

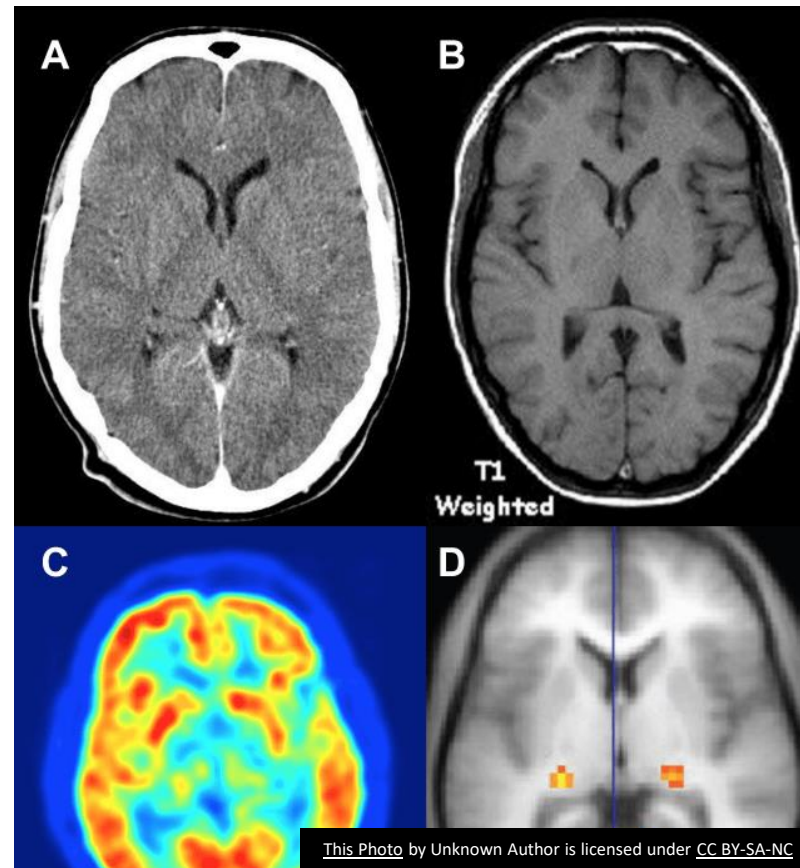
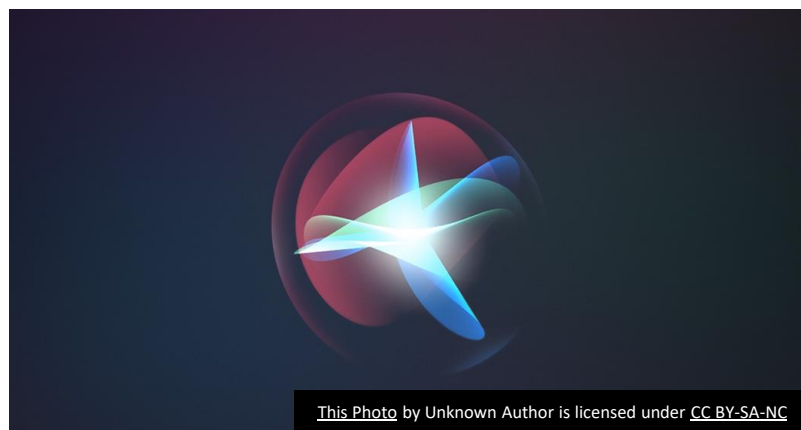
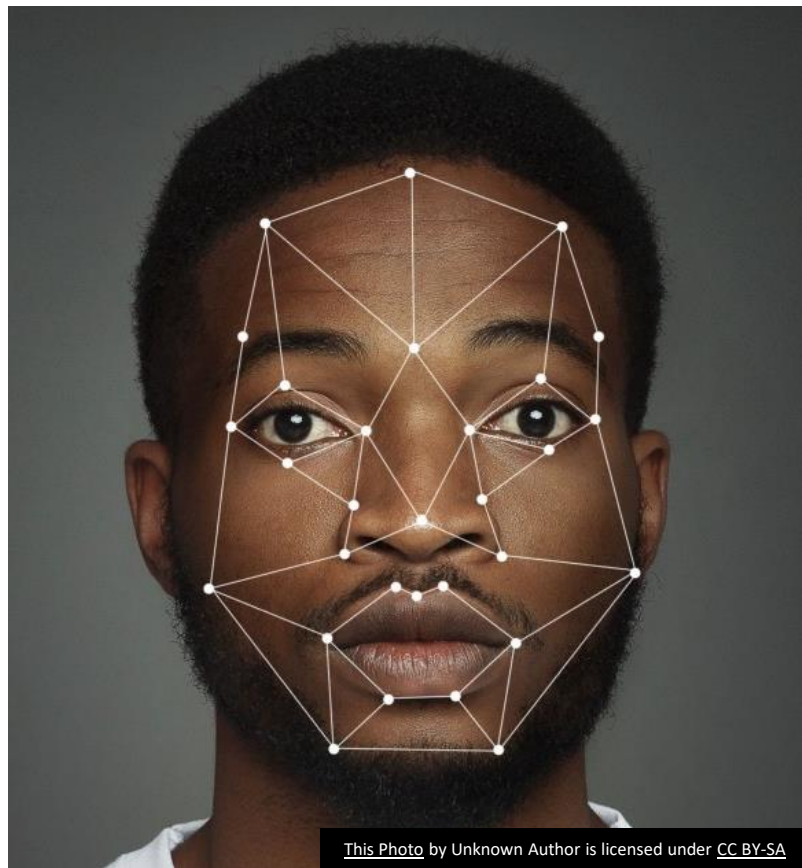
SOC 121D: People Analytics

Austin van Loon

Machine Learning 1

Overview for this week

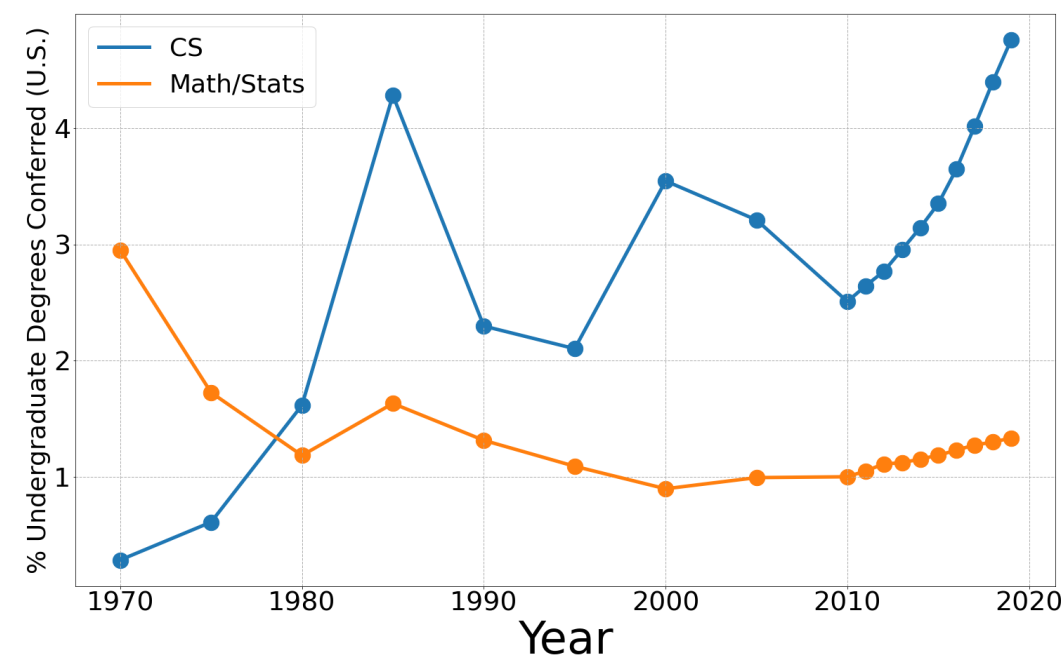
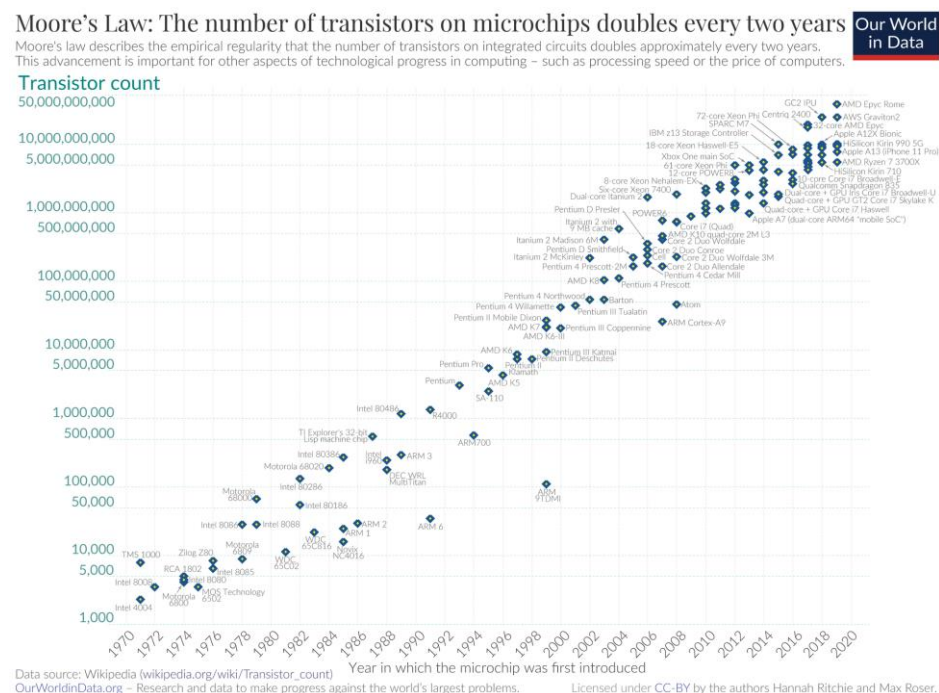
- Today – Machine learning 1
 - What is machine learning?
 - Machine learning basics and terminology
 - Introduce three families of algorithms
- Thursday – Machine Learning 2
 - Sources of bias in machine learning
 - Conceptualizing/Measuring Fairness
 - Predictive algorithms and inequality
 - In-class exercise (hopefully)
- Methods Module 1
 - Available now on Canvas
 - Provides more technical walk-through of methods and hands-on examples
 - Coding requirements are minimal (mostly copy and pasting)
 - Turn in before start of class next Tuesday, 7/5



What is machine learning (ML)?

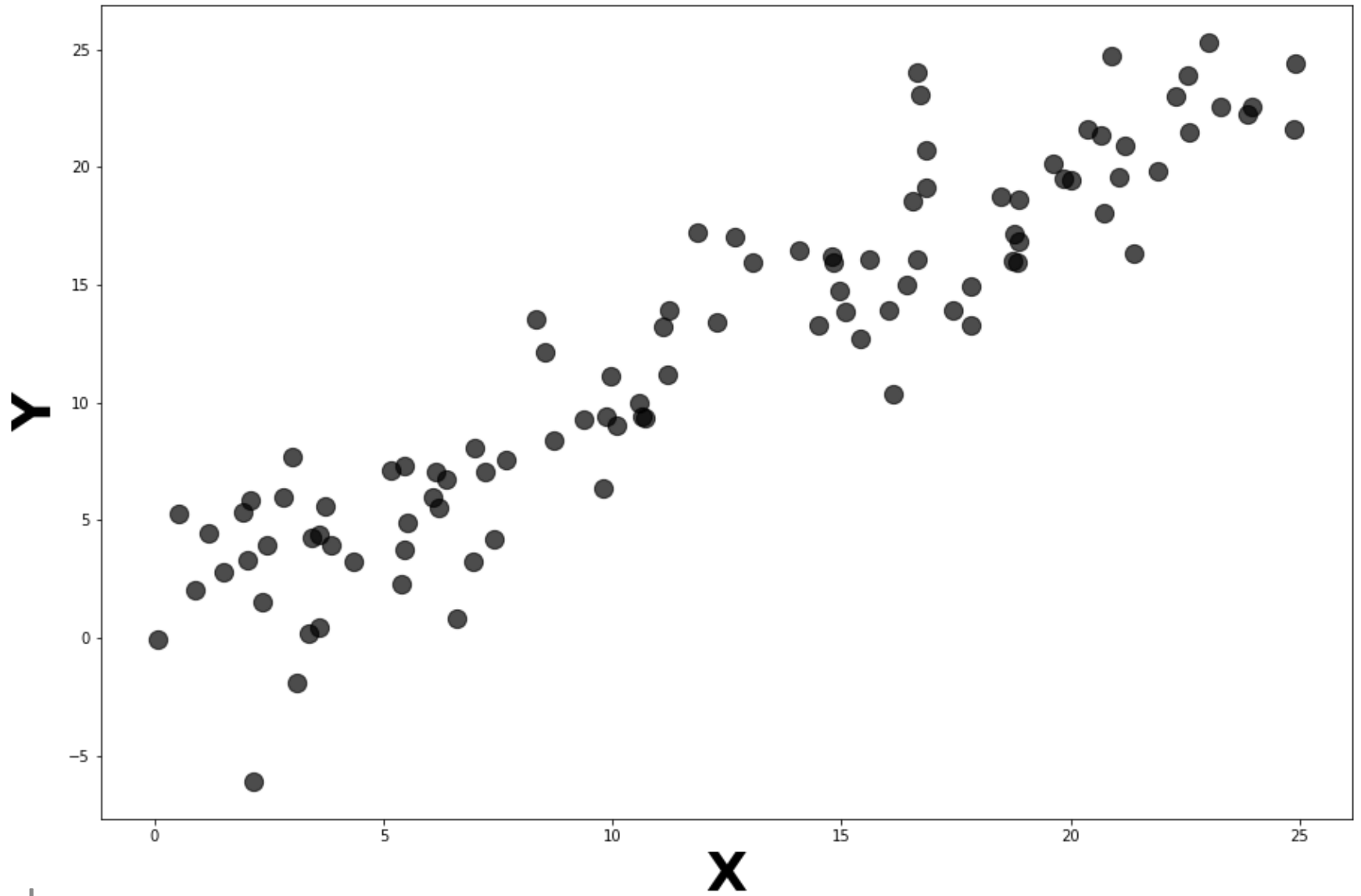
Why the rise of machine learning?

- Increased computational resources (Moore's Law)
- Prevalence of specialized knowledge (CS departments, YouTube)
- Strong economic incentives
- Highly legitimate (maybe too legitimate...)

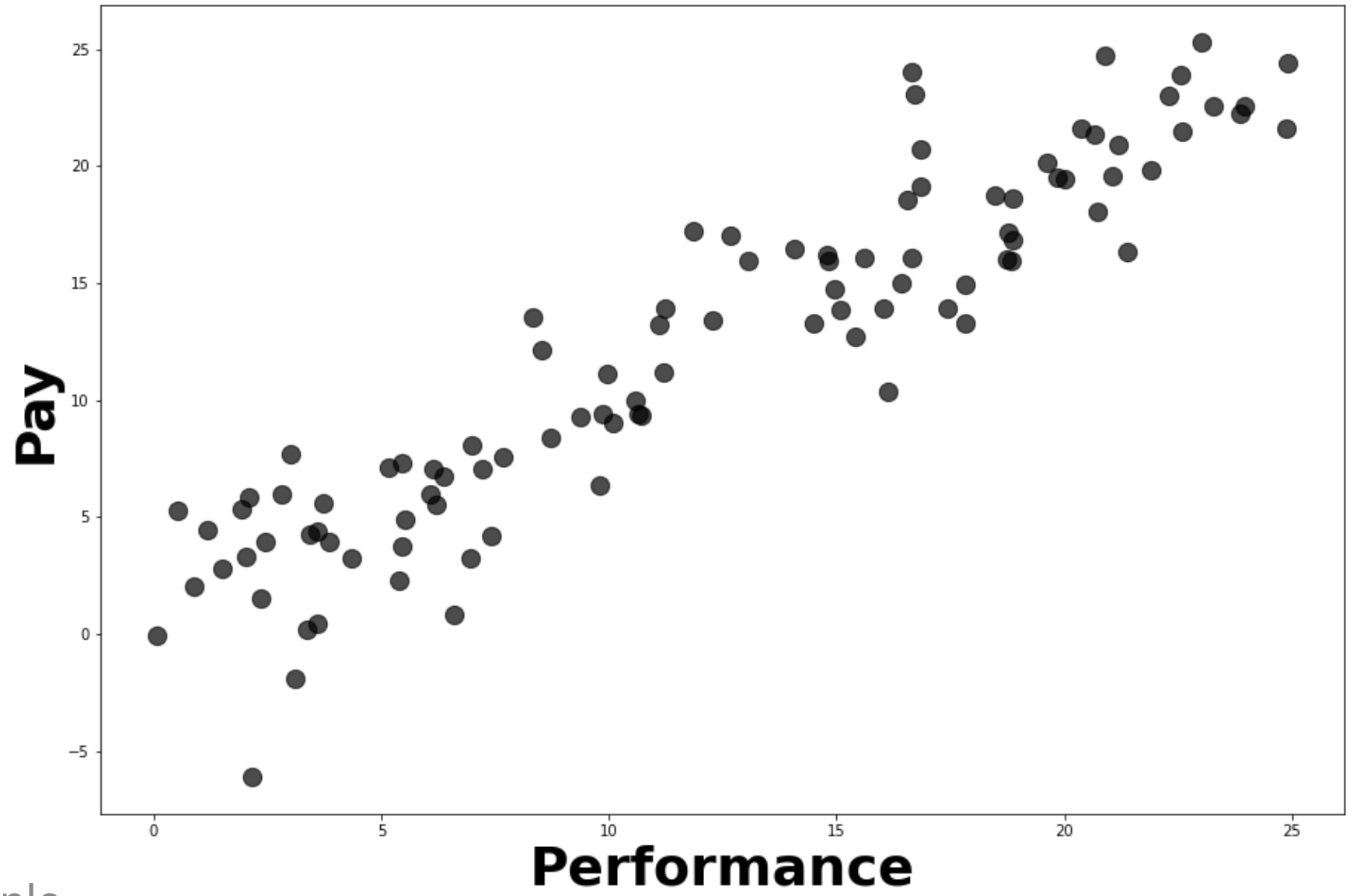


Philosophy of ML

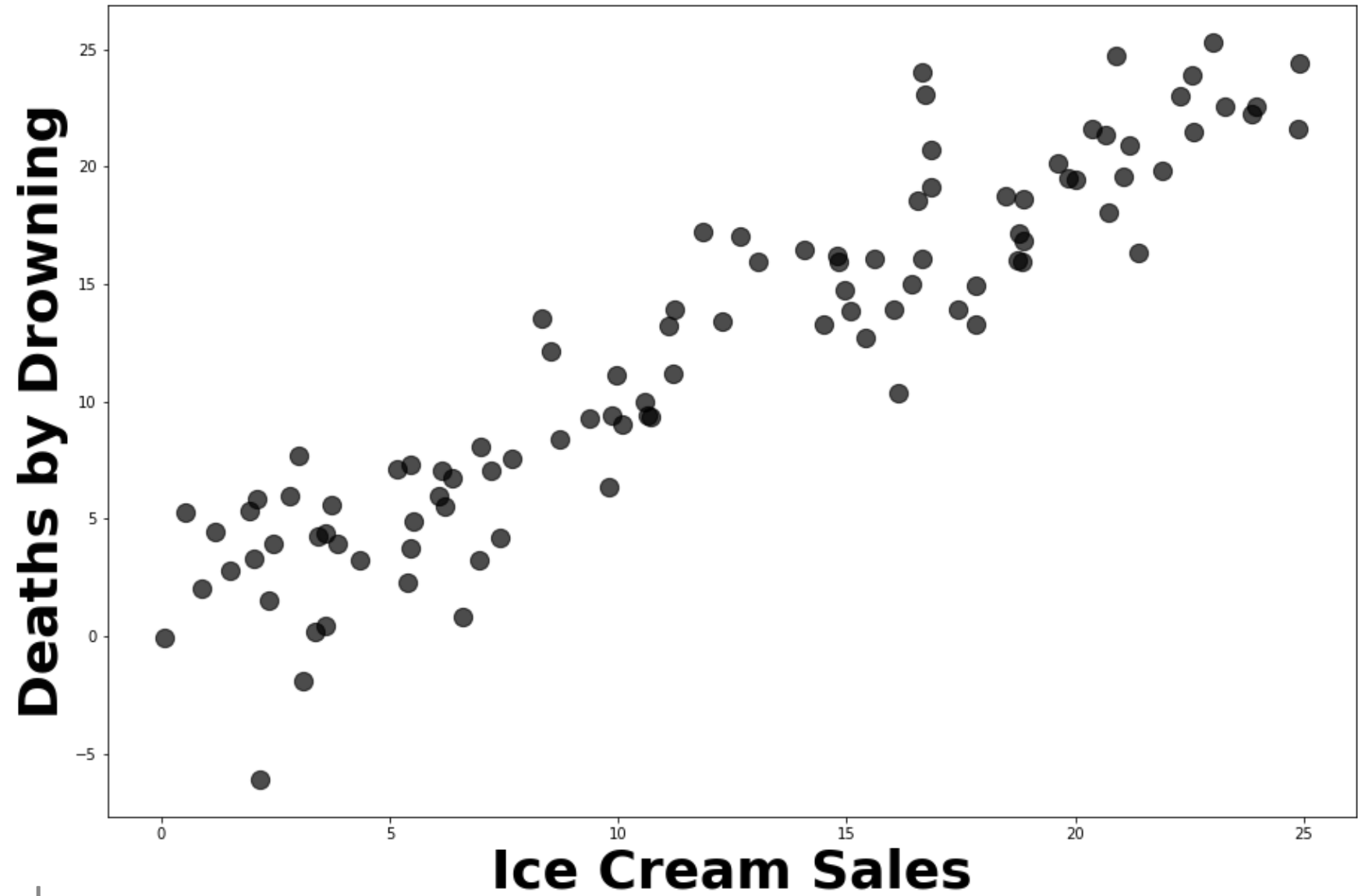
- **Correlation** is the empirical tendency for two or more variables to “move together” in a specific way
- **Causation** is the relationship between two variables that changing one will change the other.
- **Prediction** is the ability to reliably estimate the value of one variable from another (better than chance).



