x) = w_1x + w_0$$

$$h(x) = w_1x + \sin(w_0x)$$

$$h_i(x) = w_{i1}x + w_{i0},$$

$1 \leq i \leq 12$  such that  
 $h_j(x_{j+1}) = h_{j+1}(x_{j+1})$   
for  $1 \leq j \leq 11$ .

$$h(x) = \sum_{i=0}^{12} w_i x^i$$

# Fitting

(Least-squares fitting: <https://faculty.sites.iastate.edu/jia/files/inline-files/data-fit.pdf>)

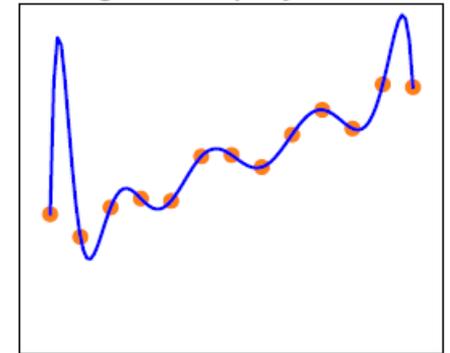
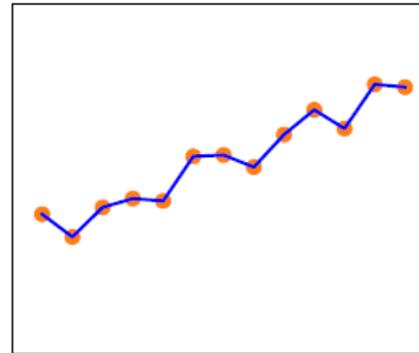
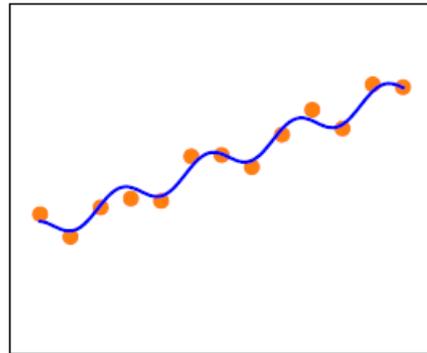
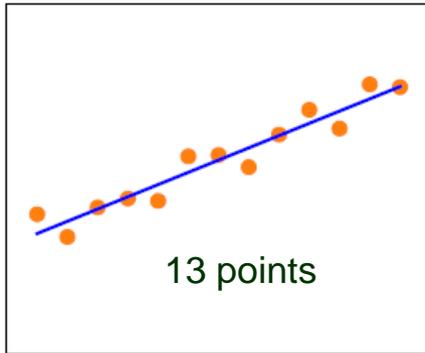
Linear

Sinusoidal

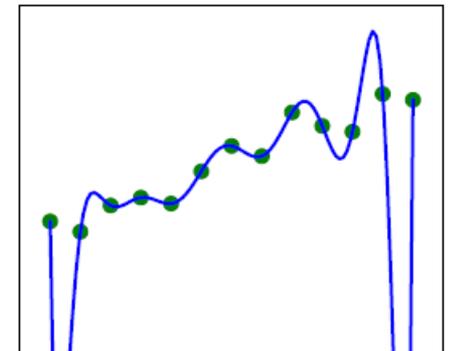
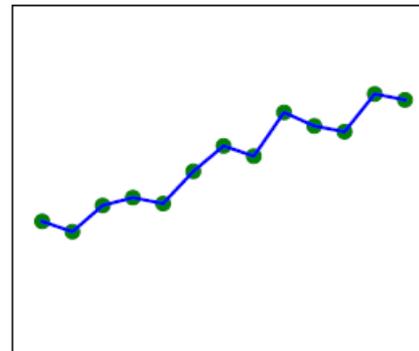
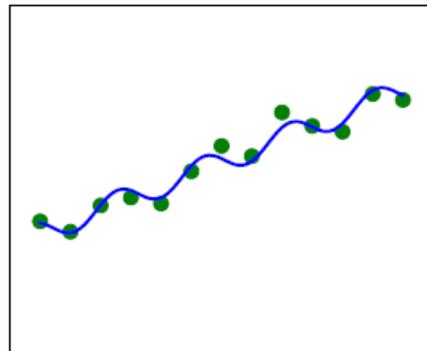
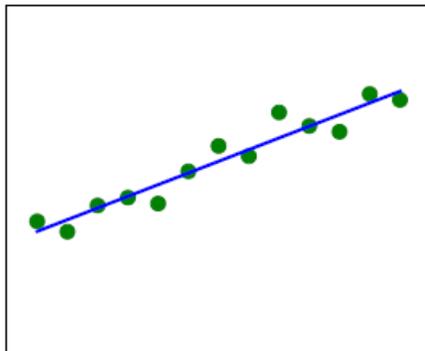
Piecewise linear

Degree-12 polynomial

Data set 1



Data set 2



$$h(x) = w_1x + w_0$$

$$h(x) = w_1x + \sin(w_0x)$$

$$h_i(x) = w_{i1}x + w_{i0},$$
$$1 \leq i \leq 12 \text{ such that}$$
$$h_j(x_{j+1}) = h_{j+1}(x_{j+1})$$
$$\text{for } 1 \leq j \leq 11.$$

$$h(x) = \sum_{i=0}^{12} w_i x^i$$

# Fitting

(Least-squares fitting: <https://faculty.sites.iastate.edu/jia/files/inline-files/data-fit.pdf>)

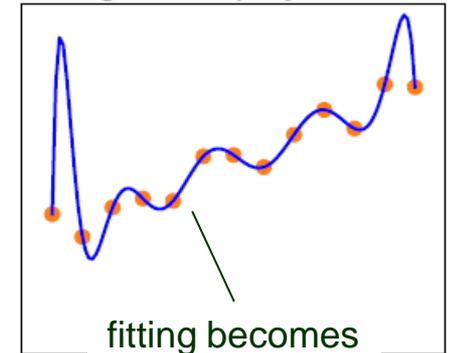
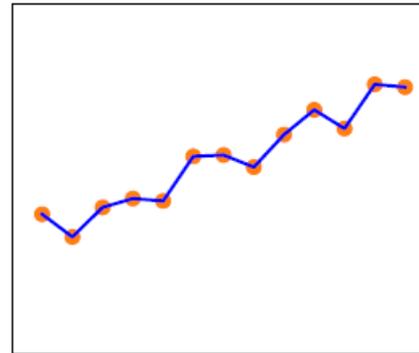
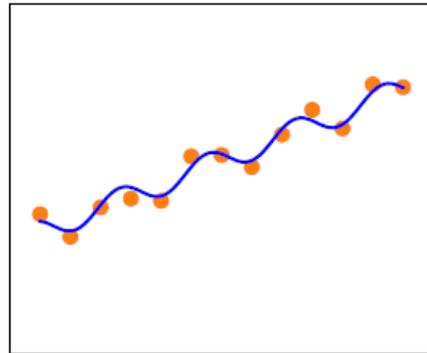
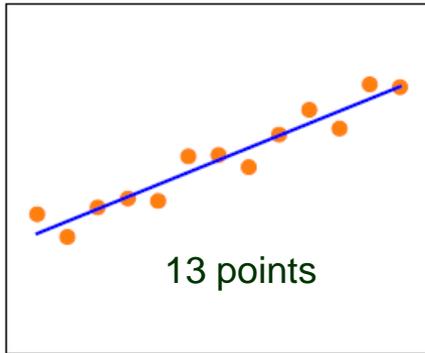
Linear