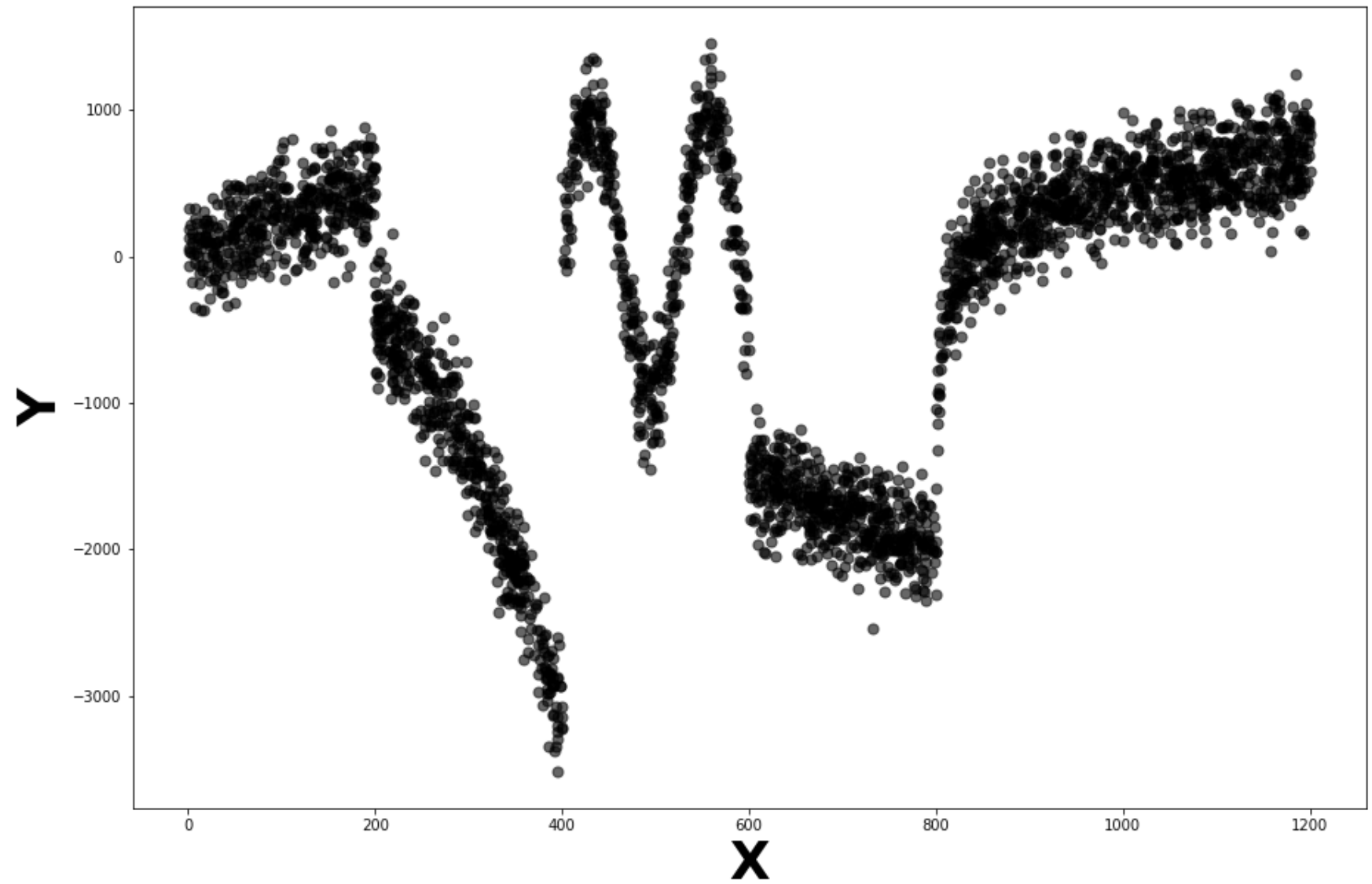
Correlation example



Correlation example



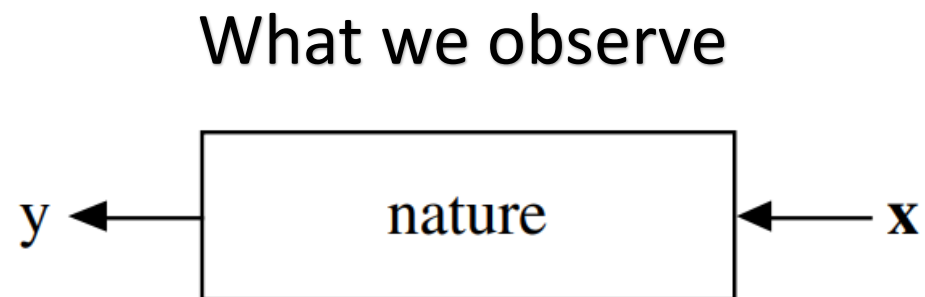
Correlation example



Prediction example

Statistical Modeling: The Two Cultures

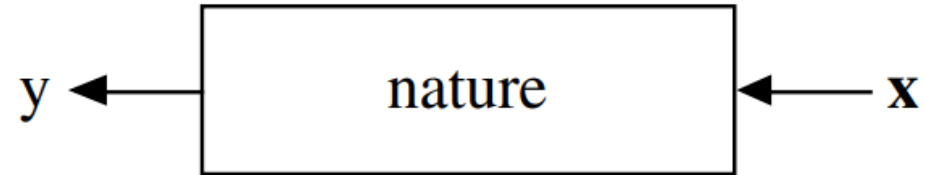
Leo Breiman



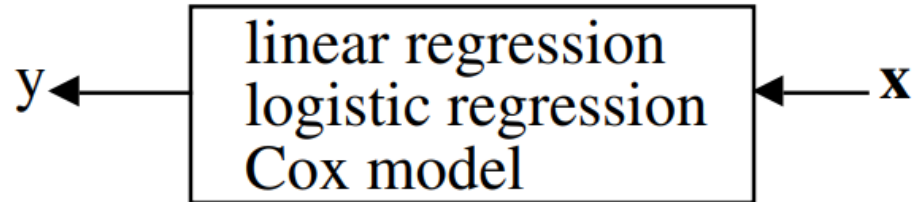
Statistical Modeling: The Two Cultures

Leo Breiman

What we observe



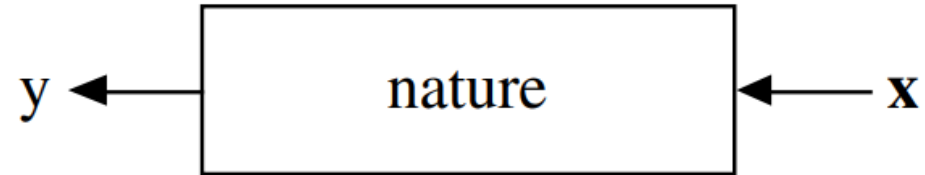
Data modeling culture



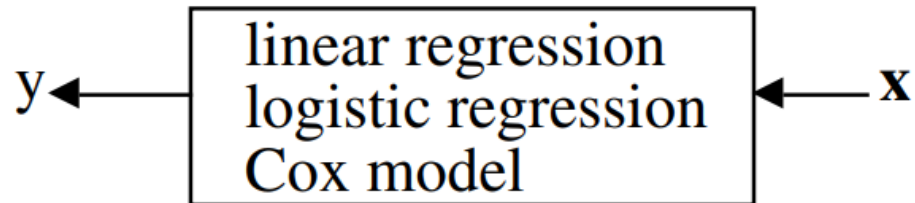
Statistical Modeling: The Two Cultures

Leo Breiman

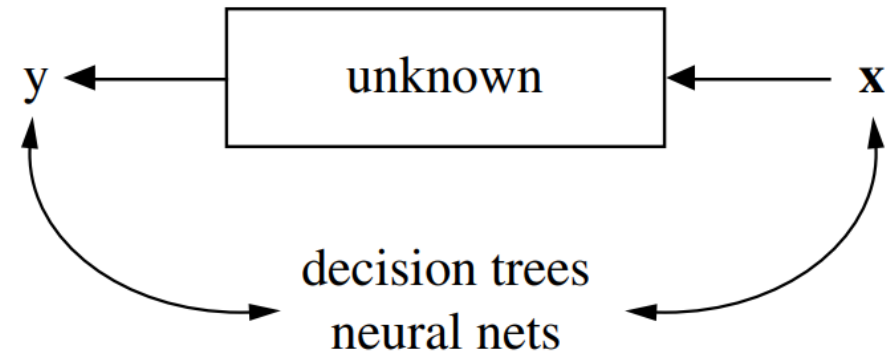
What we observe



Data modeling culture



Algorithmic modeling culture



Some ML vocabulary

- **Observation** – a record or datum (e.g., a single employee)
- **Dataset** – A collection of observations
- **Feature(s)** – the characteristic(s) of observations used to predict something
- **Outcome(s)** – the characteristic(s) of observations we want to predict

Some ML vocabulary

Name	Pay	Married?	Job	Age
Rachel	\$32k	No	Server	24
Ross	\$140k	No	Professor	26
Joey	\$350k	No	Actor	25
Monica	\$70k	Yes	Chef	24
Chandler	\$120k	Yes	Data Specialist	26
Phoebe	\$400k	No	Songwriter	26

Some ML vocabulary


Dataset



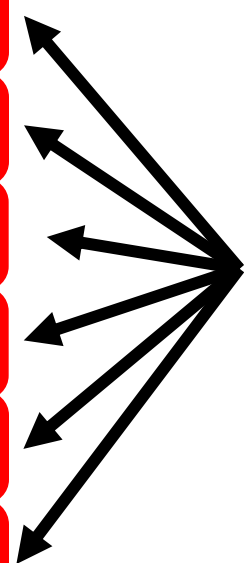
Name	Pay	Married?	Job	Age
Rachel	\$32k	No	Server	24
Ross	\$140k	No	Professor	26
Joey	\$350k	No	Actor	25
Monica	\$70k	Yes	Chef	24
Chandler	\$120k	Yes	Data Specialist	26
Phoebe	\$400k	No	Songwriter	26

Some ML vocabulary

Name	Pay	Married?	Job	Age
Rachel	\$32k	No	Server	24
Ross	\$140k	No	Professor	26
Joey	\$350k	No	Actor	25
Monica	\$70k	Yes	Chef	24
Chandler	\$120k	Yes	Data Specialist	26
Phoebe	\$400k	No	Songwriter	26



Observations



Some ML vocabulary

Features

Outcome

Name	Pay	Married?	Job	Age
Rachel	\$32k	No	Server	24
Ross	\$140k	No	Professor	26
Joey	\$350k	No	Actor	25
Monica	\$70k	Yes	Chef	24
Chandler	\$120k	Yes	Data Specialist	26
Phoebe	\$400k	No	Songwriter	26

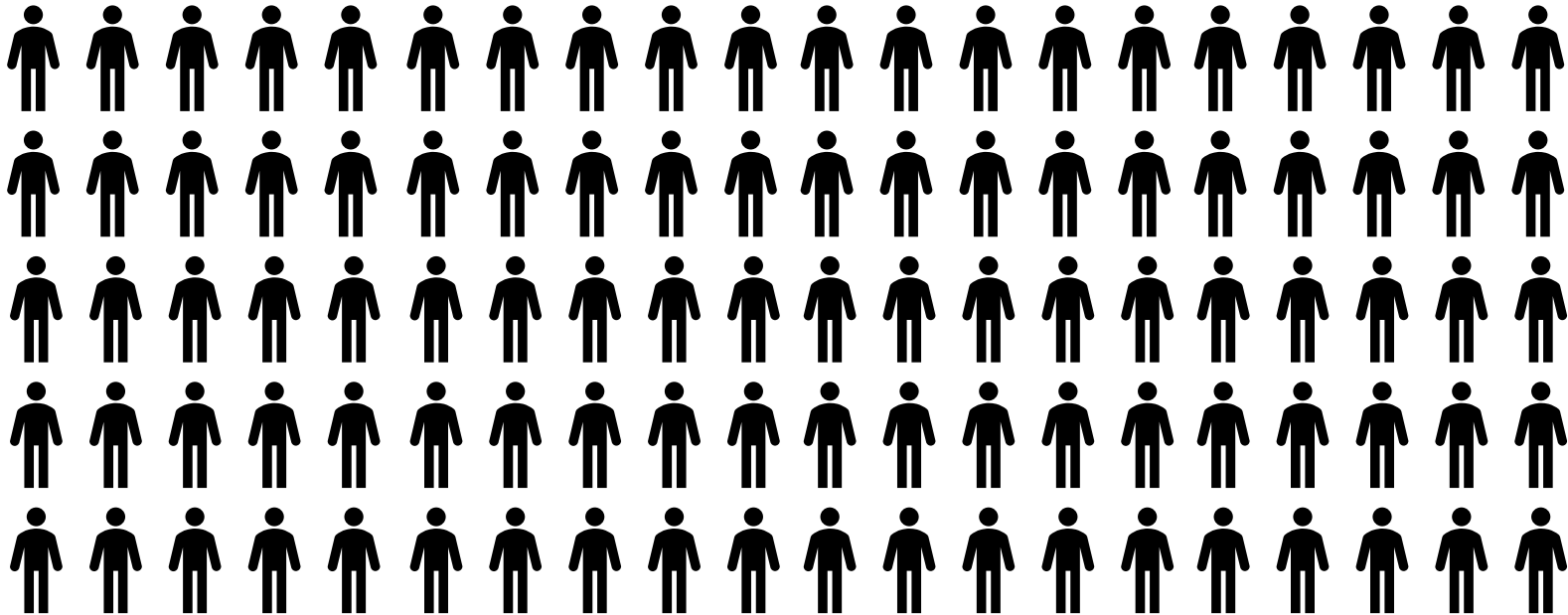
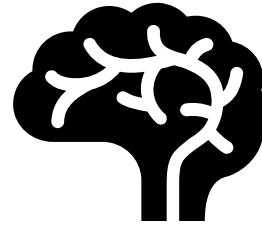
Some ML vocabulary

Features

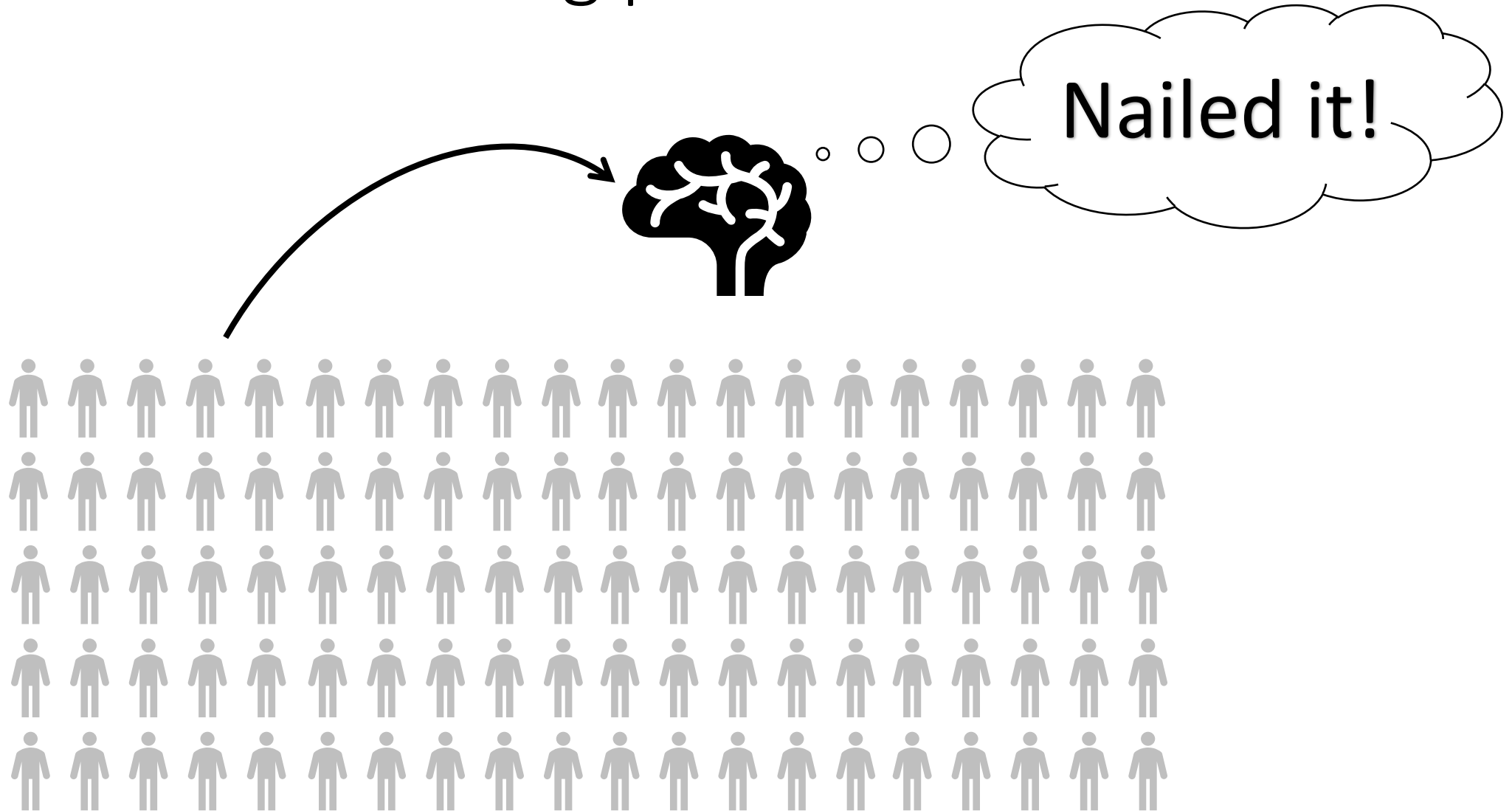
Outcome

Name	Pay	Married?	Job	Age
Rachel	\$32k	No	Server	24
Ross	\$140k	No	Professor	26
Joey	\$350k	No	Actor	25
Monica	\$70k	Yes	Chef	24
Chandler	\$120k	Yes	Data Specialist	26
Phoebe	\$400k	No	Songwriter	26

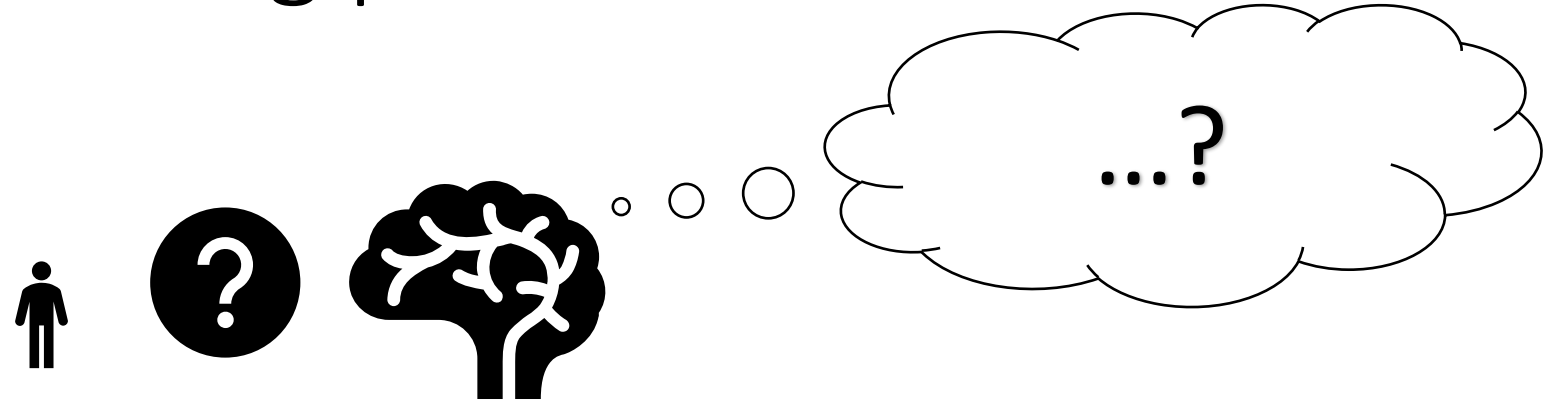
The practice: Assessing performance



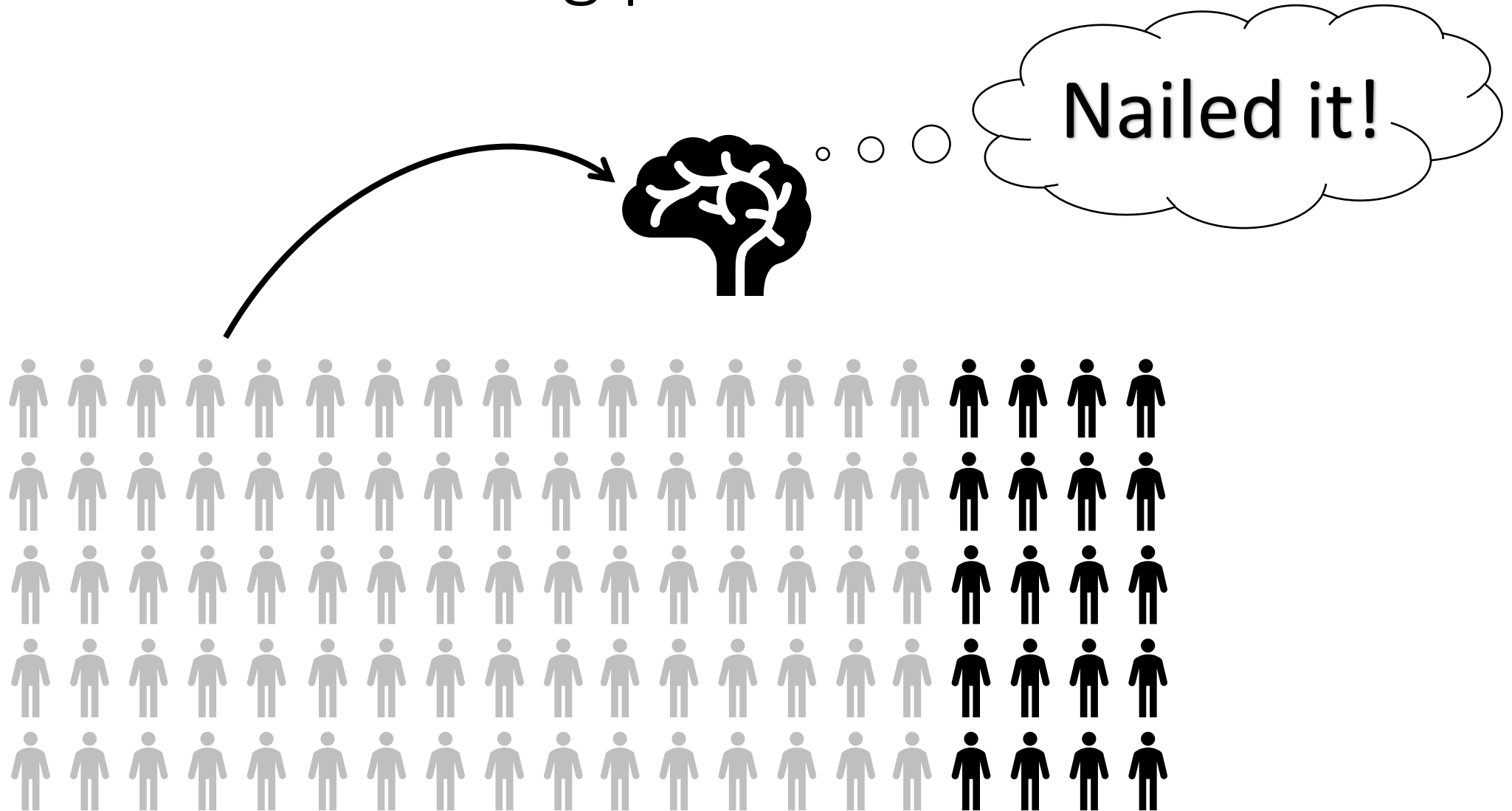
The practice: Assessing performance



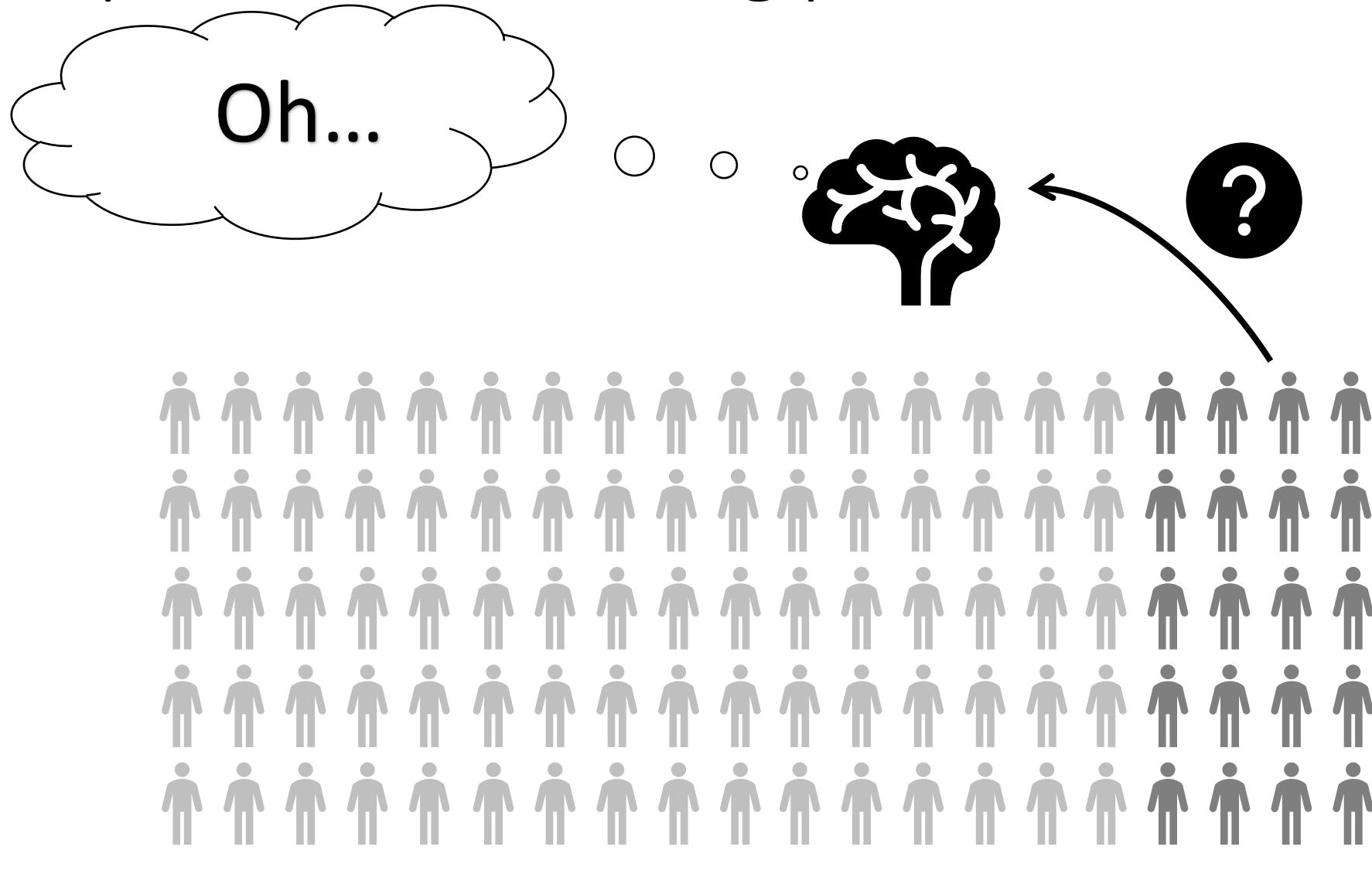
The practice: Assessing performance



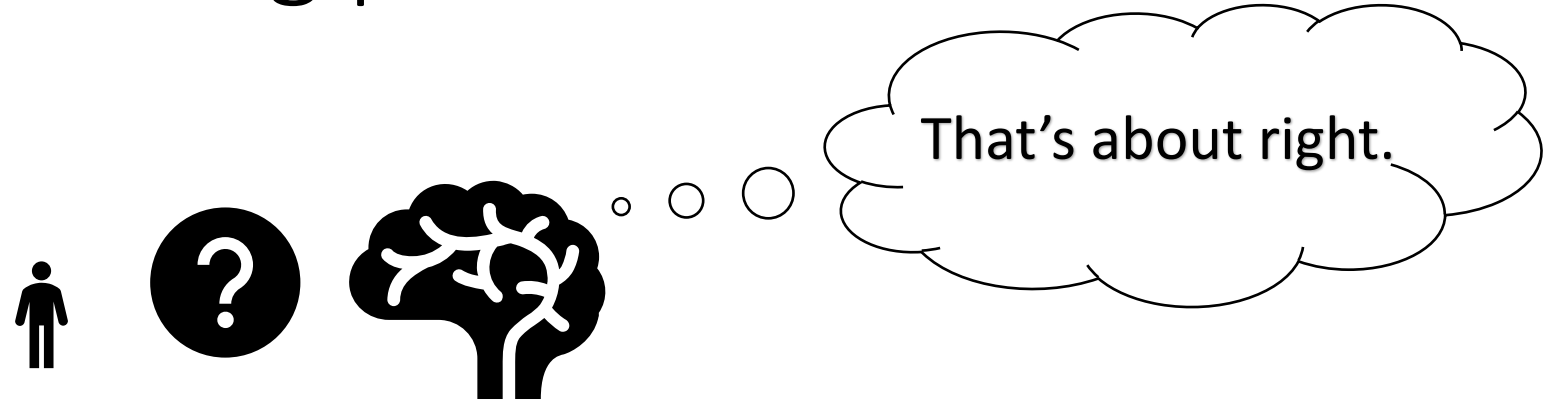
The practice: Assessing performance



The practice: Assessing performance



The practice: Assessing performance



The practice: Quantifying prediction

Binary outcomes

- Accuracy
- AUC scores

Continuous outcomes

- Mean absolute error
- Pearson correlation

The practice: Quantifying prediction

Binary outcomes

- Accuracy
- AUC scores

Continuous outcomes

- Mean absolute error and mean squared error
- Pearson correlation

Actually Positive

Actually Negative

Predicted Positive

Predicted Negative

Accuracy

	<u>Actually Positive</u>	<u>Actually Negative</u>
<u>Predicted Positive</u>	True positive	False positive
<u>Predicted Negative</u>	False negative	True negative

Accuracy

	<u>Actually Positive</u>	<u>Actually Negative</u>
<u>Predicted Positive</u>	True positive	False positive
<u>Predicted Negative</u>	False negative	True negative

Accuracy

$$Accuracy = \frac{\text{True positive} + \text{True negative}}{\text{Number of observations}}$$

	<u>Actually Positive</u>	<u>Actually Negative</u>
<u>Predicted Positive</u>	200	50
<u>Predicted Negative</u>	50	200

Accuracy

Accuracy = ?

	<u>Actually Positive</u>	<u>Actually Negative</u>
<u>Predicted Positive</u>	200	50
<u>Predicted Negative</u>	50	200

Accuracy

$$Accuracy = \frac{TP + TN}{N} = \frac{200 + 200}{200 + 50 + 50 + 200} = \frac{400}{500} = \mathbf{0.8}$$

	<u>Actually Positive</u>	<u>Actually Negative</u>
<u>Predicted Positive</u>	0	0
<u>Predicted Negative</u>	100	400

Accuracy

Accuracy = ?

	<u>Actually Positive</u>	<u>Actually Negative</u>
<u>Predicted Positive</u>	0	0
<u>Predicted Negative</u>	100	400

Accuracy

$$Accuracy = \frac{TP + TN}{N} = \frac{400 + 0}{400 + 100} = \frac{400}{500} = 0.8$$

When Accuracy seems inaccurate

- If one outcome is very common, accuracy is easy to achieve
- This can be manually examined or tested
- Other metrics (F1 score, precision or recall, chi-square) can be more useful in these situations
- Alternatively, you can compare to a “baseline model”

The practice: Quantifying prediction

Binary outcomes

- Accuracy
- AUC

Continuous outcomes

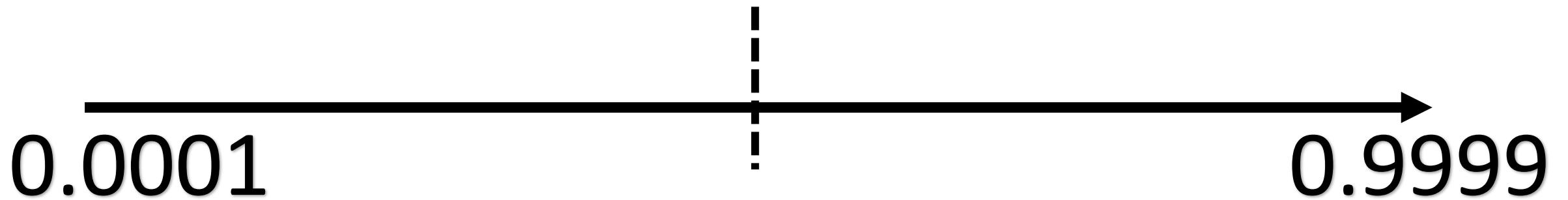
- Mean absolute error
- Pearson correlation

Motivating AUC

1. Predicting who will fail the class (and ask them to come to office hours)
2. Predicting who will get COVID (and asking them to stay home from class)

Motivating AUC

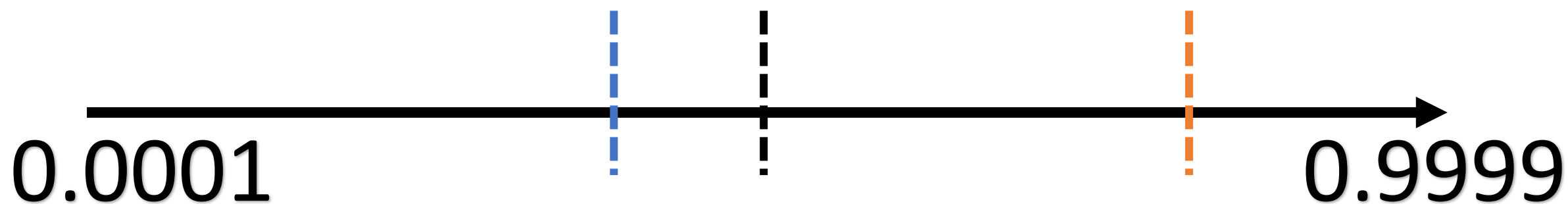
1. Predicting who will fail the class (and ask them to come to office hours)
2. Predicting who will get COVID (and asking them to stay home from class)



How sure are we that $Y = 1$?

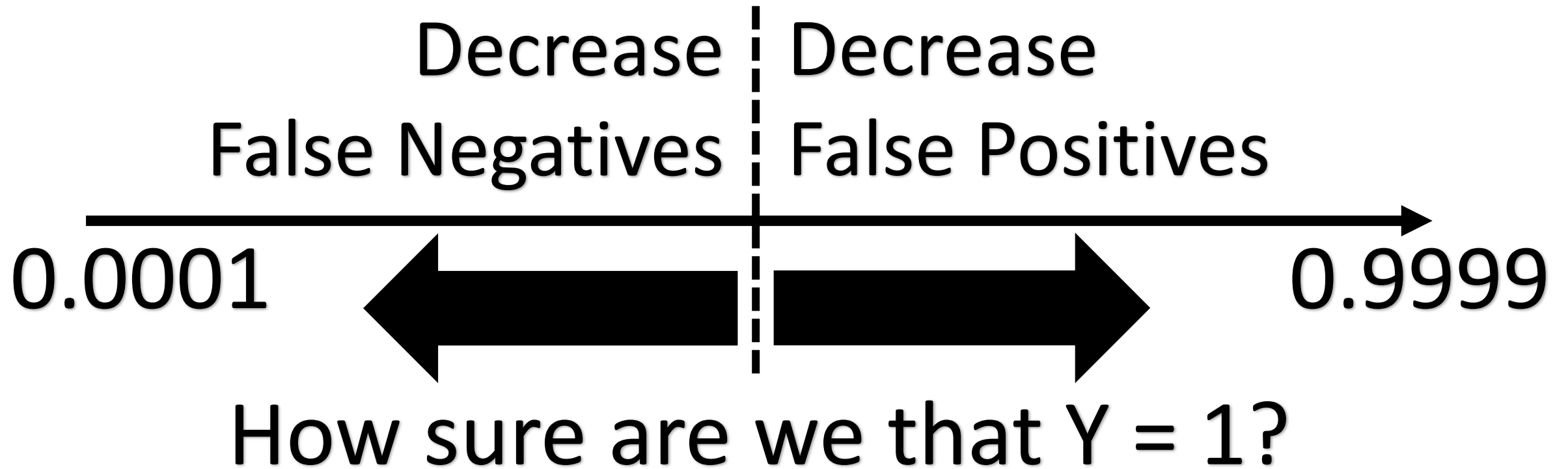
Motivating AUC

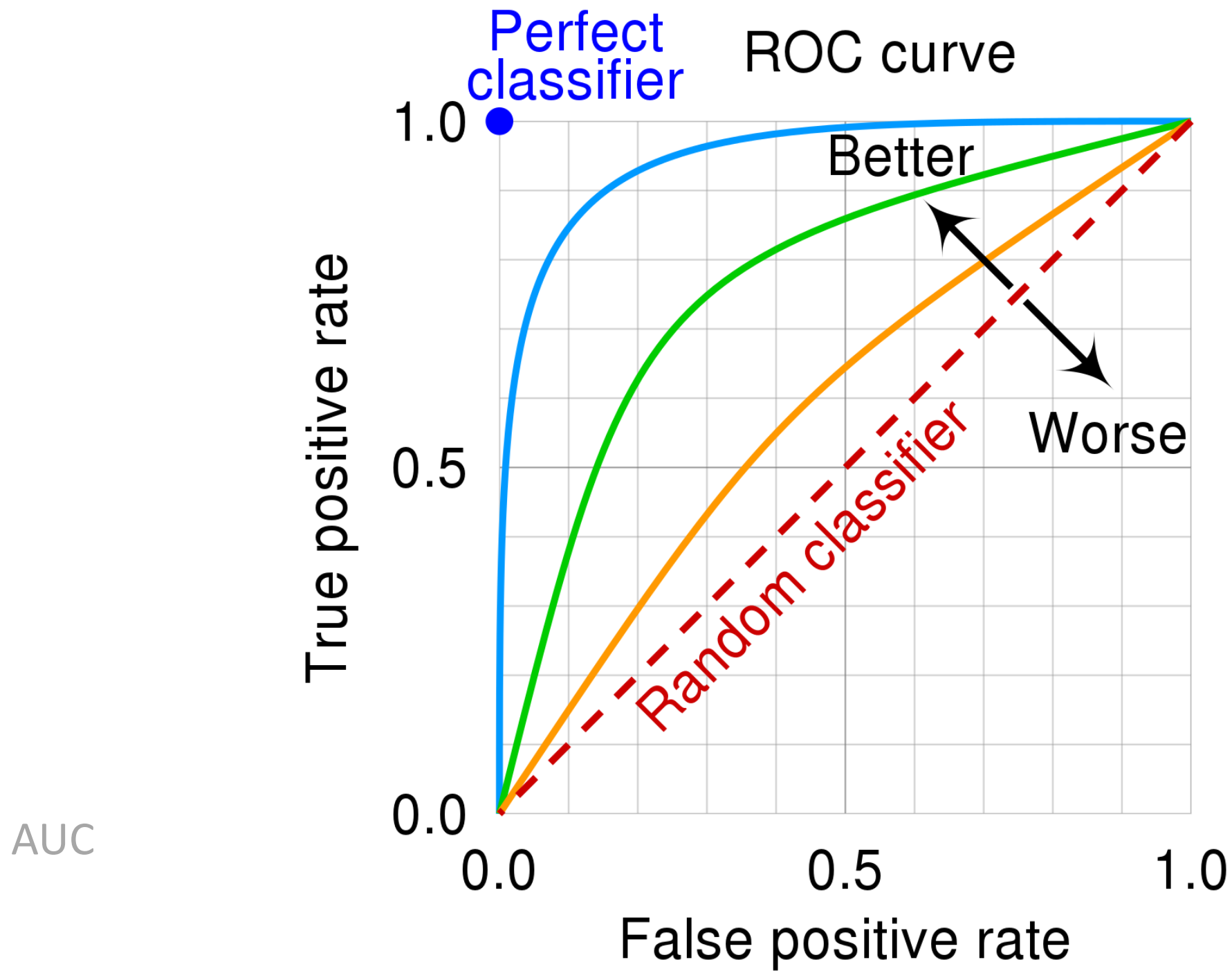
1. Predicting who will fail the class (and ask them to come to office hours)
2. Predicting who will get COVID (and asking them to stay home from class)



How sure are we that $Y = 1$?

Motivating AUC





The practice: Quantifying prediction

Binary outcomes

- Accuracy
- AUC

Continuous outcomes

- Mean absolute error
- Pearson correlation

Mean absolute error

Pay	Age	Tenure (yrs)	Bachelors'?	Predicted Pay
\$60k	35	2	Yes	\$58k
\$80k	45	27	No	\$100k
\$40k	19	1	No	\$30k
\$100k	53	20	No	\$100k
\$180k	52	3	Yes	\$200k

Mean absolute error

Pay	Age	Tenure (yrs)	Bachelors'?	Predicted Pay	Error
\$60k	35	2	Yes	\$80k	\$20k
\$80k	45	27	No	\$50k	-\$30k
\$40k	19	1	No	\$30k	-\$10k
\$100k	53	20	No	\$100k	\$0
\$180k	52	3	Yes	\$200k	\$20k

Mean absolute error

Average error = 0?