Sinusoidal

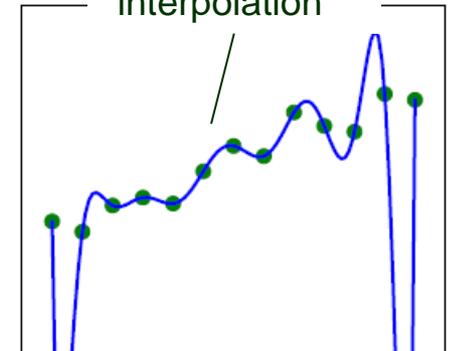
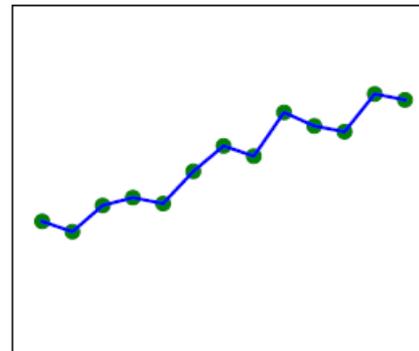
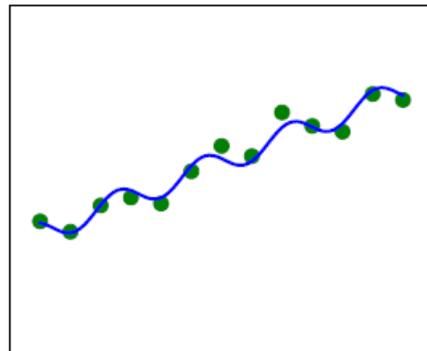
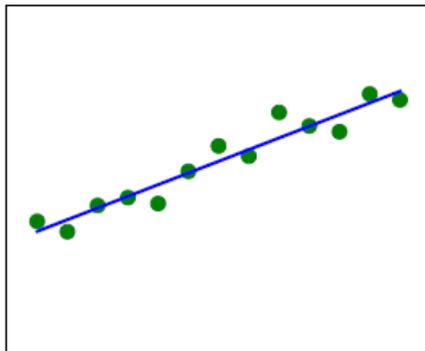
Piecewise linear

Degree-12 polynomial

Data set 1



Data set 2



$$h(x) = w_1x + w_0$$

$$h(x) = w_1x + \sin(w_0x)$$

$$h_i(x) = w_{i1}x + w_{i0},$$

$$1 \leq i \leq 12 \text{ such that}$$

$$h_j(x_{j+1}) = h_{j+1}(x_{j+1})$$

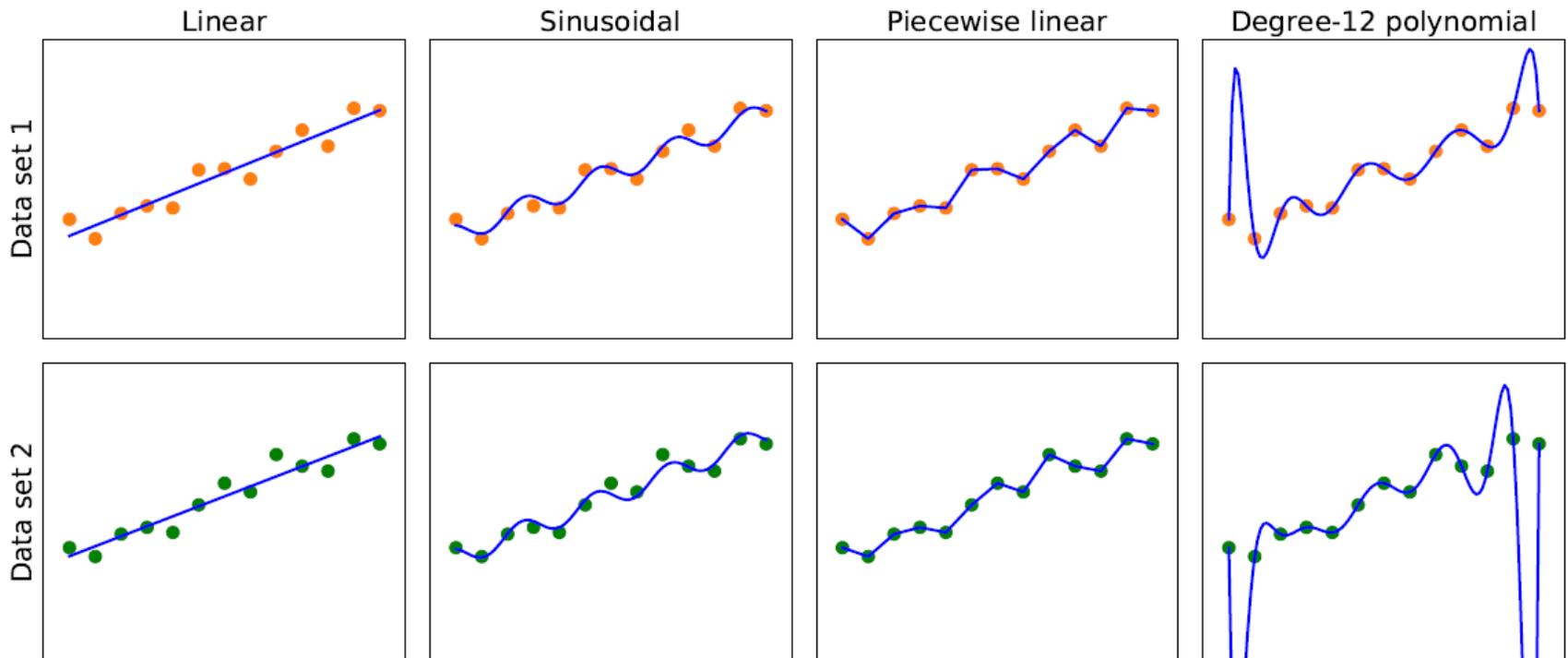
$$\text{for } 1 \leq j \leq 11.$$

$$h(x) = \sum_{i=0}^{12} w_i x^i$$

# Underfitting, Variance & Overfitting

- A hypothesis is *underfitting* when it fails to find a pattern in the data.
  - ◆ Such a hypothesis has high bias and low variance.