Pay	Age	Tenure (yrs)	Bachelors'?	Predicted Pay	Error
\$60k	35	2	Yes	\$80k	\$20k
\$80k	45	27	No	\$50k	-\$30k
\$40k	19	1	No	\$30k	-\$10k
\$100k	53	20	No	\$100k	\$0
\$180k	52	3	Yes	\$200k	\$20k

Mean absolute error

Pay	Age	Tenure (yrs)	Bachelors'?	Predicted Pay	Error	Abs(Error)
\$60k	35	2	Yes	\$80k	\$20k	\$20k
\$80k	45	27	No	\$50k	-\$30k	\$30k
\$40k	19	1	No	\$30k	-\$10k	\$10k
\$100k	53	20	No	\$100k	\$0	0
\$180k	52	3	Yes	\$200k	\$20k	\$20k

Mean absolute error

Average error = \$16k

Pay	Age	Tenure (yrs)	Bachelors'?	Predicted Pay	Error	Abs(Error)
\$60k	35	2	Yes	\$80k	\$20k	\$20k
\$80k	45	27	No	\$50k	-\$30k	\$30k
\$40k	19	1	No	\$30k	-\$10k	\$10k
\$100k	53	20	No	\$100k	\$0	0
\$180k	52	3	Yes	\$200k	\$20k	\$20k



$\text{sum}(\text{prediction} - \text{actual}) / N$



$\text{sum}(\text{abs}(\text{prediction} - \text{actual})) / N$

The practice: Quantifying prediction

Binary outcomes

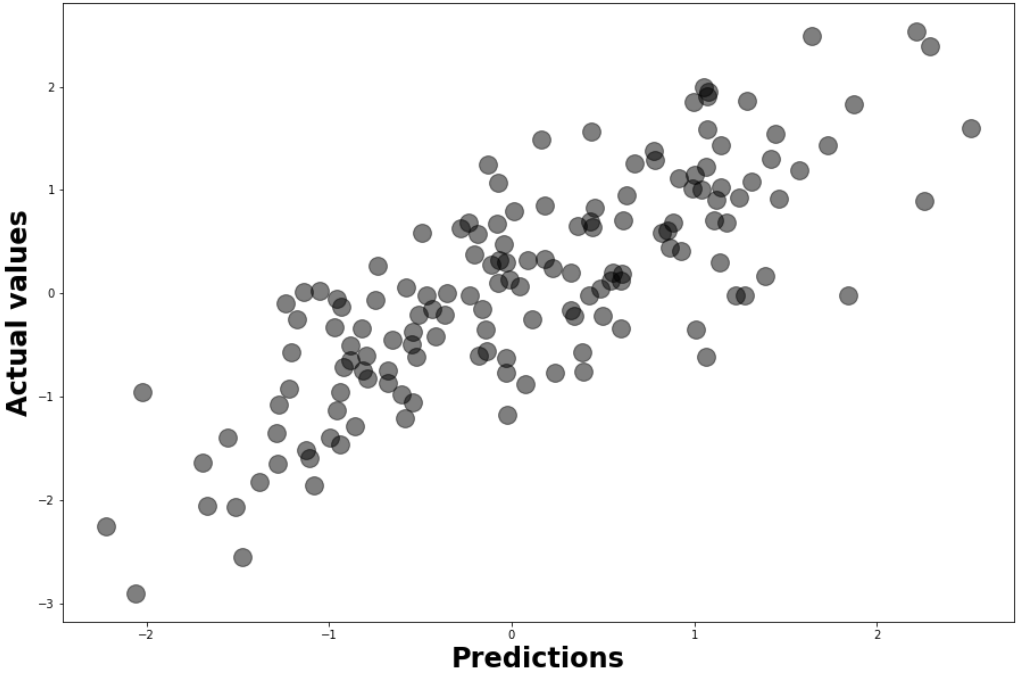
- Accuracy
- AUC

Continuous outcomes

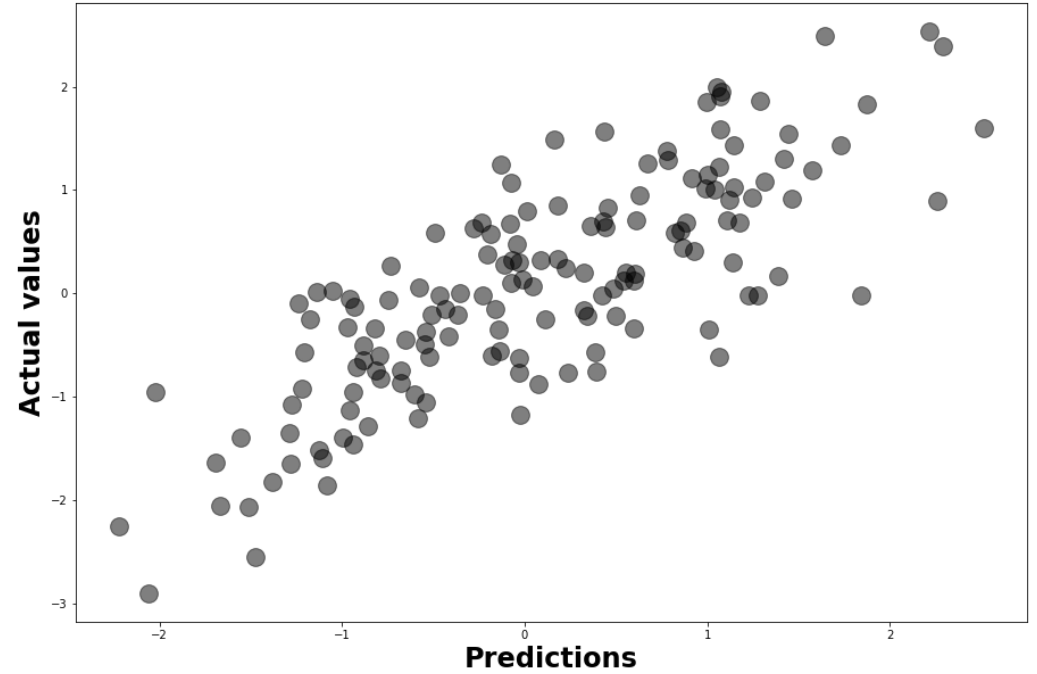
- Mean absolute error
- **Pearson correlation**

Evaluation	Lines of code	Days late	Pred. evaluation
-1	80	8	-0.8
1.3	105	3	1.5
0.1	95	4	0
-2	30	12	-2.1
...	
0.5	100	1	0.7

Evaluation	Lines of code	Days late	Pred. evaluation
-1	80	8	-0.8
1.3	105	3	1.5
0.1	95	4	0
-2	30	12	-2.1
...	
0.5	100	1	0.7

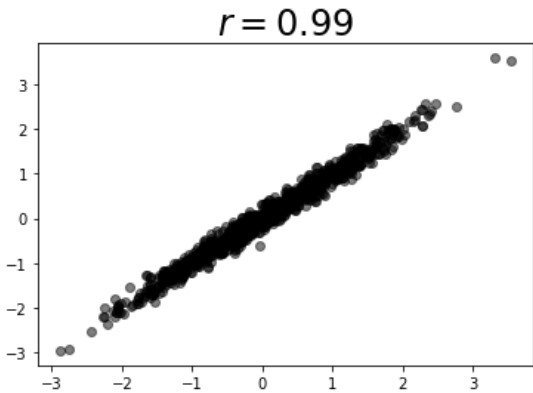
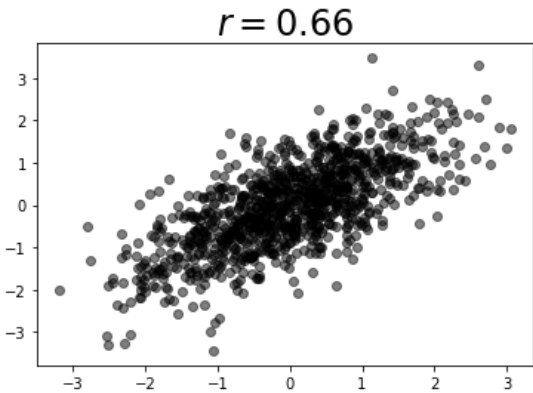
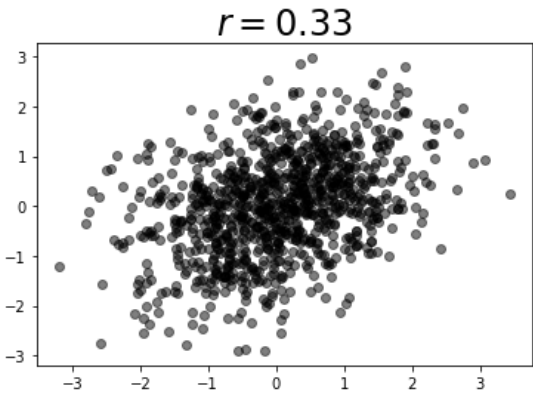
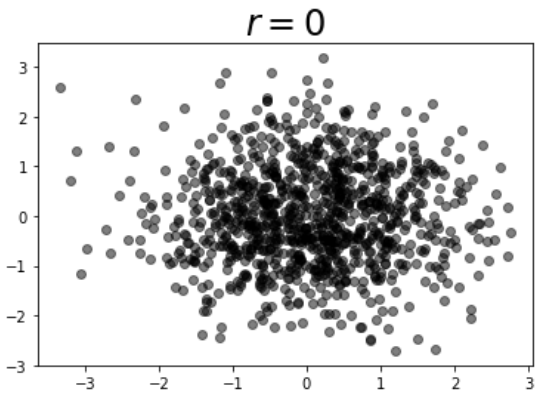
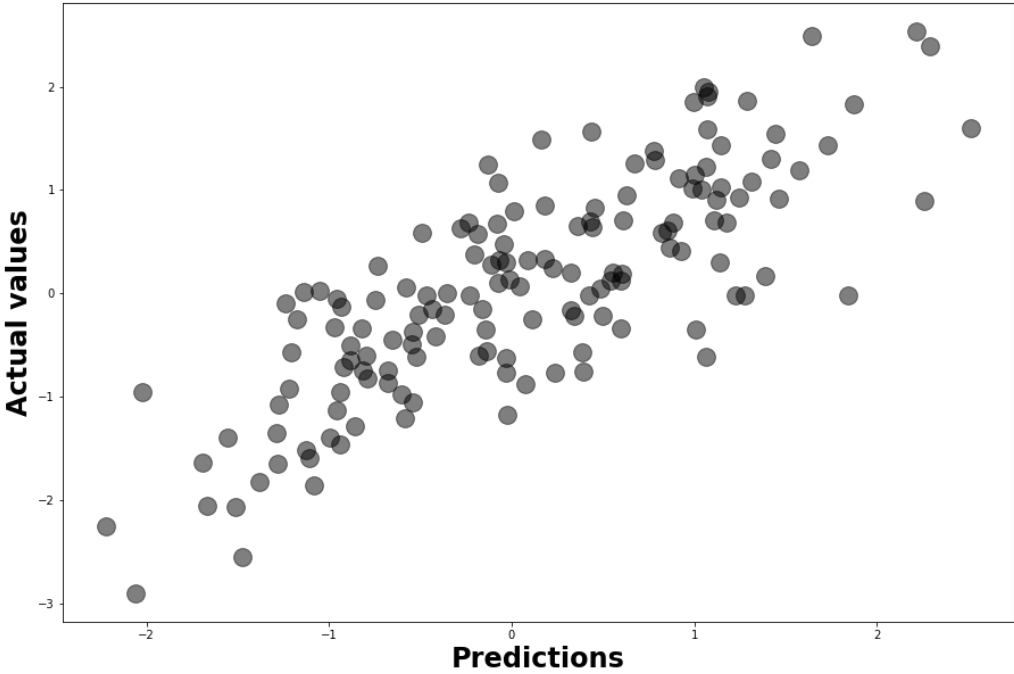


Evaluation	Lines of code	Days late	Pred. evaluation
-1	80	8	-0.8
1.3	105	3	1.5
0.1	95	4	0
-2	30	12	-2.1
...	
0.5	100	1	0.7

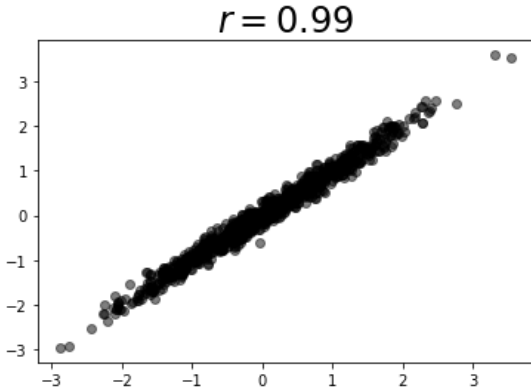
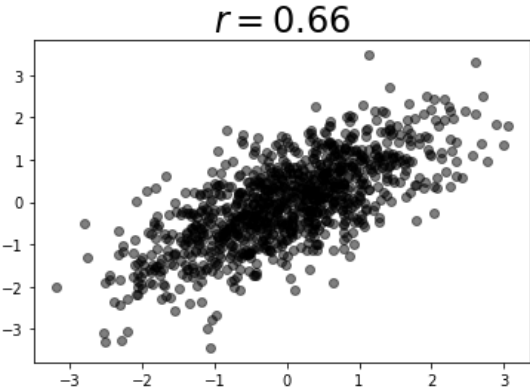
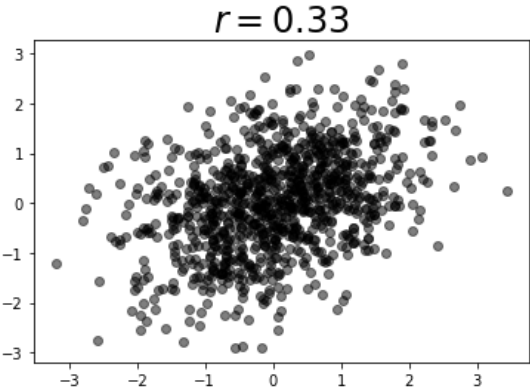
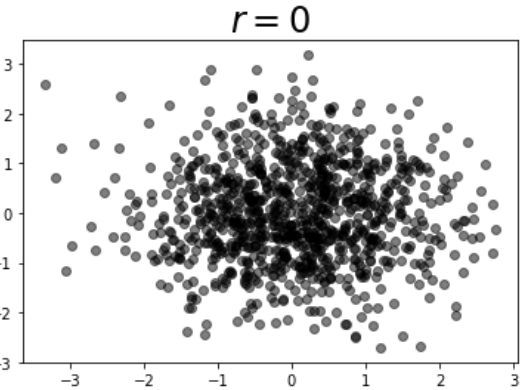
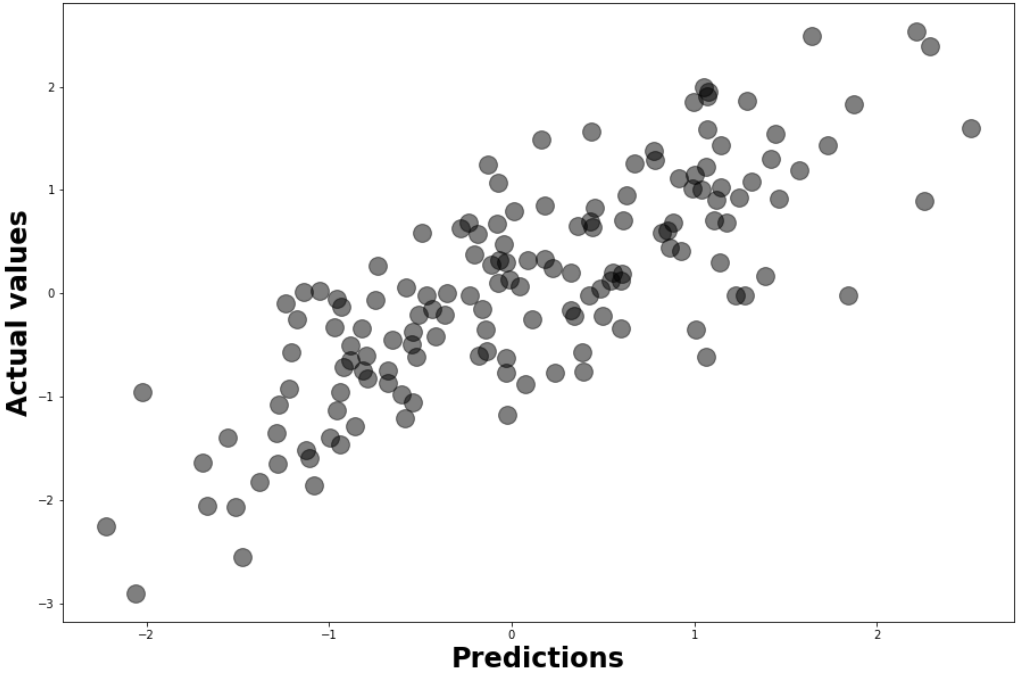


$$r = \frac{\sum_{i=1}^N ([x_i - \bar{x}] * [y_i - \bar{y}])}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2 * \sum_{i=1}^N (y_i - \bar{y})^2}}$$

Evaluation	Lines of code	Days late	Pred. evaluation
-1	80	8	-0.8
1.3	105	3	1.5
0.1	95	4	0
-2	30	12	-2.1
...	
0.5	100	1	0.7



Evaluation	Lines of code	Days late	Pred. evaluation
-1	80	8	-0.8
1.3	105	3	1.5
0.1	95	4	0
-2	30	12	-2.1
...	
0.5	100	1	0.7



Meet the Algorithms



LINEAR REGRESSION

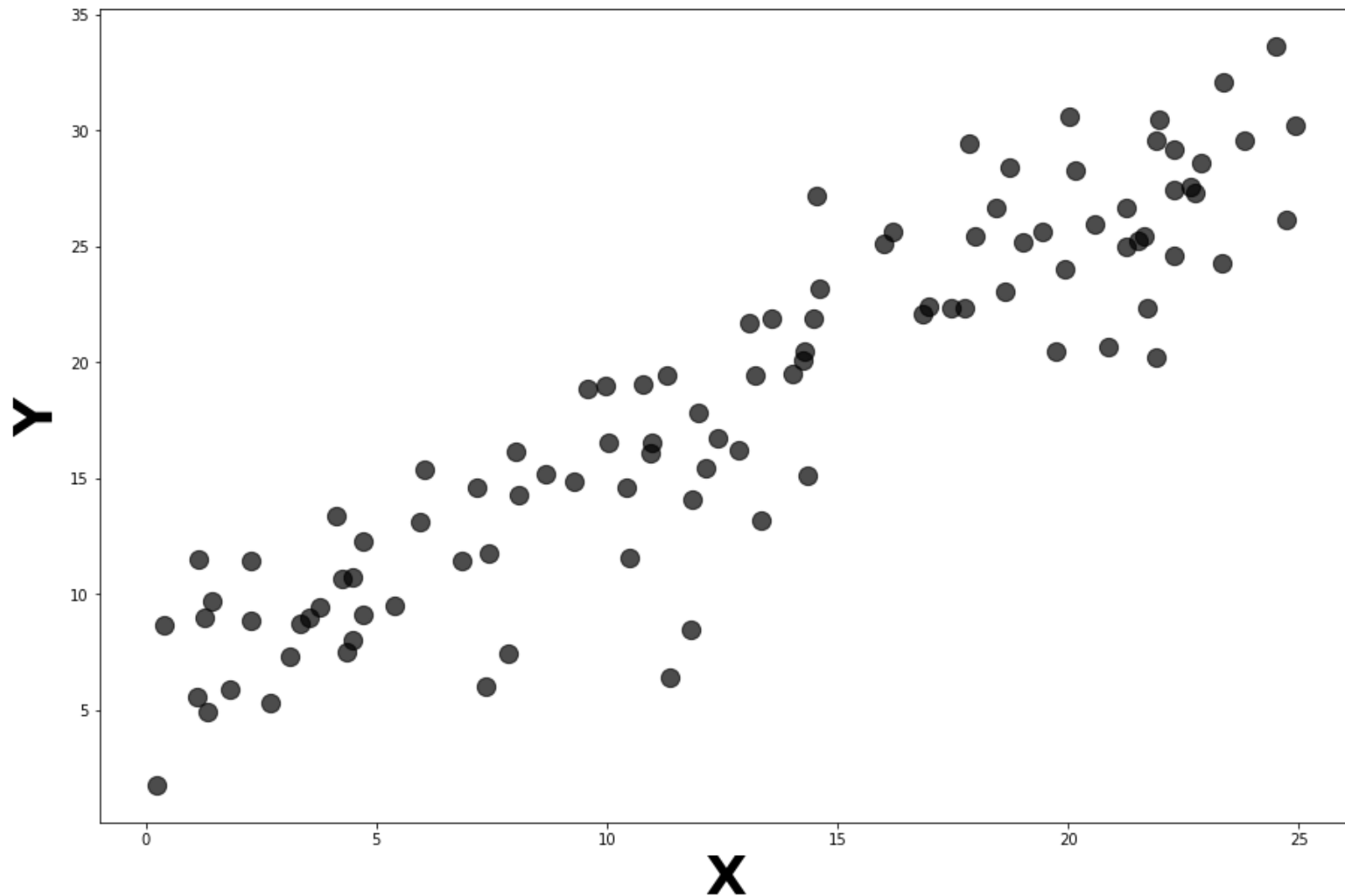


TREE-BASED METHODS

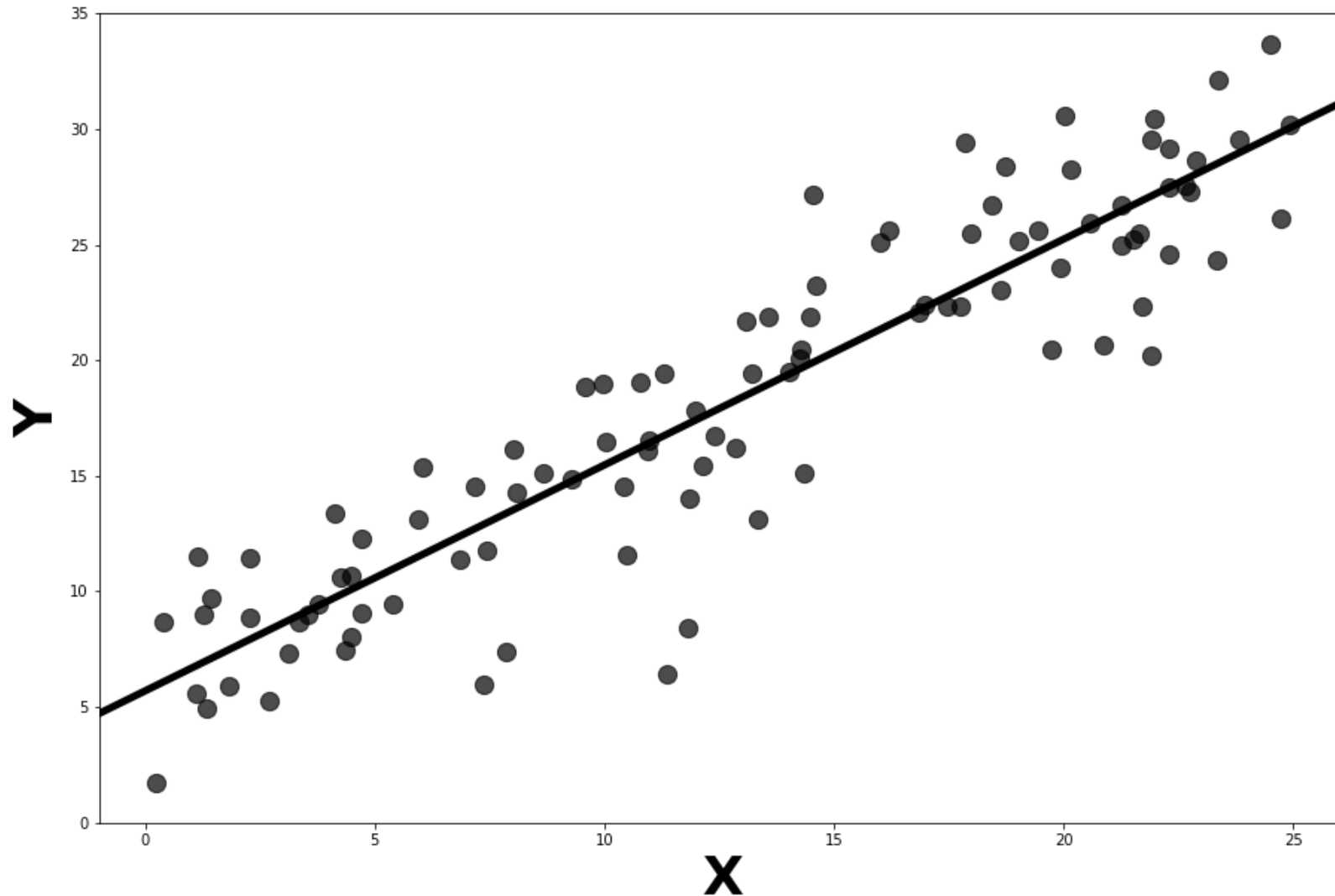


DEEP LEARNING

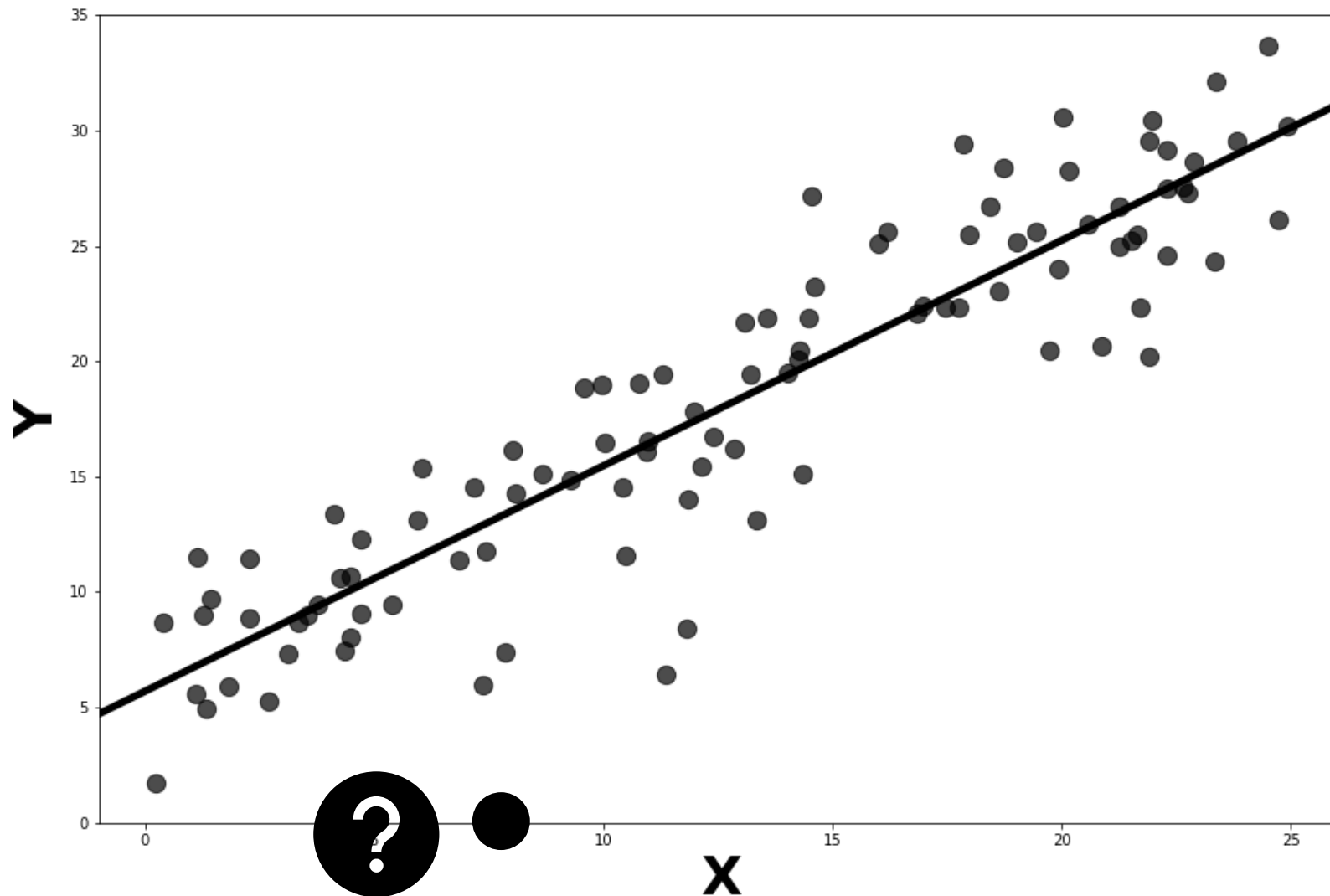
Meet the algorithms: Linear regression



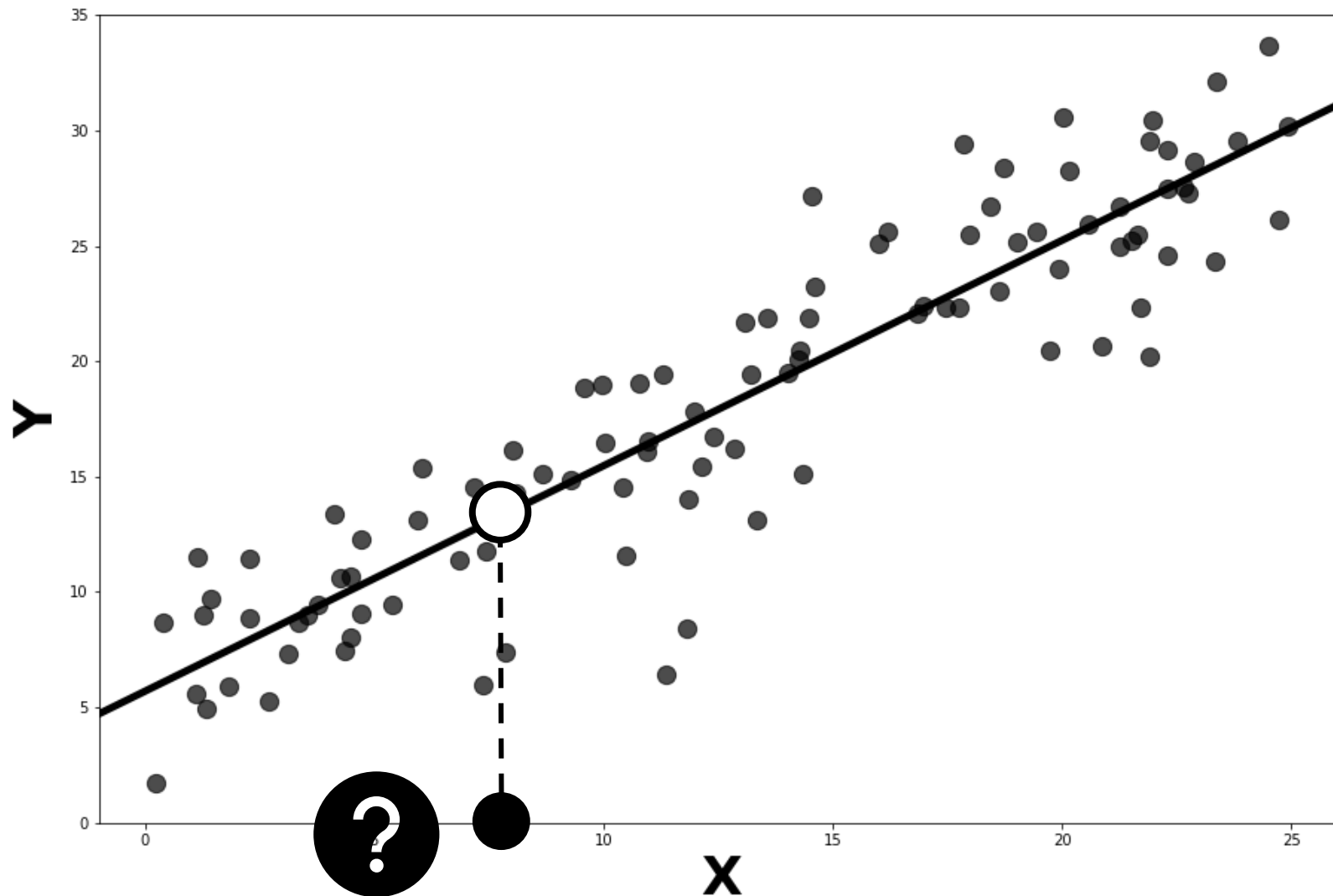
Meet the algorithms: Linear regression



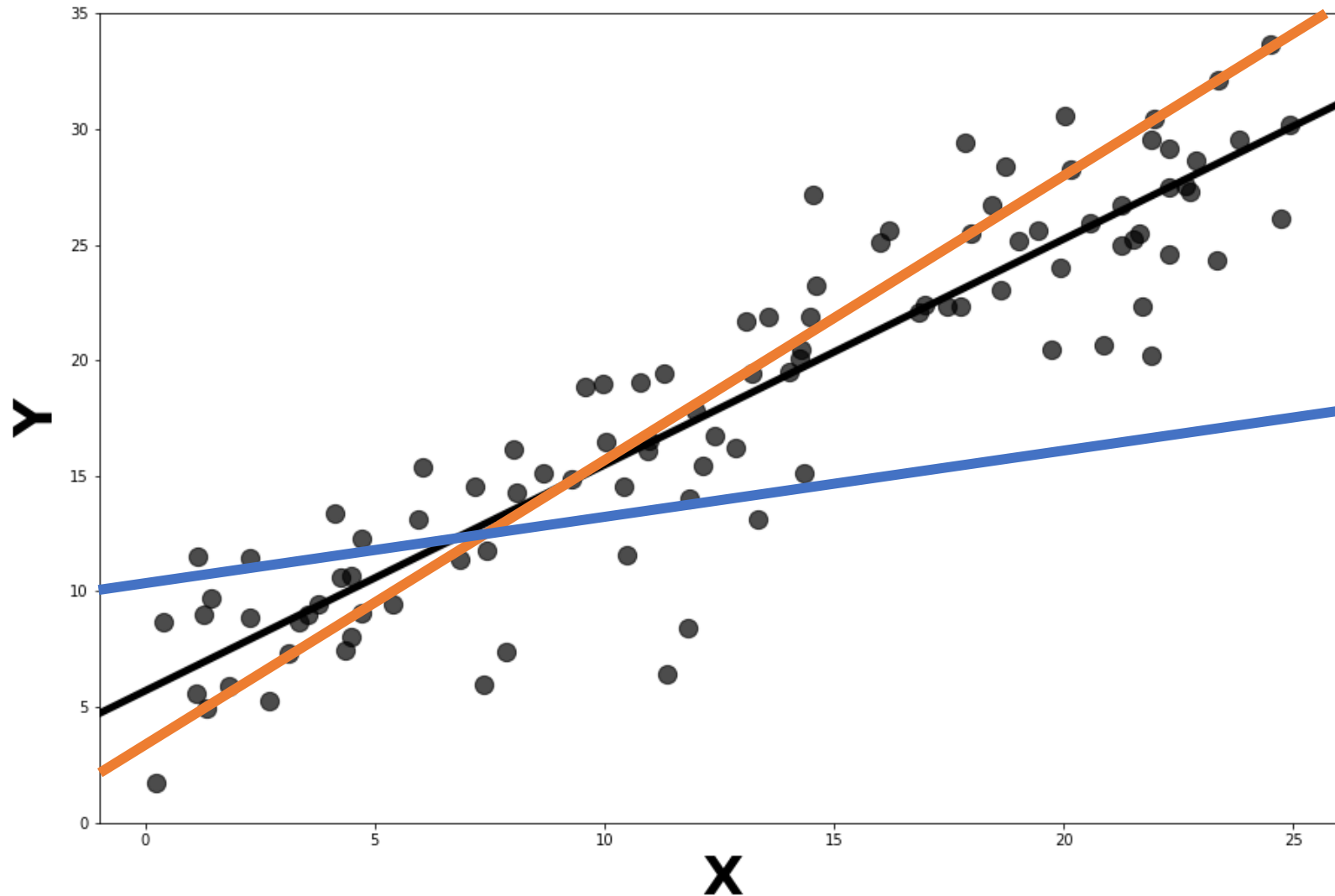
Meet the algorithms: Linear regression



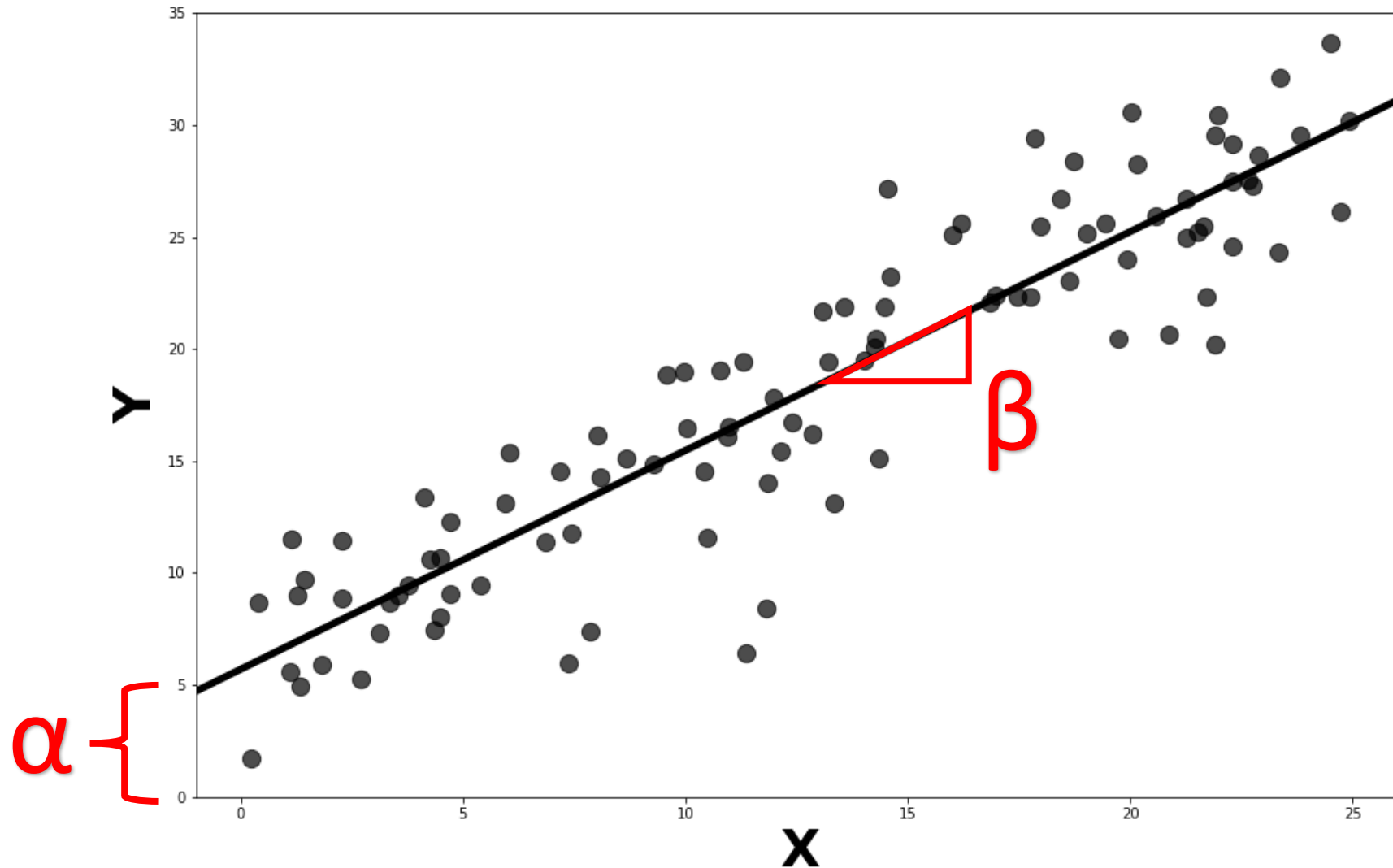
Meet the algorithms: Linear regression



Meet the algorithms: Linear regression



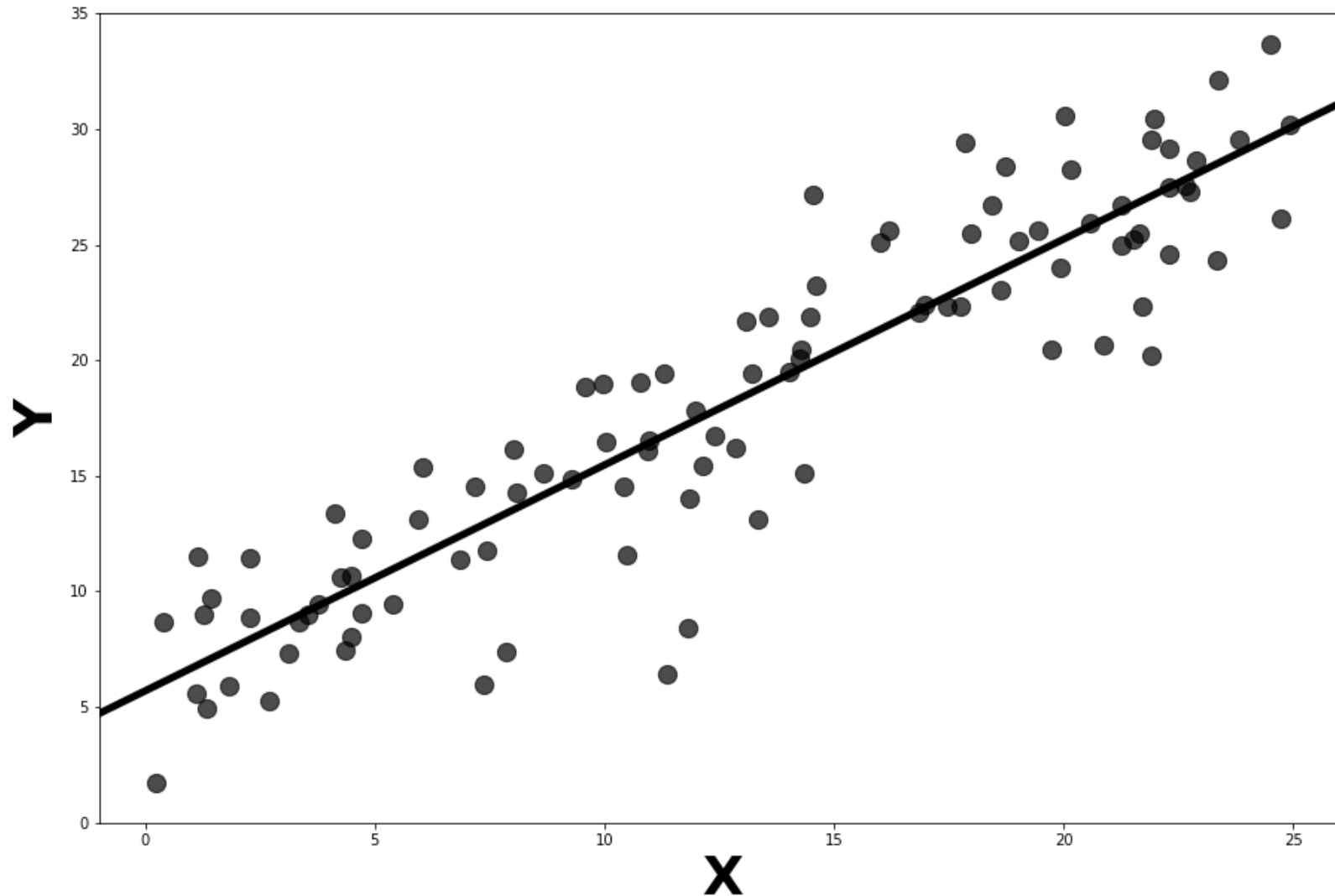
Meet the algorithms: Linear regression



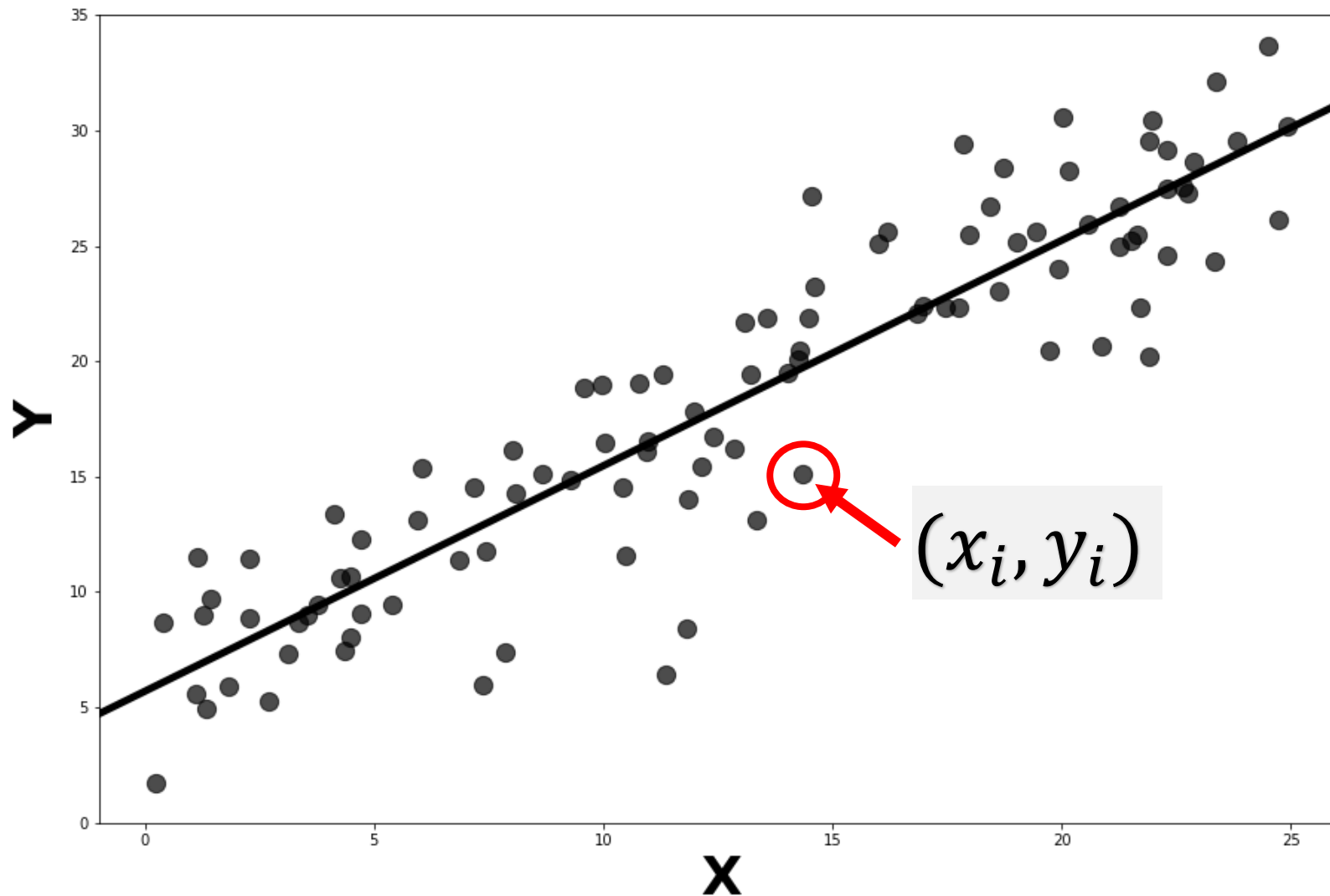
Meet the algorithms: Linear regression

$$SS_{\alpha, \beta} = \sum_{i=1}^N (y_i - [\alpha + \beta x_i])^2$$