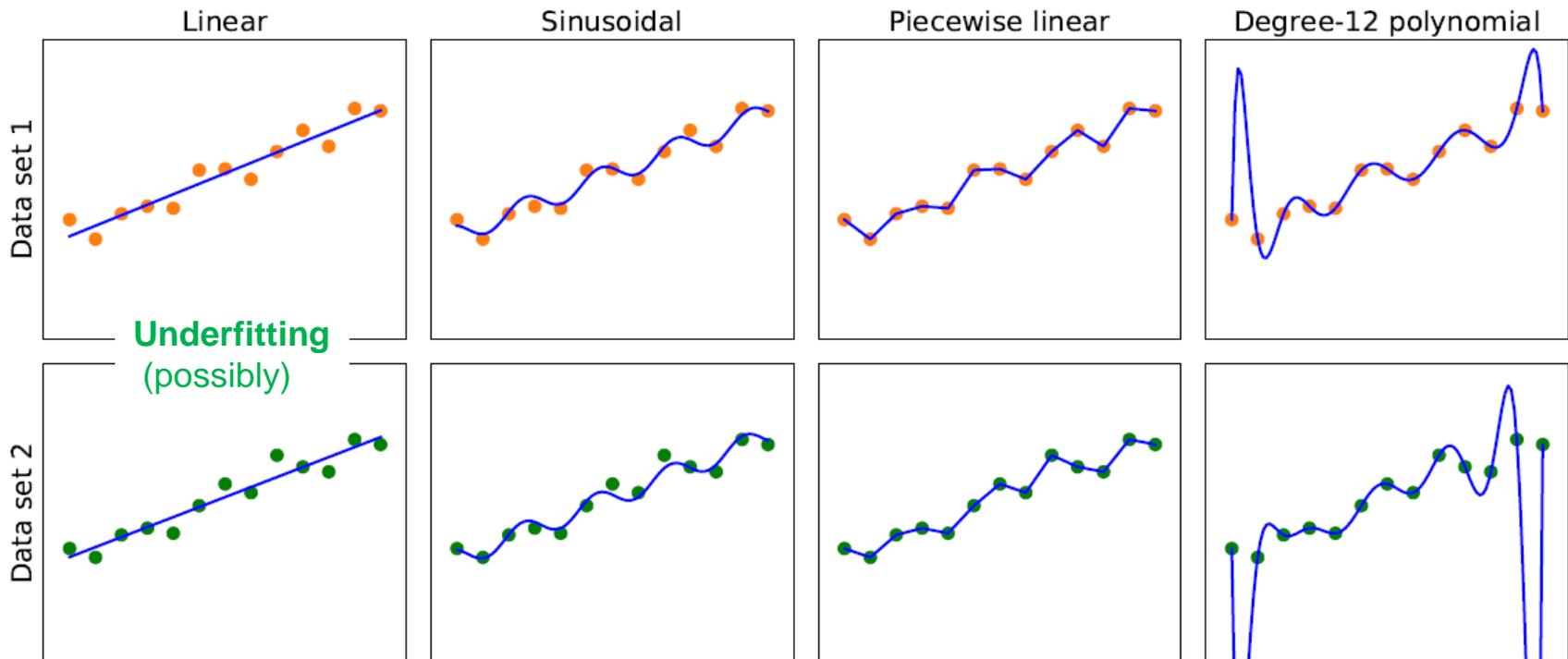
- ◆ Such a hypothesis has low bias and high variance.



# Underfitting, Variance & Overfitting

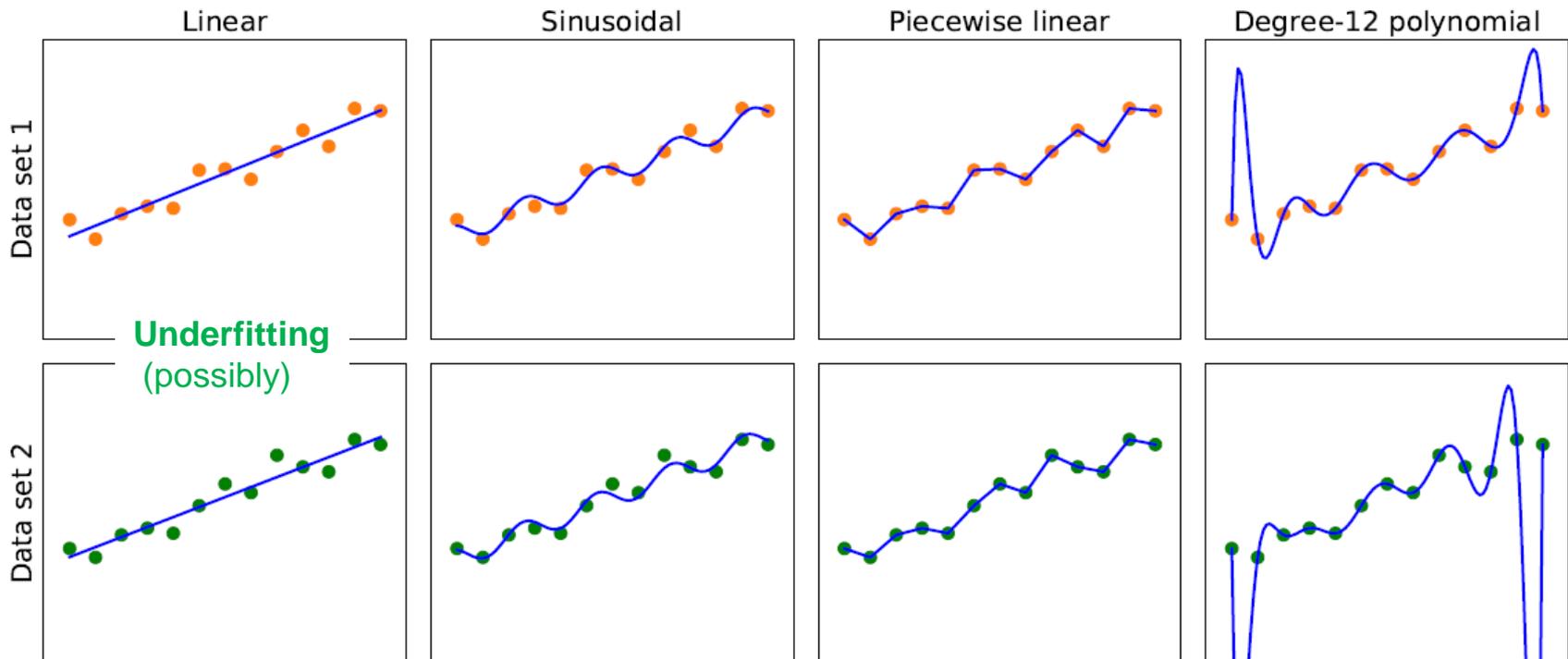
- A hypothesis is *underfitting* when it fails to find a pattern in the data.
  - ◆ Such a hypothesis has high bias and low variance.

- ◆ Such a hypothesis has low bias and high variance.



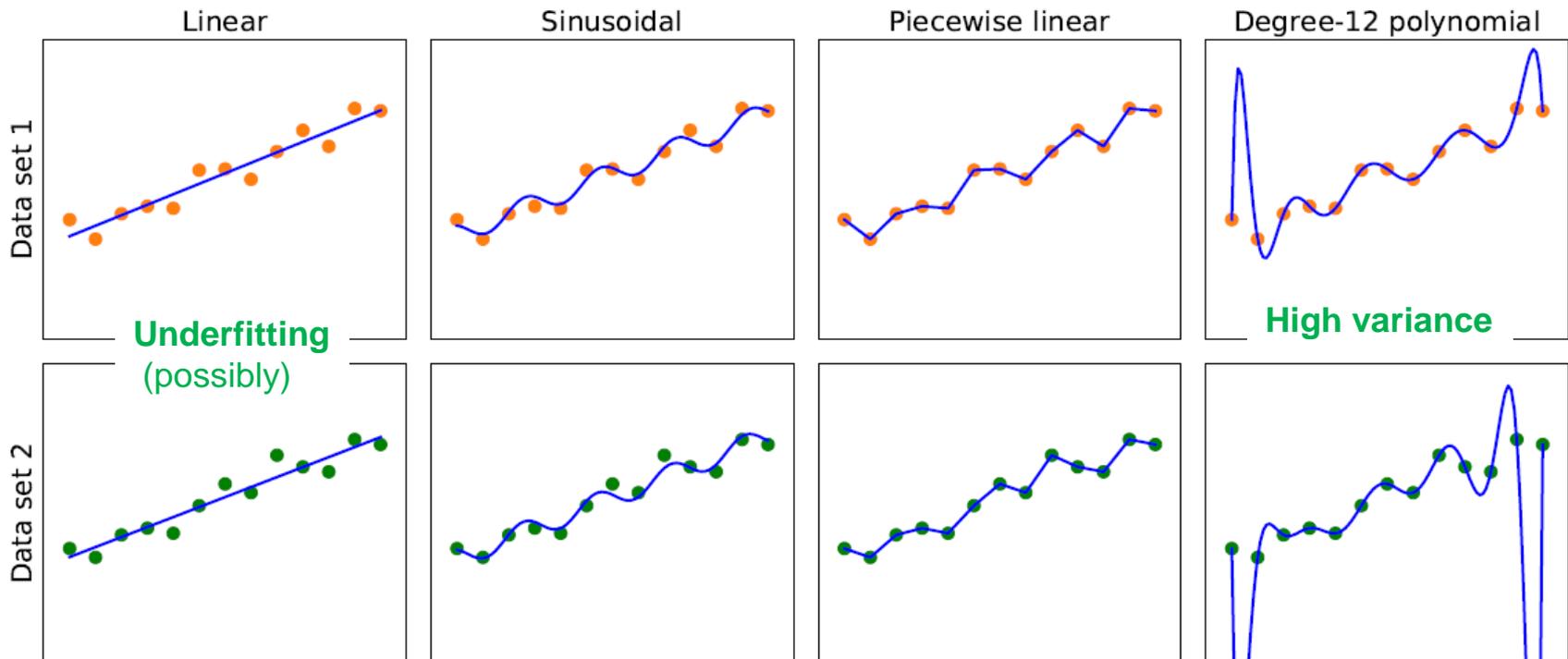
# Underfitting, Variance & Overfitting

- A hypothesis is *underfitting* when it fails to find a pattern in the data.
  - ◆ Such a hypothesis has high bias and low variance.
- *Variance* characterizes the amount of change in the hypothesis due to fluctuation in the training data.
  - ◆ Such a hypothesis has low bias and high variance.



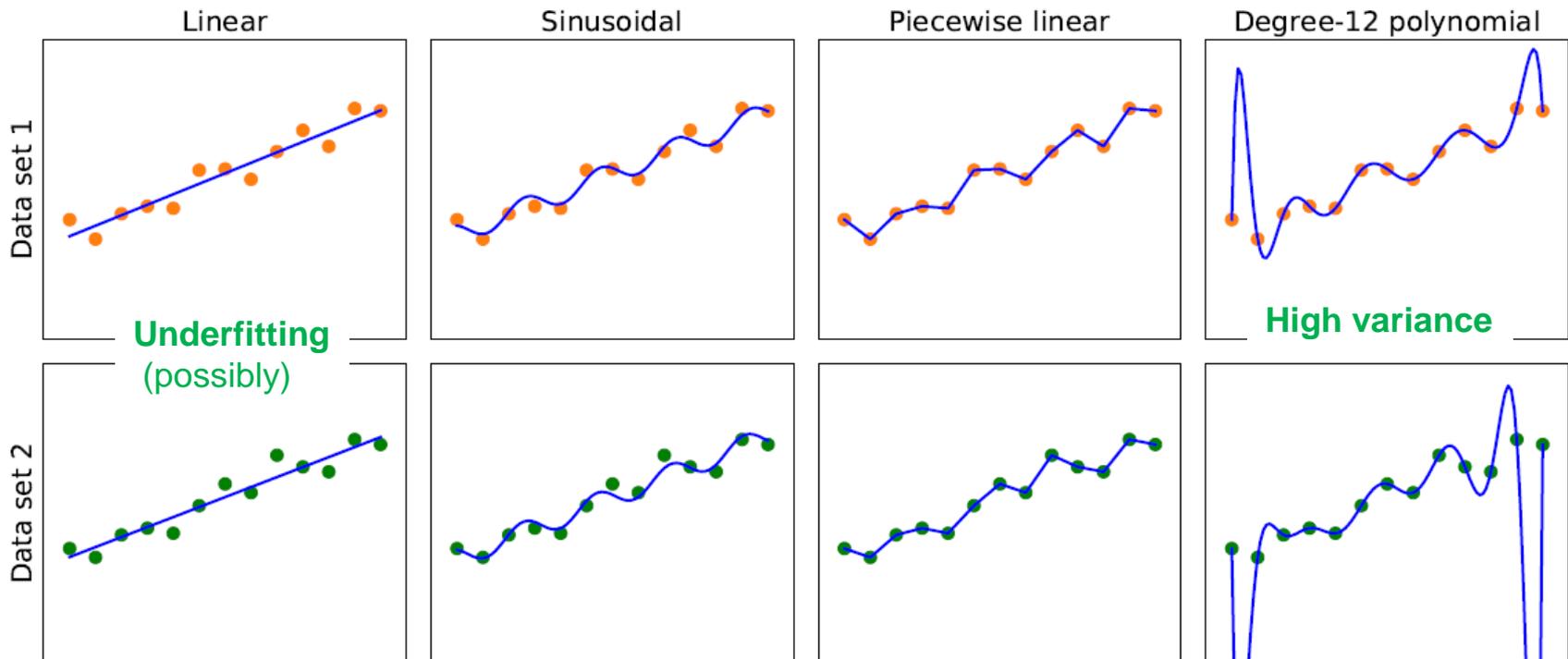
# Underfitting, Variance & Overfitting

- A hypothesis is *underfitting* when it fails to find a pattern in the data.
  - ◆ Such a hypothesis has high bias and low variance.
- *Variance* characterizes the amount of change in the hypothesis due to fluctuation in the training data.
  - ◆ Such a hypothesis has low bias and high variance.