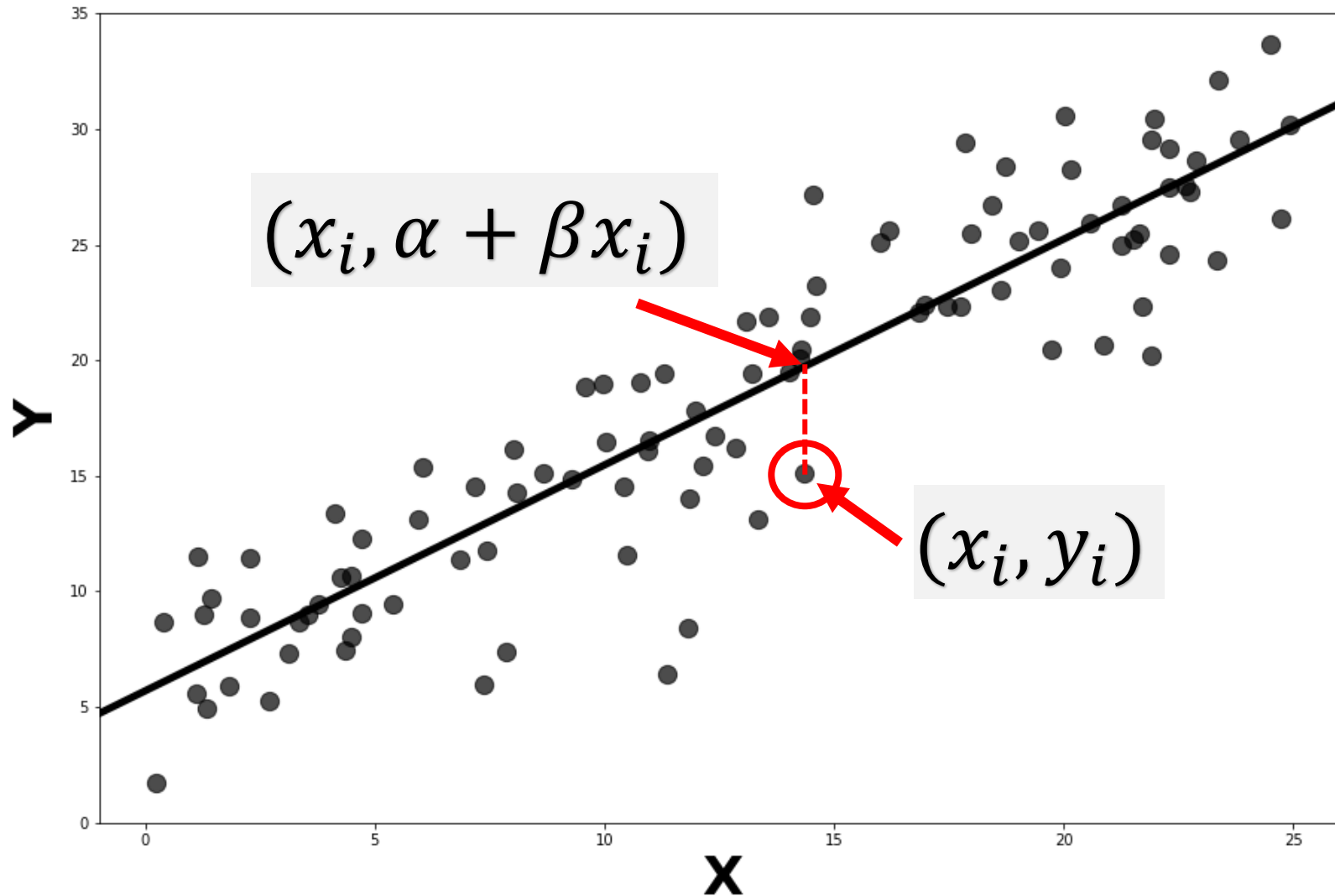
Meet the algorithms: Linear regression



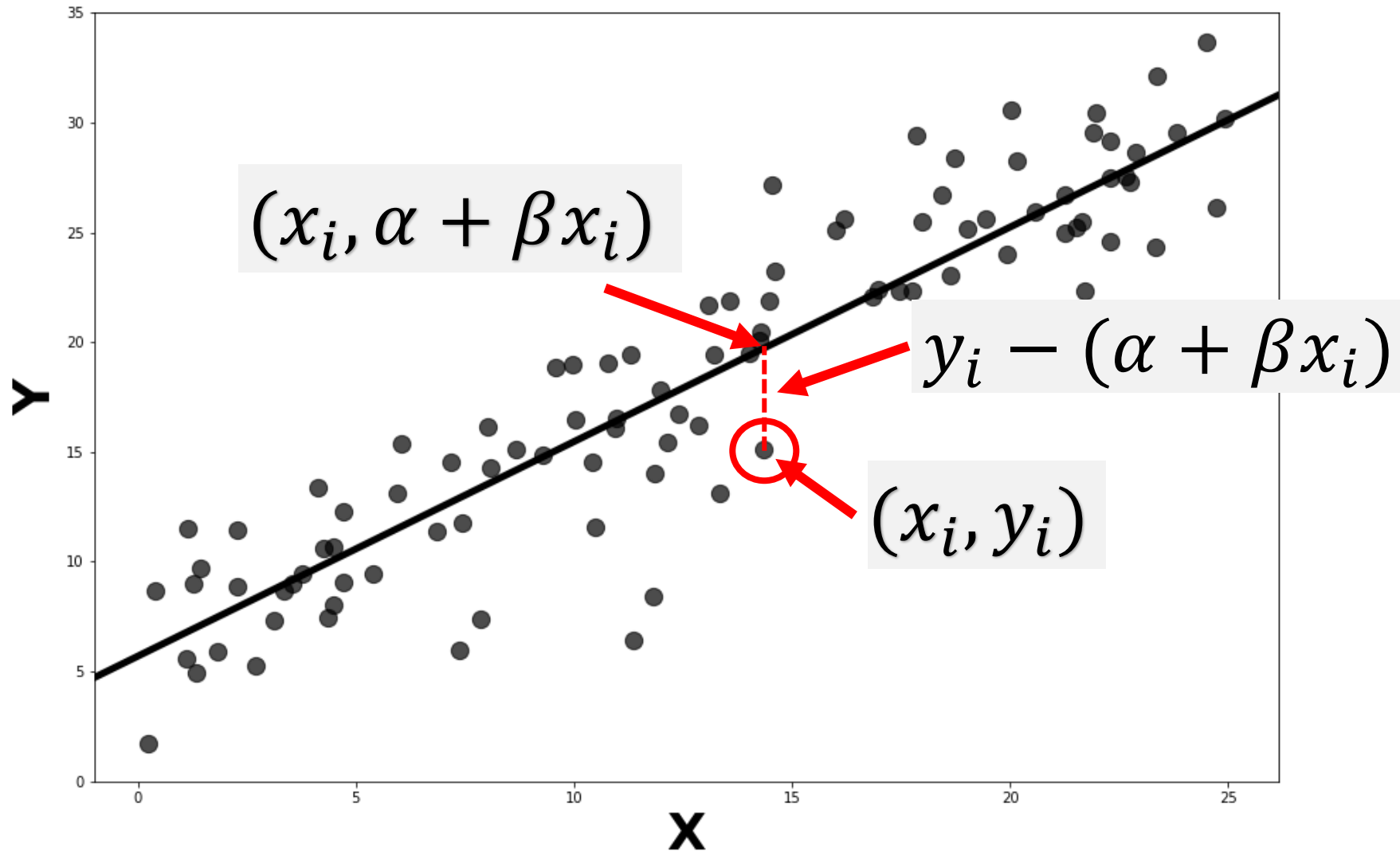
Meet the algorithms: Linear regression



Meet the algorithms: Linear regression



Meet the algorithms: Linear regression



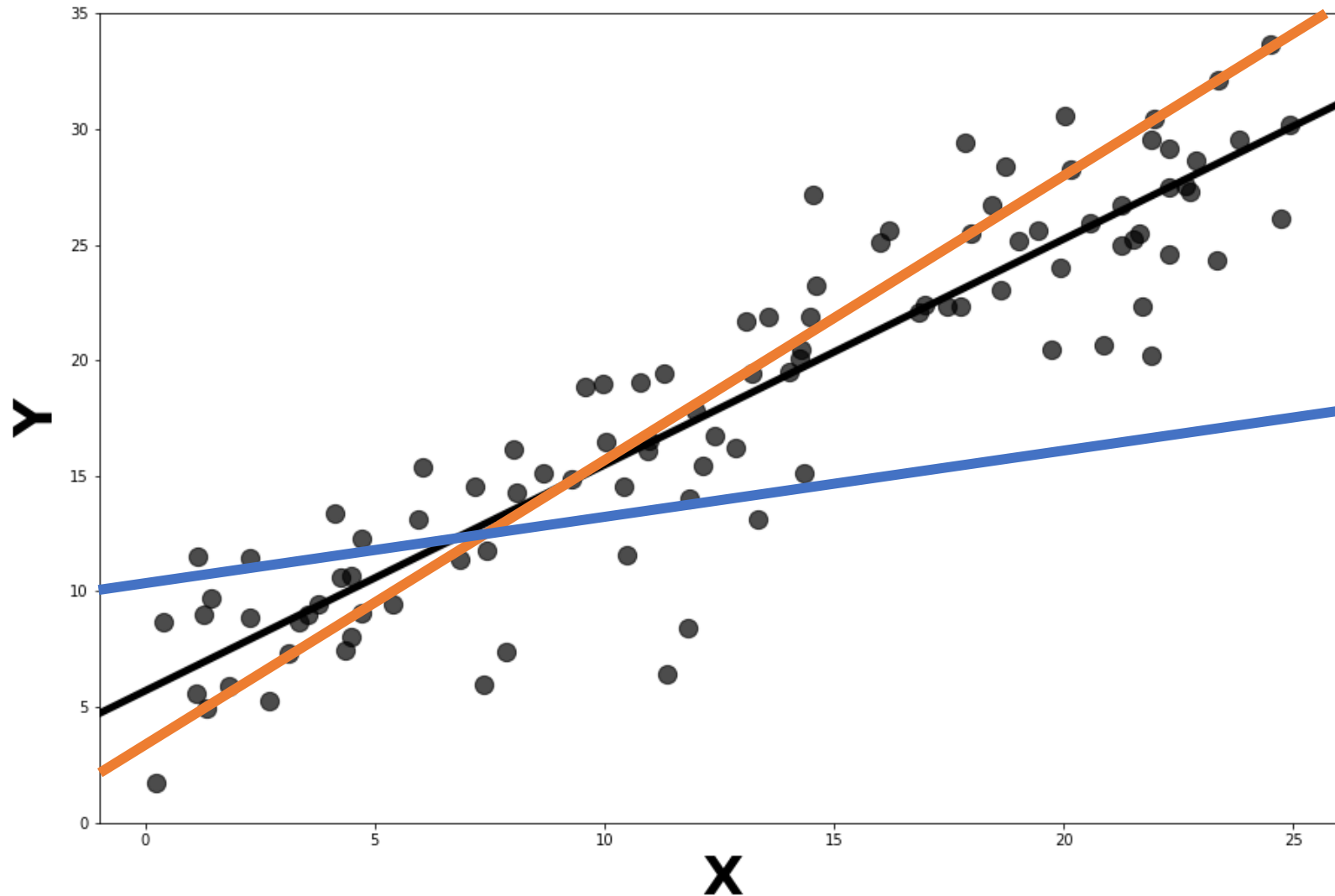
Meet the algorithms: Linear regression

$$SS_{\alpha, \beta} = \sum_{i=1}^N (y_i - [\alpha + \beta x_i])^2$$

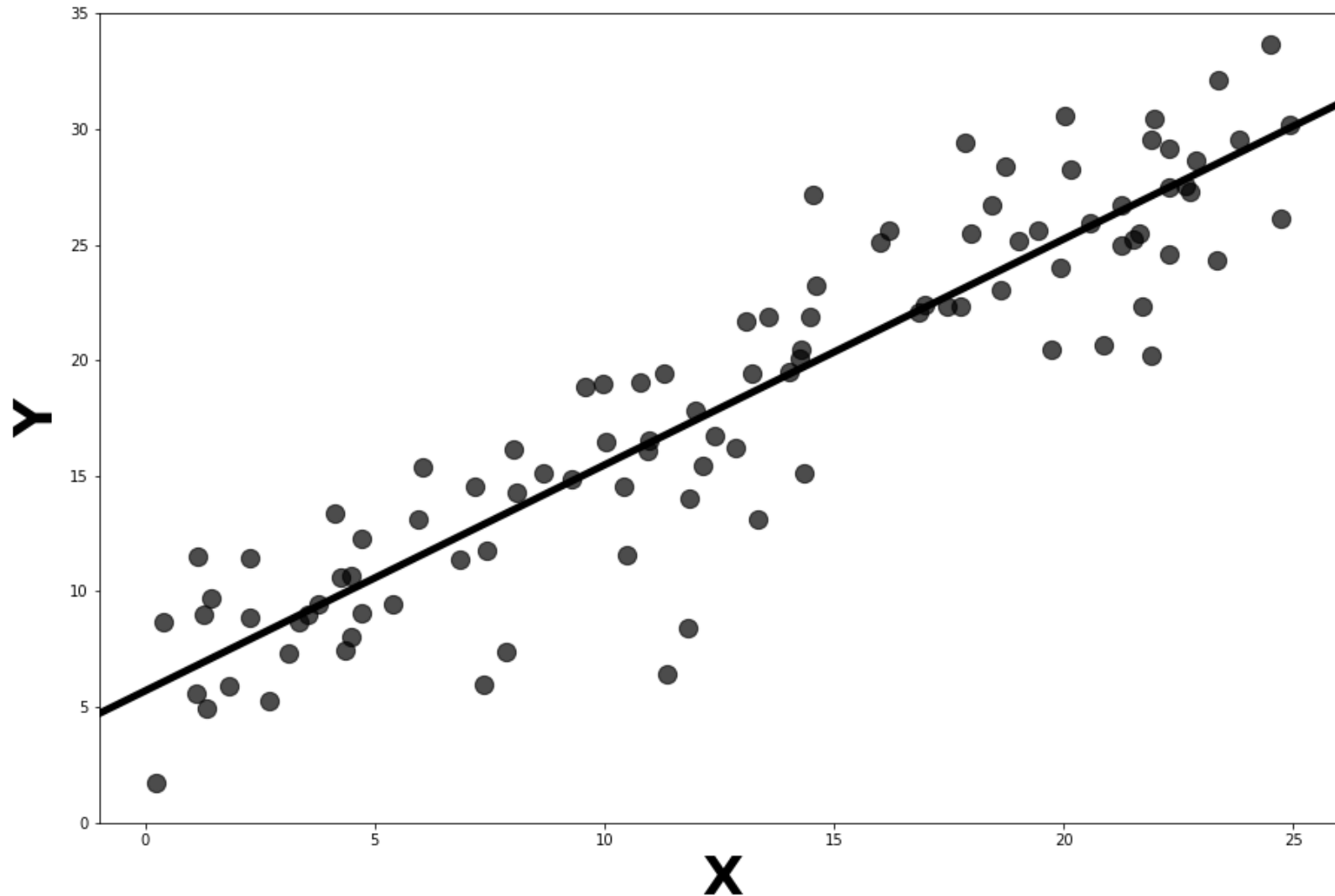
Diagram illustrating the components of the linear regression cost function $SS_{\alpha, \beta}$:

- Actual y** : Points to y_i .
- Our guess of y based on x** : Points to the expression $[\alpha + \beta x_i]$.
- How off the line is**: Points to the difference $y_i - [\alpha + \beta x_i]$.
- The sum of...**: Points to the summation index $i=1$ to N .