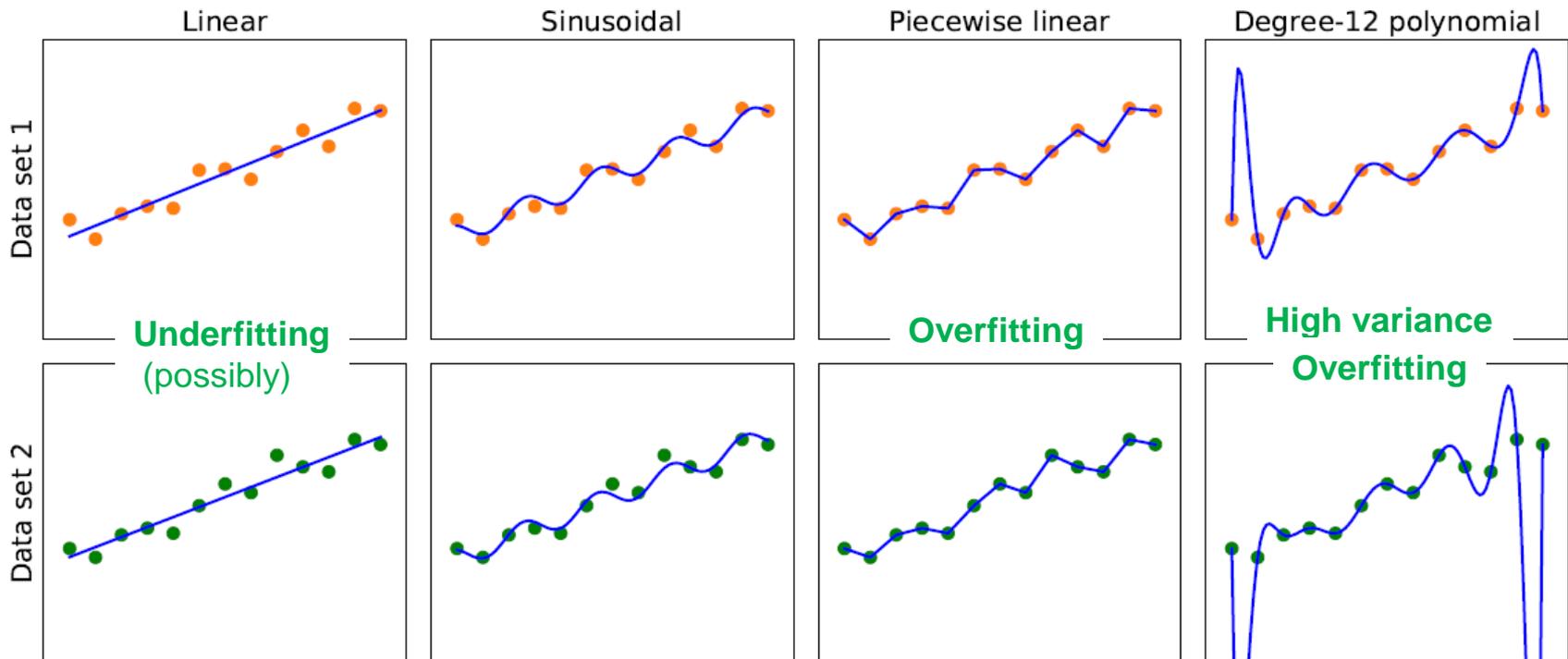
# Underfitting, Variance & Overfitting

- A hypothesis is *underfitting* when it fails to find a pattern in the data.
  - ◆ Such a hypothesis has high bias and low variance.
- *Variance* characterizes the amount of change in the hypothesis due to fluctuation in the training data.
- A hypothesis is *overfitting* when it pays too much attention to the particular training set. ◆ Such a hypothesis has low bias and high variance.



# Underfitting, Variance & Overfitting

- A hypothesis is *underfitting* when it fails to find a pattern in the data.
  - ♦ Such a hypothesis has high bias and low variance.
- *Variance* characterizes the amount of change in the hypothesis due to fluctuation in the training data.
- A hypothesis is *overfitting* when it pays too much attention to the particular training set. ♦ Such a hypothesis has low bias and high variance.



# Best Hypothesis

---

It depends on what we knew was represented by the data, or, what we were expecting from the data ...

# Best Hypothesis

---

It depends on what we knew was represented by the data, or, what we were expecting from the data ...

Choosing the hypothesis that is *most probable* given the data:

# Best Hypothesis

---

It depends on what we knew was represented by the data, or, what we were expecting from the data ...

Choosing the hypothesis that is *most probable* given the data:

$$h^* = \operatorname{argmax}_{h \in \mathcal{H}} P(h \mid \text{data})$$

# Best Hypothesis

---

It depends on what we knew was represented by the data, or, what we were expecting from the data ...

Choosing the hypothesis that is *most probable* given the data:

$$h^* = \operatorname{argmax}_{h \in \mathcal{H}} P(h \mid \text{data})$$

$$= \operatorname{argmax}_{h \in \mathcal{H}} P(\text{data} \mid h)P(h)$$

(Bayes' rule)

# Best Hypothesis

---

It depends on what we knew was represented by the data, or, what we were expecting from the data ...

Choosing the hypothesis that is *most probable* given the data:

$$h^* = \operatorname{argmax}_{h \in \mathcal{H}} P(h \mid \text{data})$$

$$= \operatorname{argmax}_{h \in \mathcal{H}} P(\text{data} \mid h)P(h) \quad (\text{Bayes' rule})$$

A simple hypothesis space  $\mathcal{H}$  is often preferred:

- ♠ The more expressiveness of  $\mathcal{H}$ , the higher the computational cost of finding a good hypothesis within that space.
- ♣ We will likely be using  $h$  for evaluations after we have learned it.

# Restaurant Waiting

---

**Problem** Decide whether to wait for a table at a restaurant.

# Restaurant Waiting

---

**Problem** Decide whether to wait for a table at a restaurant.

**Output:** a Boolean variable *WillWait* (true where we do wait for a table).

# Restaurant Waiting

---

**Problem** Decide whether to wait for a table at a restaurant.

**Output:** a Boolean variable *WillWait* (true where we do wait for a table).

**Input:** a vector of ten attributes, each with discrete values:

# Restaurant Waiting

---

**Problem** Decide whether to wait for a table at a restaurant.

**Output:** a Boolean variable *WillWait* (true where we do wait for a table).

**Input:** a vector of ten attributes, each with discrete values:

1. *Alternate*: whether there is a suitable alternative restaurant nearby.
2. *Bar*: whether the restaurant has a comfortable bar area to wait in.
3. *Fri/Sat*: true on Fridays and Saturdays.
4. *Hungry*: whether we are hungry right now.
5. *Patrons*: how many people are in the restaurant (values: *None*, *Some*, and *Full*).
6. *Price*: the restaurant's price range (\$, \$\$, \$\$\$).
7. *Raining*: whether it is raining outside.
8. *Reservation*: whether we made a reservation.
9. *Type*: the kind of restaurant (*French*, *Italian*, *Thai*, or *burger*).
10. *WaitEstimate*: host's wait estimate: 0-10, 10-30, 30-60, or >60 minutes.

# Restaurant Waiting

---

**Problem** Decide whether to wait for a table at a restaurant.

**Output:** a Boolean variable *WillWait* (true where we do wait for a table).

**Input:** a vector of ten attributes, each with discrete values:

1. *Alternate*: whether there is a suitable alternative restaurant nearby.
2. *Bar*: whether the restaurant has a comfortable bar area to wait in.
3. *Fri/Sat*: true on Fridays and Saturdays.
4. *Hungry*: whether we are hungry right now.
5. *Patrons*: how many people are in the restaurant (values: *None*, *Some*, and *Full*).
6. *Price*: the restaurant's price range (\$, \$\$, \$\$\$).
7. *Raining*: whether it is raining outside.
8. *Reservation*: whether we made a reservation.
9. *Type*: the kind of restaurant (*French*, *Italian*, *Thai*, or *burger*).
10. *WaitEstimate*: host's wait estimate: 0-10, 10-30, 30-60, or >60 minutes.

$2^6 \times 3^2 \times 4^2 = 9,216$  possible combinations of attribute values.

# Training Examples

Example	Input Attributes										Output
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
<b>x<sub>1</sub></b>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>0–10</i>	<i>y<sub>1</sub> = Yes</i>
<b>x<sub>2</sub></b>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>30–60</i>	<i>y<sub>2</sub> = No</i>
<b>x<sub>3</sub></b>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Some</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>0–10</i>	<i>y<sub>3</sub> = Yes</i>
<b>x<sub>4</sub></b>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Thai</i>	<i>10–30</i>	<i>y<sub>4</sub> = Yes</i>
<b>x<sub>5</sub></b>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>&gt;60</i>	<i>y<sub>5</sub> = No</i>
<b>x<sub>6</sub></b>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Italian</i>	<i>0–10</i>	<i>y<sub>6</sub> = Yes</i>
<b>x<sub>7</sub></b>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>0–10</i>	<i>y<sub>7</sub> = No</i>
<b>x<sub>8</sub></b>	<i>No</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Thai</i>	<i>0–10</i>	<i>y<sub>8</sub> = Yes</i>
<b>x<sub>9</sub></b>	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>&gt;60</i>	<i>y<sub>9</sub> = No</i>
<b>x<sub>10</sub></b>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>Italian</i>	<i>10–30</i>	<i>y<sub>10</sub> = No</i>
<b>x<sub>11</sub></b>	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>0–10</i>	<i>y<sub>11</sub> = No</i>
<b>x<sub>12</sub></b>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>30–60</i>	<i>y<sub>12</sub> = Yes</i>

♣ The correct output is given for only 12 out of 9,216 examples.

# Training Examples

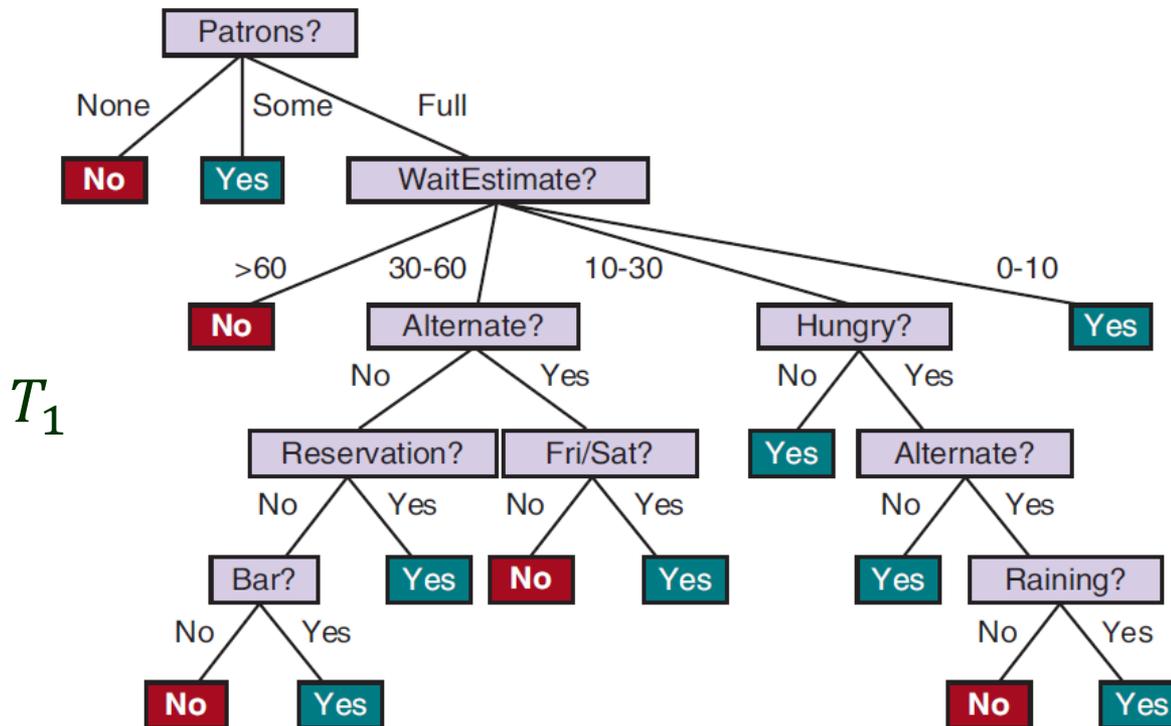
Example	Input Attributes										Output
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
$x_1$	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>0-10</i>	$y_1 = \text{Yes}$
$x_2$	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>30-60</i>	$y_2 = \text{No}$
$x_3$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Some</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>0-10</i>	$y_3 = \text{Yes}$
$x_4$	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Thai</i>	<i>10-30</i>	$y_4 = \text{Yes}$
$x_5$	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>&gt;60</i>	$y_5 = \text{No}$
$x_6$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Italian</i>	<i>0-10</i>	$y_6 = \text{Yes}$
$x_7$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>0-10</i>	$y_7 = \text{No}$
$x_8$	<i>No</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Thai</i>	<i>0-10</i>	$y_8 = \text{Yes}$
$x_9$	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>&gt;60</i>	$y_9 = \text{No}$
$x_{10}$	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>Italian</i>	<i>10-30</i>	$y_{10} = \text{No}$
$x_{11}$	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>0-10</i>	$y_{11} = \text{No}$
$x_{12}$	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>30-60</i>	$y_{12} = \text{Yes}$

- ♣ The correct output is given for only 12 out of 9,216 examples.
- ♣ We need to make our best guess at the missing 9,204 output values.

# III. What is a Decision Tree?

A *decision tree* maps a set of attribute values to a “decision”

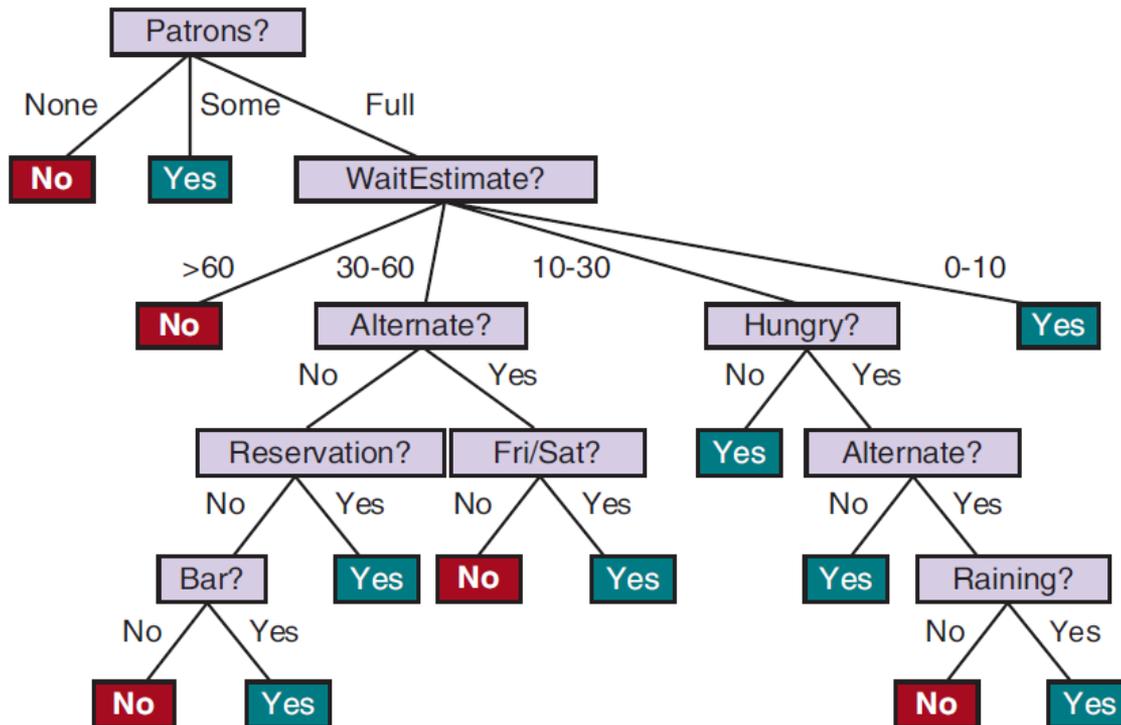
- ◆ It performs a sequence of tests, starting at the root and descending down to a leaf (which specifies the output value).
- ◆ Each internal node corresponds to a test of the value of an attribute.



A decision tree for the restaurant waiting problem

# Boolean Classification

- Inputs are discrete values.
- Outputs are either *true* (a positive example) or *false* (a negative one).