Meet the algorithms: Linear regression



Meet the algorithms: Linear regression



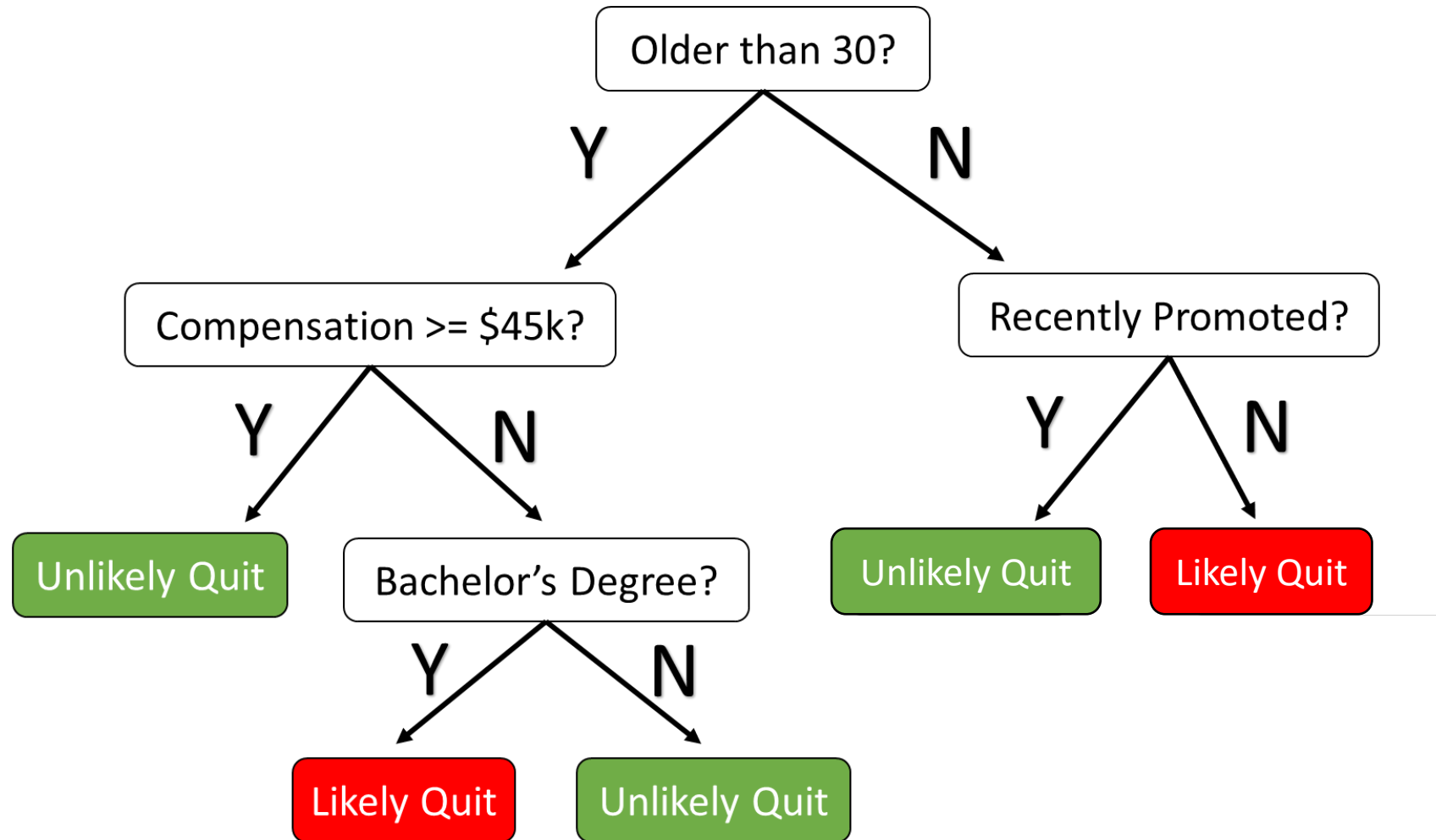
Going further: Multiple regression

$$SS_{\alpha, \beta} = \sum_{i=1}^N \left(y_i - \left[\alpha + \sum_{k=1}^K \beta_k x_{ik} \right] \right)^2$$

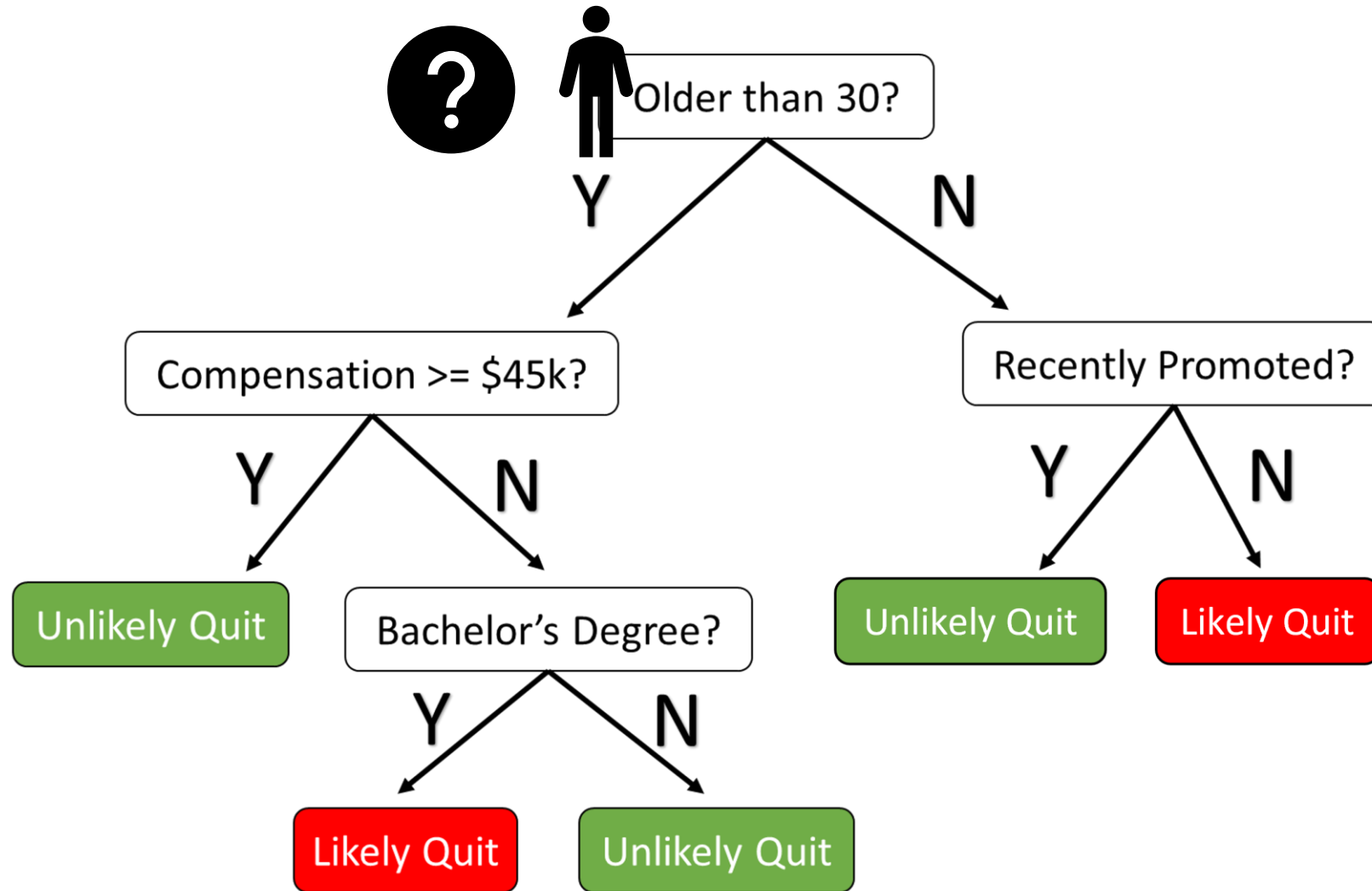
Going further: Penalized regression (LASSO)

$$SS_{\alpha,\beta} = \sum_{i=1}^N \left(y_i - \left[\alpha + \sum_{k=1}^K \beta_k x_{ik} \right] \right)^2 + \lambda \sum_{k=1}^K |\beta_k|$$

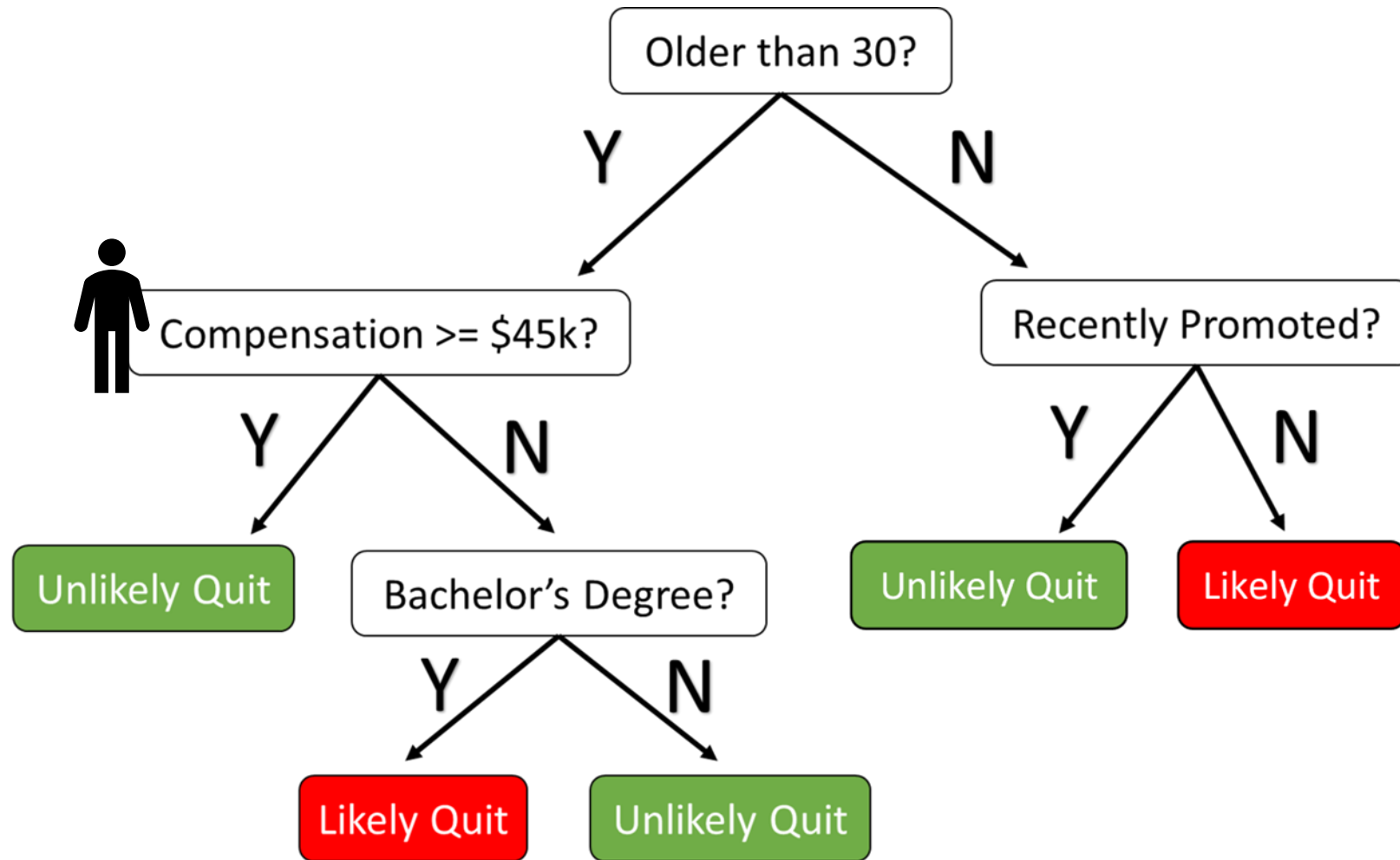
Meet the algorithms: Tree-based learning



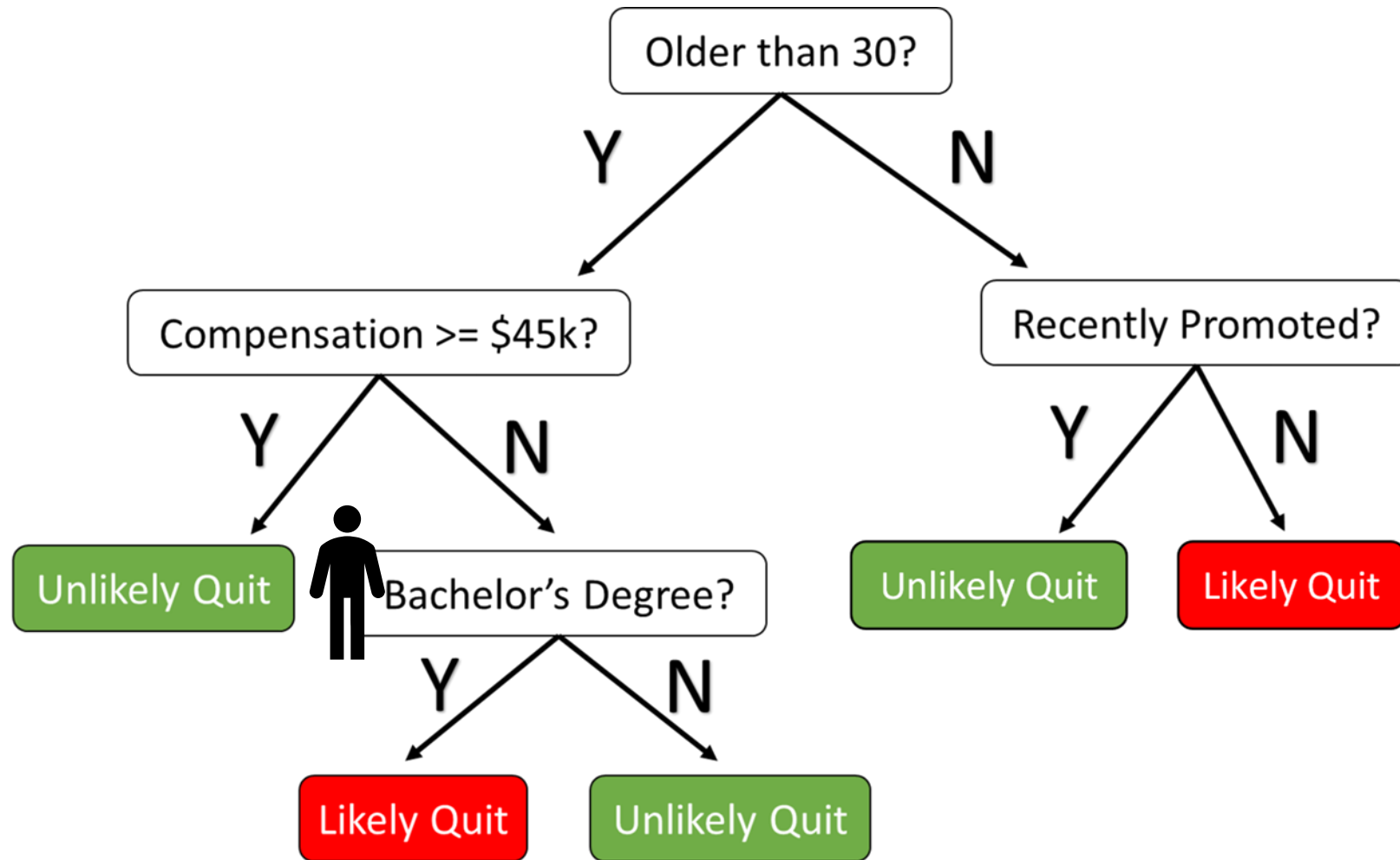
Meet the algorithms: Tree-based learning



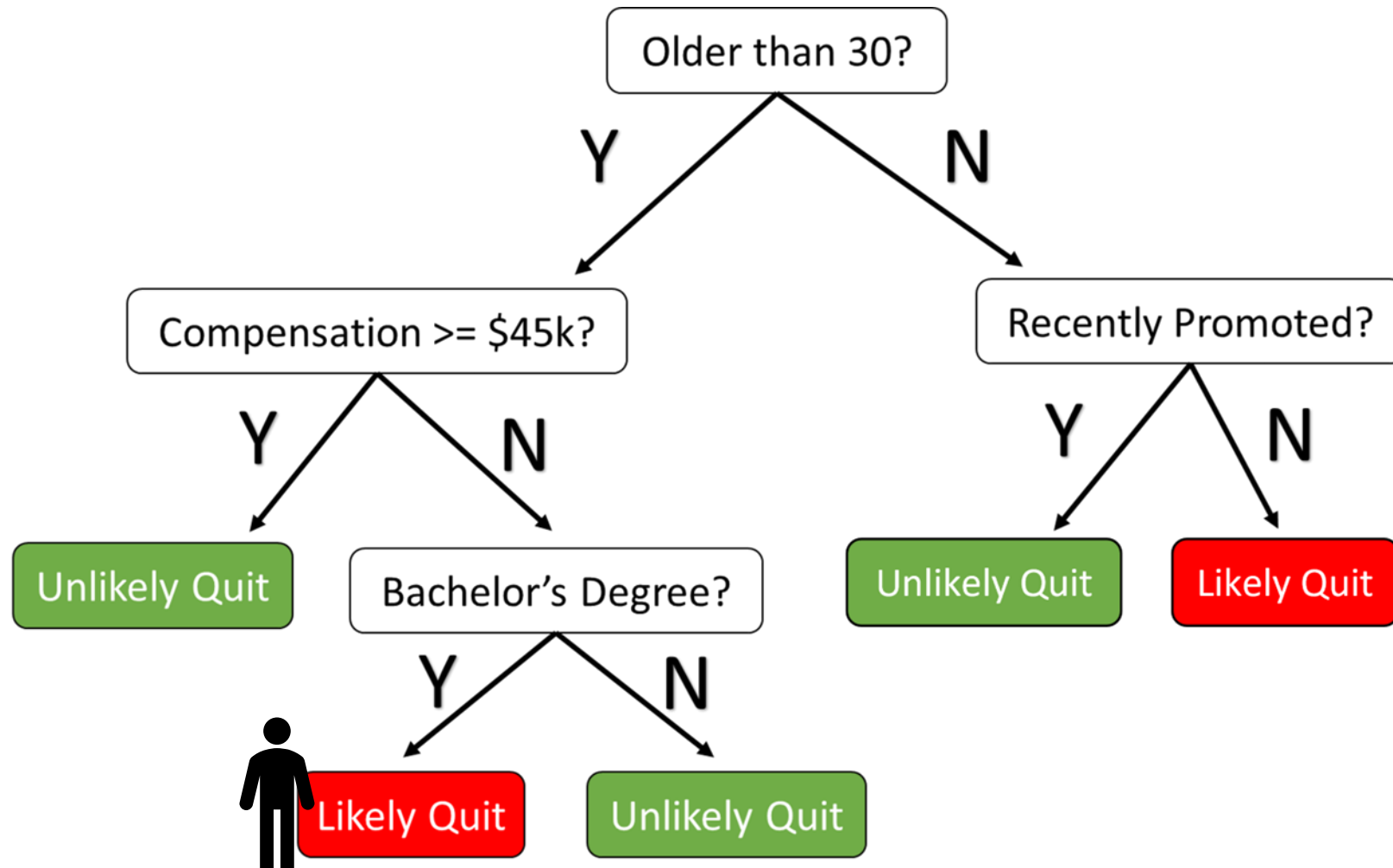
Meet the algorithms: Tree-based learning



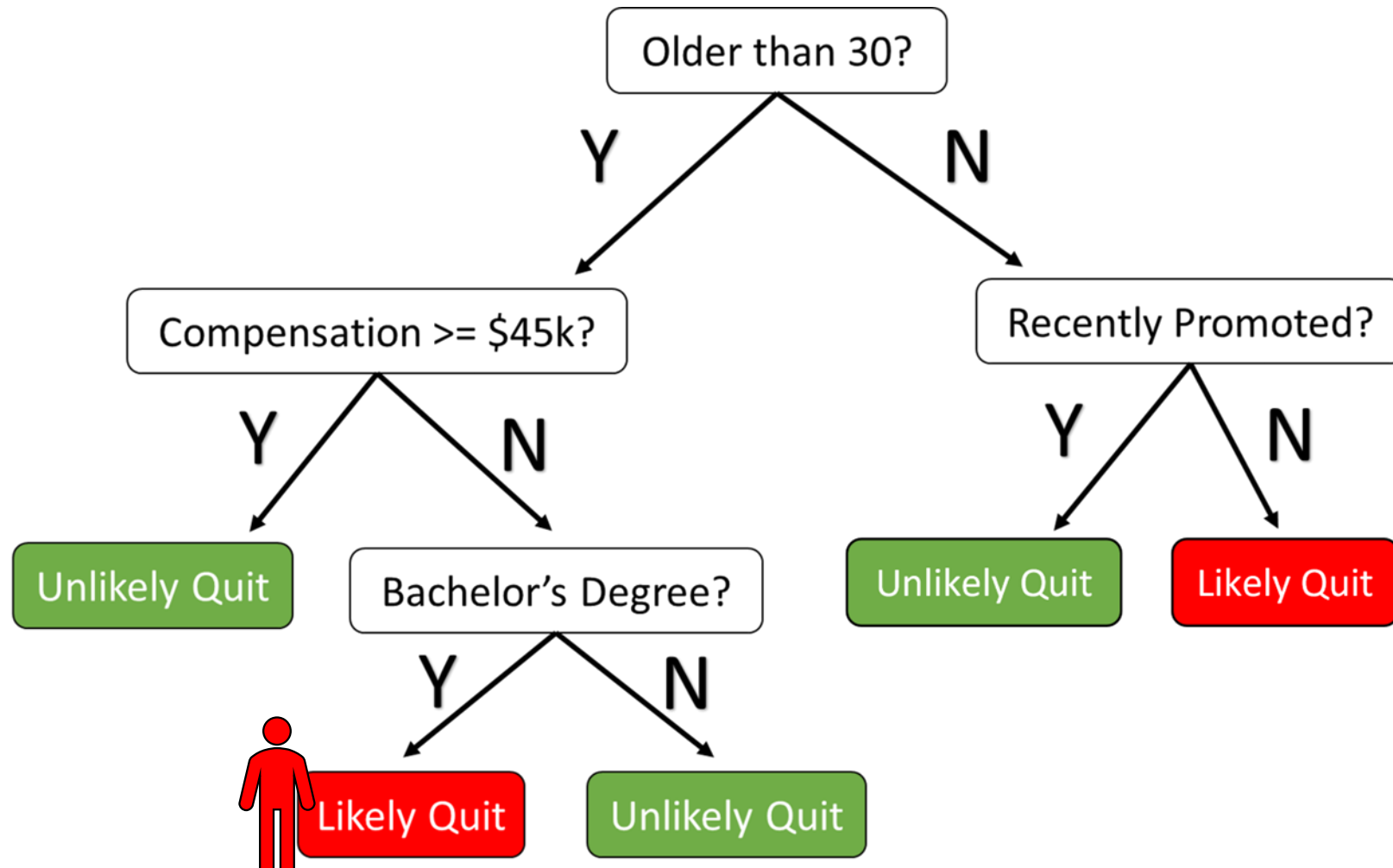
Meet the algorithms: Tree-based learning