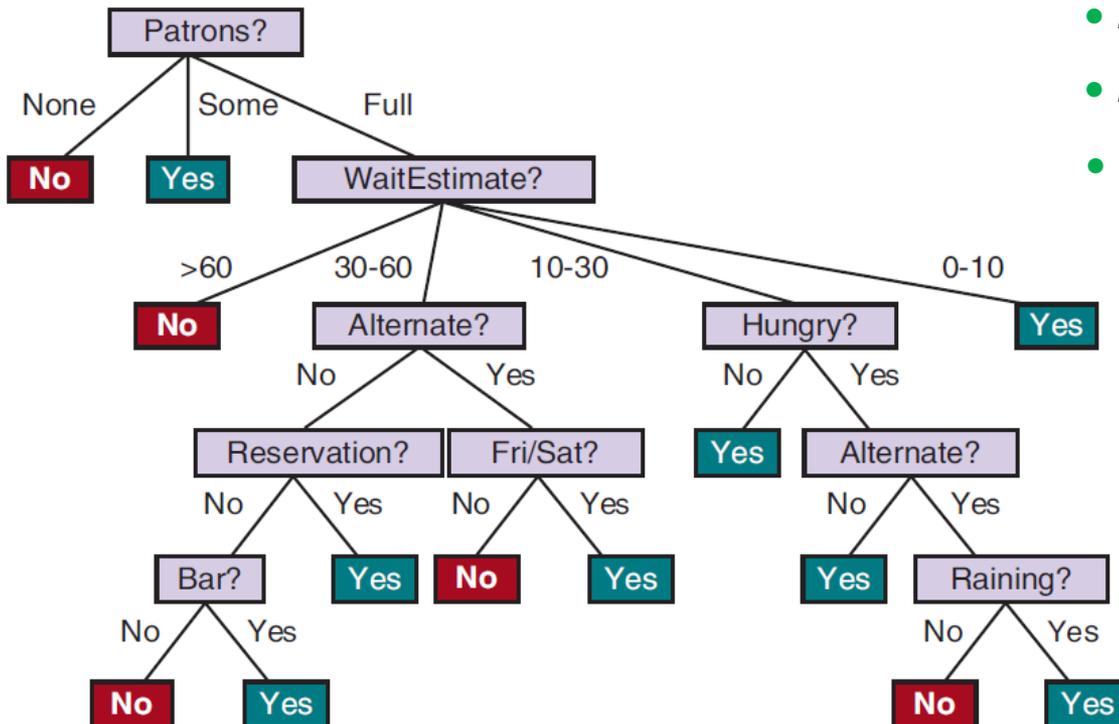


# Boolean Classification

- Inputs are discrete values.
- Outputs are either *true* (a positive example) or *false* (a negative one).

Positive examples:

- *Patrons = Some*
- *Patrons = Full*  $\wedge$  *WaitEstimate = 0-10*
- *Patrons = Full*  
 $\wedge$  *WaitEstimate = 10-30*  
 $\wedge$  *Hungry = Yes*  
 $\wedge$  *Alternate = No*

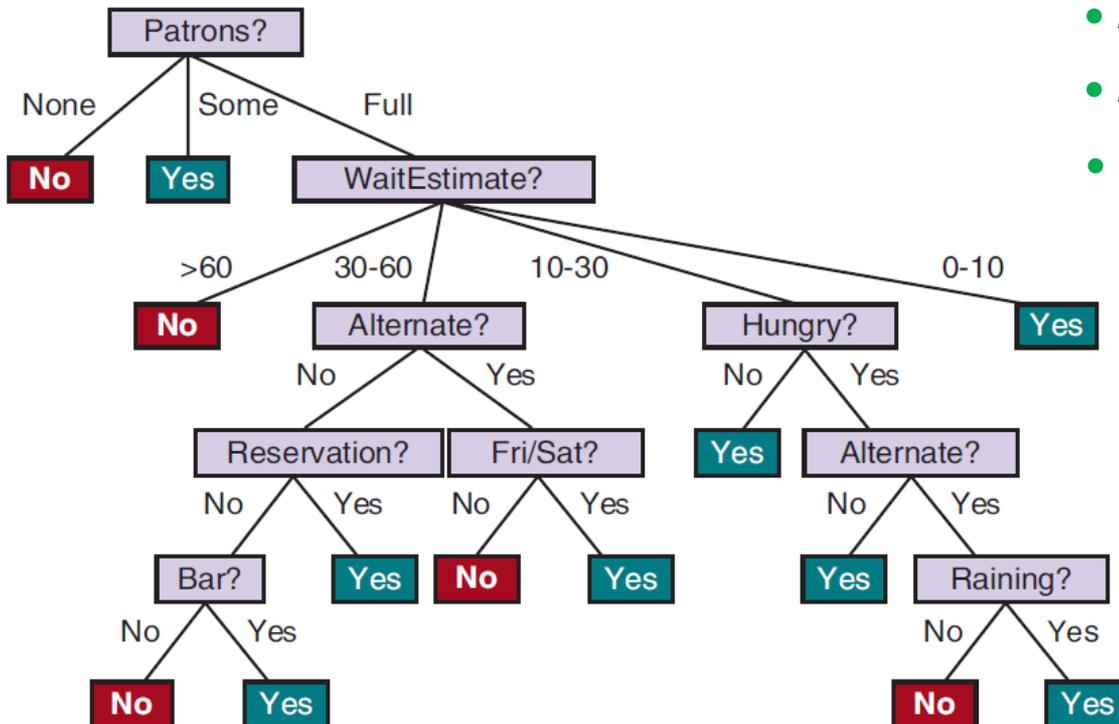


# Boolean Classification

- Inputs are discrete values.
- Outputs are either *true* (a positive example) or *false* (a negative one).

## Positive examples:

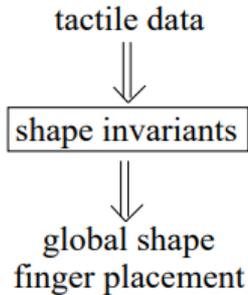
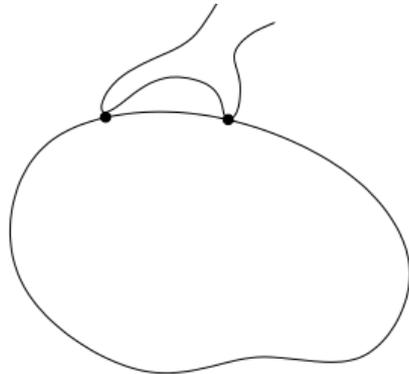
- *Patrons = Some*
- *Patrons = Full*  $\wedge$  *WaitEstimate = 0-10*
- *Patrons = Full*  
 $\wedge$  *WaitEstimate = 10-30*  
 $\wedge$  *Hungry = Yes*  
 $\wedge$  *Alternate = No*



## Negative examples:

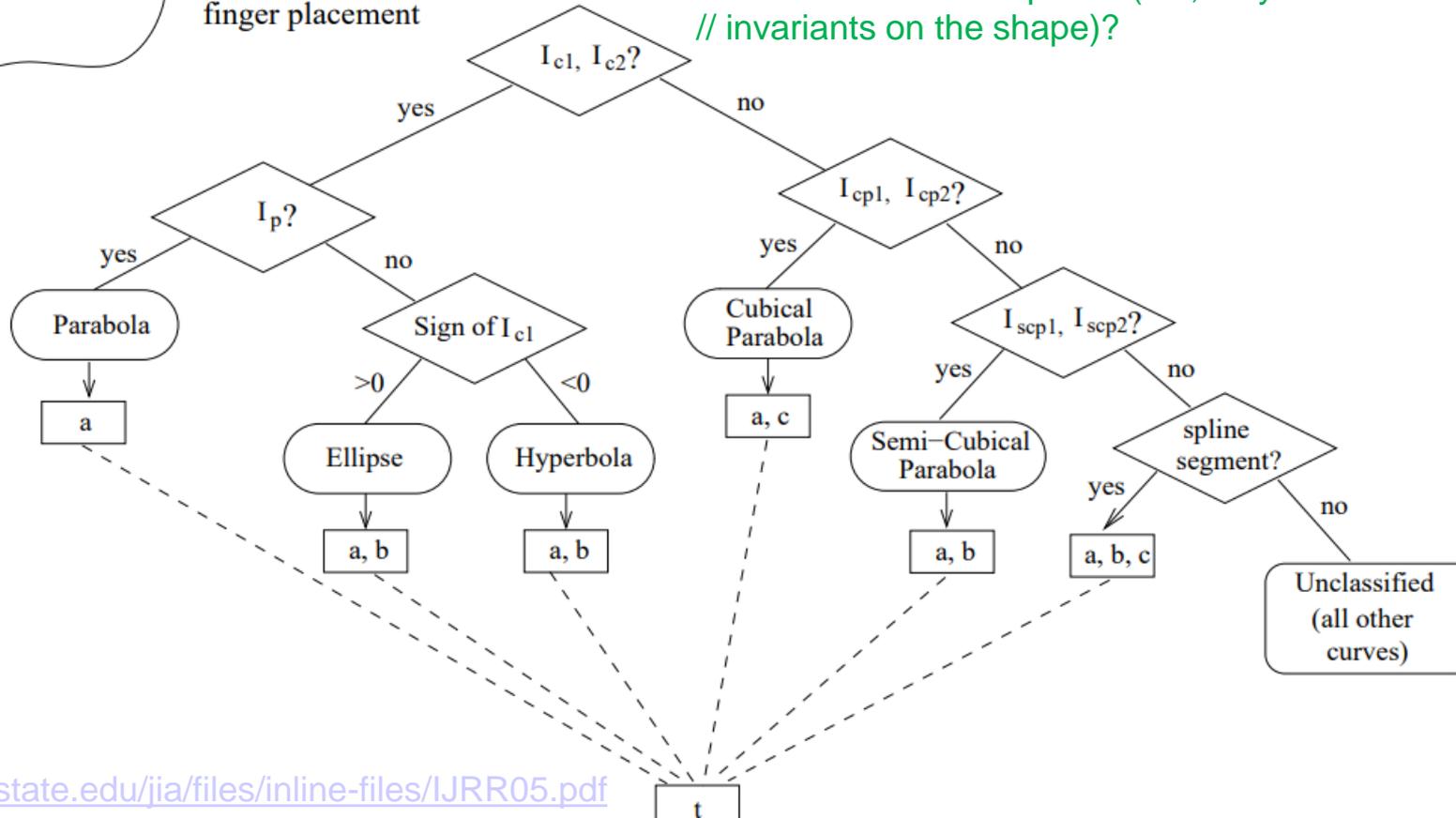
- *Patrons = No*
- *Patrons = Full*  
 $\wedge$  *WaitEstimate > 60*

# Shape Recognition Tree



A robotic hand can recognize curved shapes from differential invariants (expressions of curvature, torsion, etc.) constructed over tactile data.

// Values of  $I_{c1}$  and  $I_{c2}$  do not vary  
 // at different contour points (i.e., they are  
 // invariants on the shape)?



# Why Decision Trees?

---

- ◆ They yield nice, concise, and understandable results.
- ◆ They represent simple rules for classifying instances that are described by discrete attribute values.
- ◆ Decision tree learning algorithms are relatively efficient
  - *linear* in the size of the decision tree and the size of the data set.
- ◆ Decision trees are often among the first to be tried on a new data set.
- ♠ They are not good with real-valued attributes.

# Expressiveness of a Decision Tree

---

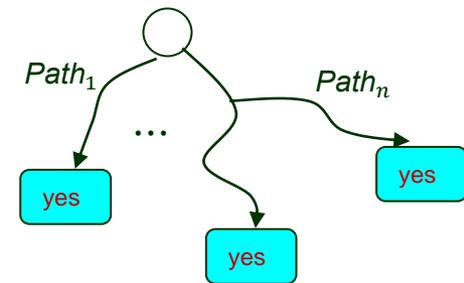
Equivalent logical statement of a decision tree:

$$\text{Output} \Leftrightarrow \text{Path}_1 \vee \text{Path}_2 \vee \dots$$

(one statement for Output = Yes  
and another for Output = No)

where

$$\text{Path}_i = (A_{i1} = v_{i1} \wedge A_{i2} = v_{i2} \wedge \dots)$$



# Expressiveness of a Decision Tree

Equivalent logical statement of a decision tree:

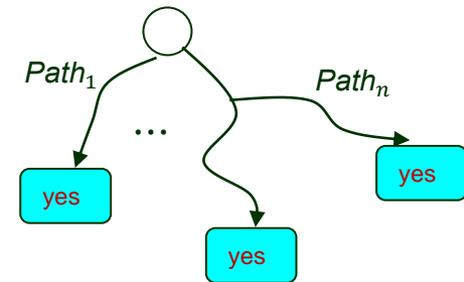
$$\text{Output} \Leftrightarrow \text{Path}_1 \vee \text{Path}_2 \vee \dots$$

(one statement for Output = Yes  
and another for Output = No)

where

$$\text{Path}_i = (A_{i1} = v_{i1} \wedge A_{i2} = v_{i2} \wedge \dots)$$

|            |  
attribute value



# Expressiveness of a Decision Tree

Equivalent logical statement of a decision tree:

$$\text{Output} \Leftrightarrow \text{Path}_1 \vee \text{Path}_2 \vee \dots$$

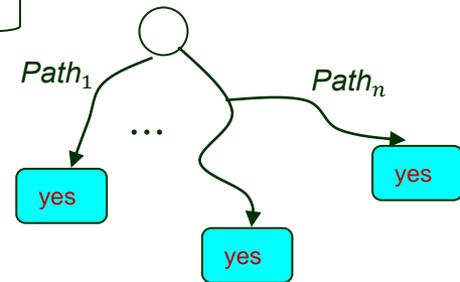
(one statement for Output = Yes  
and another for Output = No)

where

$$\text{Path}_i = (A_{i1} = v_{i1} \wedge A_{i2} = v_{i2} \wedge \dots)$$

|            |  
attribute value

Disjunctive normal form (DNF)



# Expressiveness of a Decision Tree

Equivalent logical statement of a decision tree:

$$\text{Output} \Leftrightarrow \text{Path}_1 \vee \text{Path}_2 \vee \dots$$

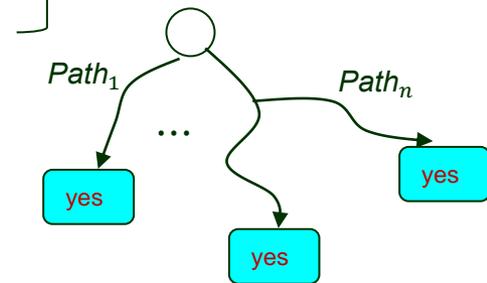
(one statement for Output = Yes  
and another for Output = No)

where

$$\text{Path}_i = (A_{i1} = v_{i1} \wedge A_{i2} = v_{i2} \wedge \dots)$$

|            |  
attribute value

Disjunctive normal form (DNF)



Any sentence in propositional logic can be converted into a DNF.

# Expressiveness of a Decision Tree

Equivalent logical statement of a decision tree:

$$\text{Output} \Leftrightarrow \text{Path}_1 \vee \text{Path}_2 \vee \dots$$

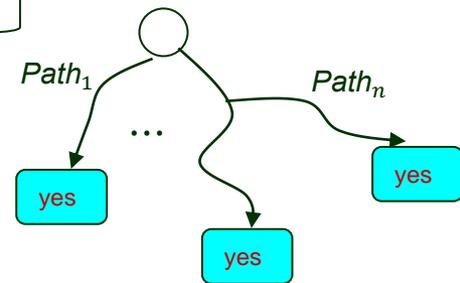
(one statement for Output = Yes  
and another for Output = No)

where

$$\text{Path}_i = (A_{i1} = v_{i1} \wedge A_{i2} = v_{i2} \wedge \dots)$$

attribute value

Disjunctive normal form (DNF)



Any sentence in propositional logic can be converted into a DNF.



Any sentence in propositional logic can be expressed as a decision tree.