Meet the algorithms: Tree-based learning

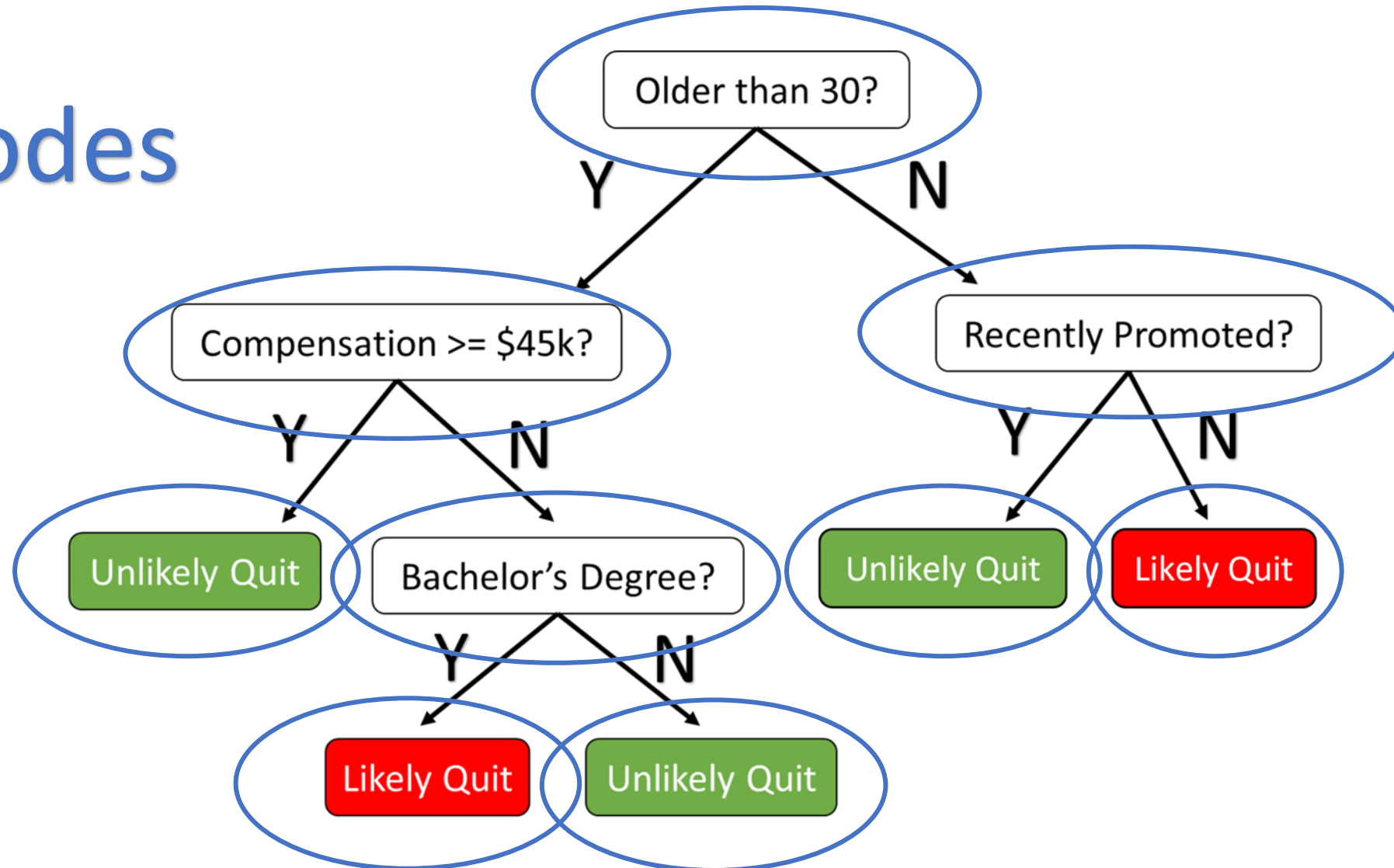


Meet the algorithms: Tree-based learning



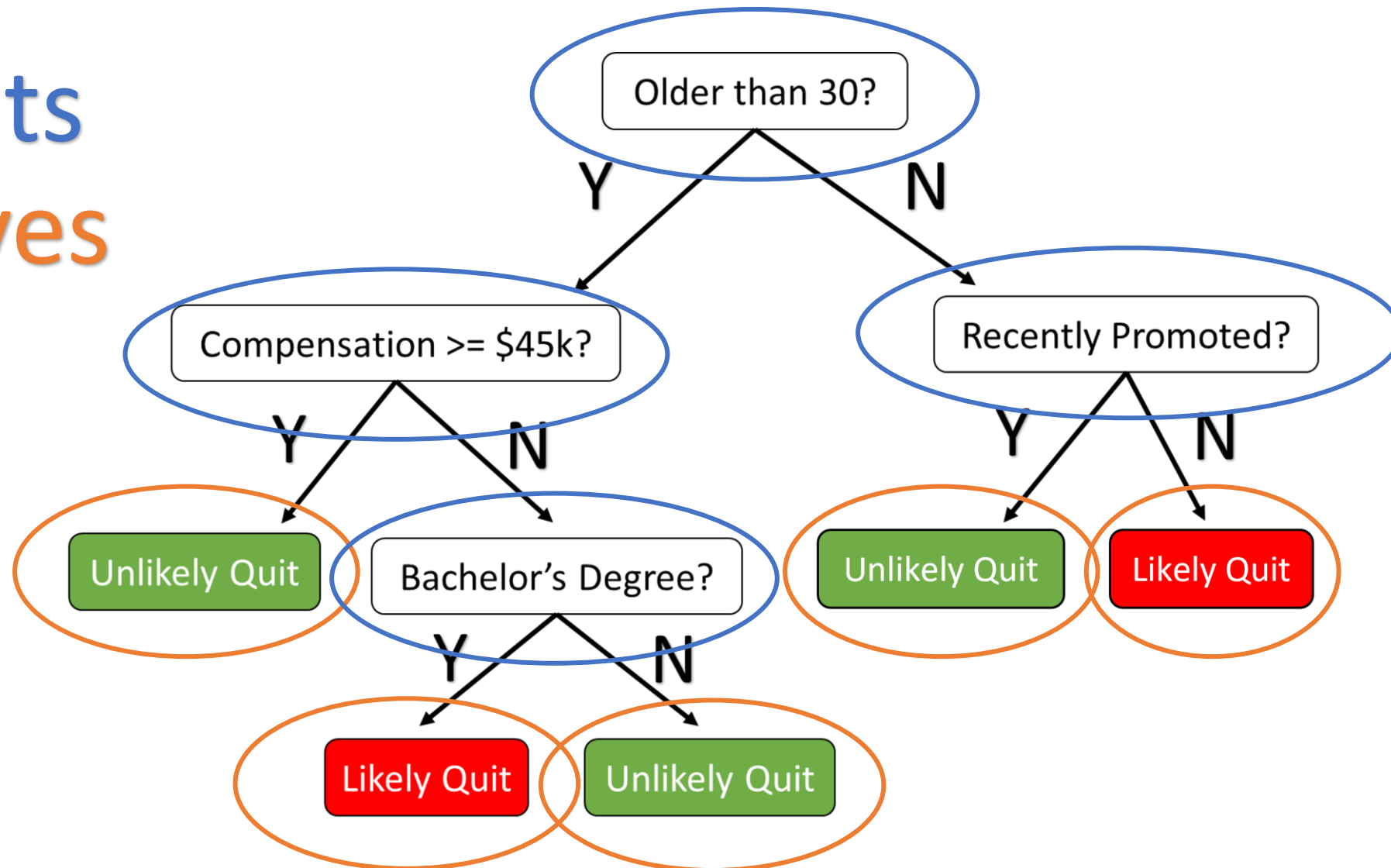
Meet the algorithms: Tree-based learning

Nodes



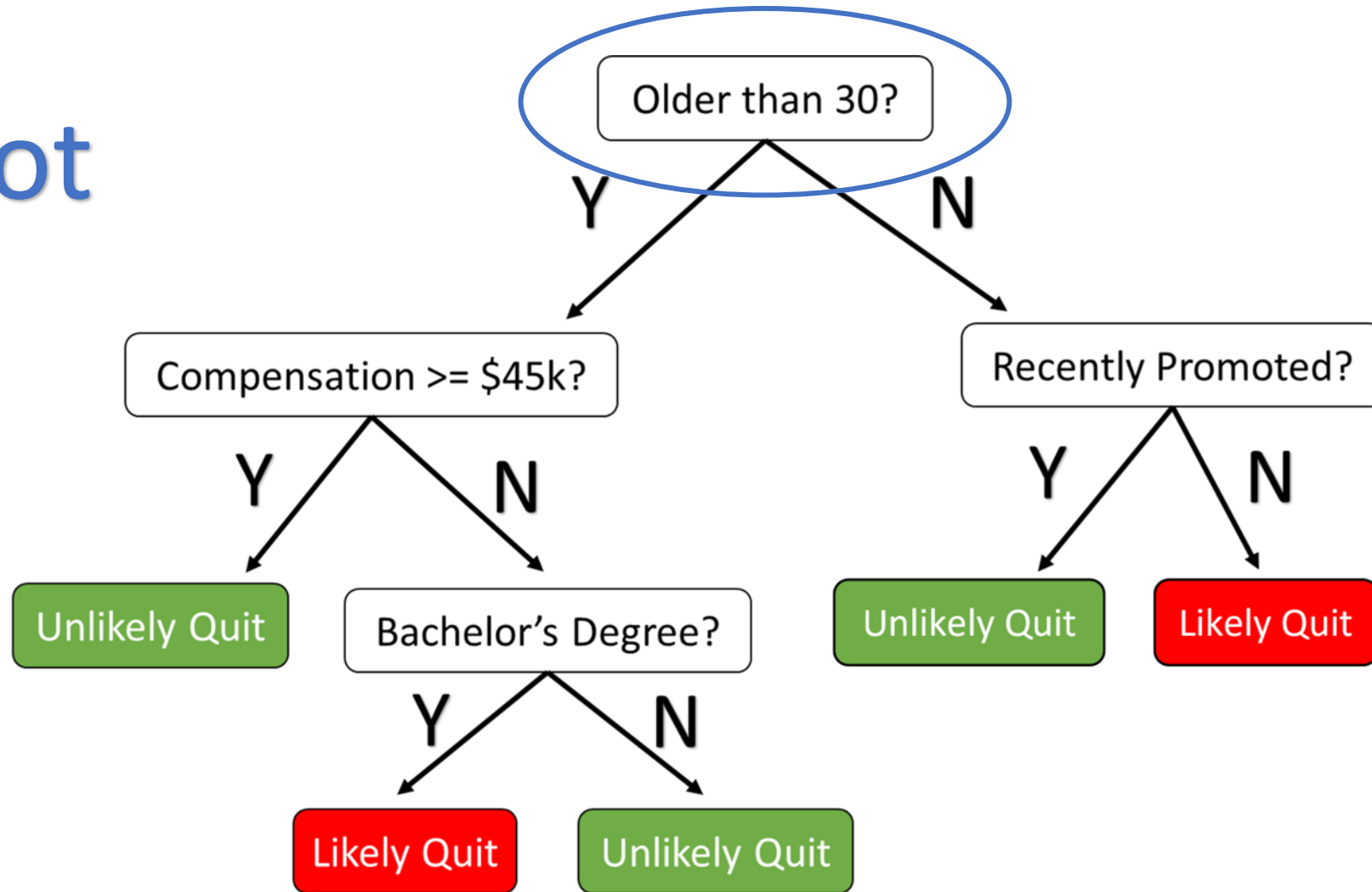
Meet the algorithms: Tree-based learning

Splits
Leaves



Meet the algorithms: Tree-based learning

Root

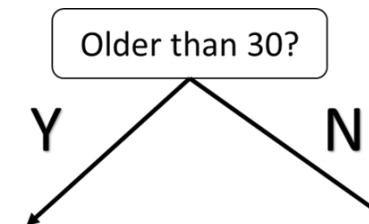


Meet the algorithms: Tree-based learning

Pay	Age	Promoted?	College?	Quit?
30	60	No	No	No
34	64	No	No	No
42	21	Yes	Yes	No
32	35	Yes	Yes	Yes
40	50	No	Yes	Yes
44	23	No	No	Yes
44	27	No	No	Yes
46	33	No	Yes	No
80	40	No	Yes	No
160	18	Yes	No	No
60	19	No	Yes	Yes
120	25	No	Yes	Yes

Meet the algorithms: Tree-based learning

Pay	Age	Promoted?	College	Quit?
30	60	No	No	No
34	64	No	No	No
32	35	Yes	Yes	Yes
40	50	No	Yes	Yes
46	33	No	Yes	No
80	40	No	Yes	No

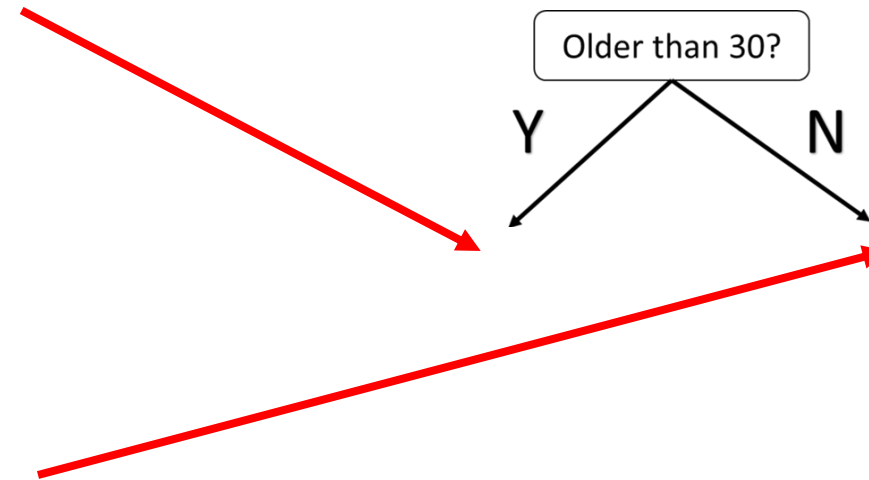


Pay	Age	Promoted?	College	Quit?
42	21	Yes	Yes	No
44	23	No	No	Yes
44	27	No	No	Yes
160	18	Yes	No	No
60	19	No	Yes	Yes
120	25	No	Yes	Yes

Meet the algorithms: Tree-based learning

Pay	Age	Promoted?	College	Quit?
30	60	No	No	No
34	64	No	No	No
32	35	Yes	Yes	Yes
40	50	No	Yes	Yes
46	33	No	Yes	No
80	40	No	Yes	No

Pay	Age	Promoted?	College	Quit?
42	21	Yes	Yes	No
44	23	No	No	Yes
44	27	No	No	Yes
160	18	Yes	No	No
60	19	No	Yes	Yes
120	25	No	Yes	Yes



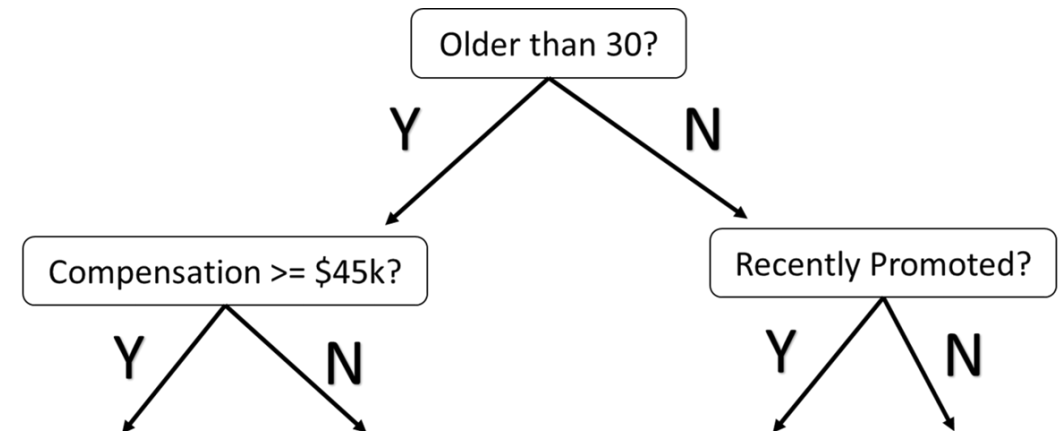
Meet the algorithms: Tree-based learning

Pay	Age	Promoted?	College	Quit?
30	60	No	No	No
34	64	No	No	No
32	35	Yes	Yes	Yes
40	50	No	Yes	Yes

Pay	Age	Promoted?	College	Quit?
46	33	No	Yes	No
80	40	No	Yes	No

Pay	Age	Promoted?	College	Quit?
42	21	Yes	Yes	No
160	18	Yes	No	No

Pay	Age	Promoted?	College	Quit?
44	23	No	No	Yes
44	27	No	No	Yes
60	19	No	Yes	Yes
120	25	No	Yes	Yes



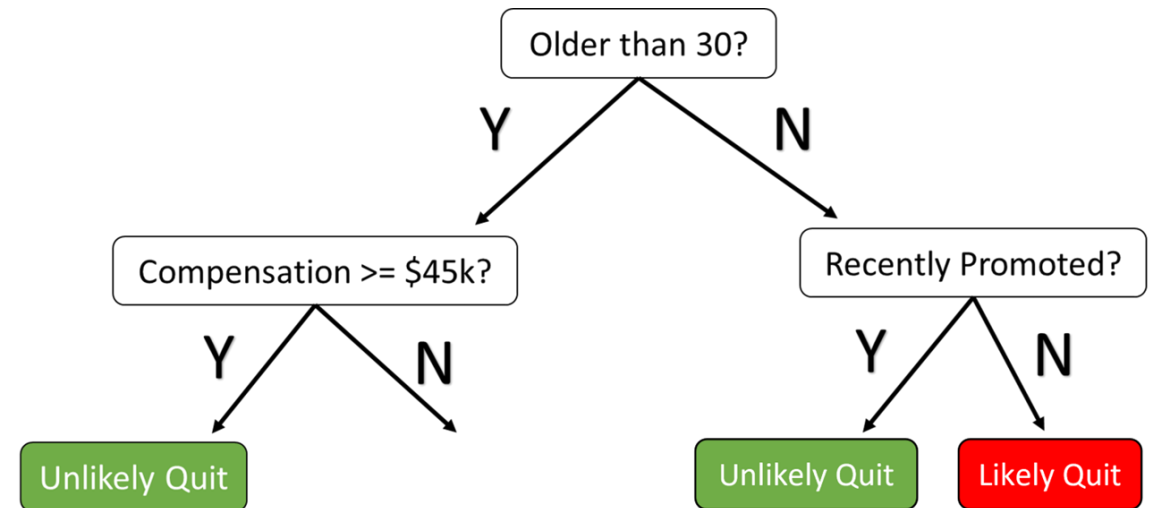
Meet the algorithms: Tree-based learning

Pay	Age	Promoted?	College	Quit?
30	60	No	No	No
34	64	No	No	No
32	35	Yes	Yes	Yes
40	50	No	Yes	Yes

Pay	Age	Promoted?	College	Quit?
46	33	No	Yes	No
80	40	No	Yes	No

Pay	Age	Promoted?	College	Quit?
42	21	Yes	Yes	No
160	18	Yes	No	No

Pay	Age	Promoted?	College	Quit?
44	23	No	No	Yes
44	27	No	No	Yes
60	19	No	Yes	Yes
120	25	No	Yes	Yes



Meet the algorithms: Tree-based learning

Pay	Age	Promoted?	College	Quit?
32	35	Yes	Yes	Yes
40	50	No	Yes	Yes

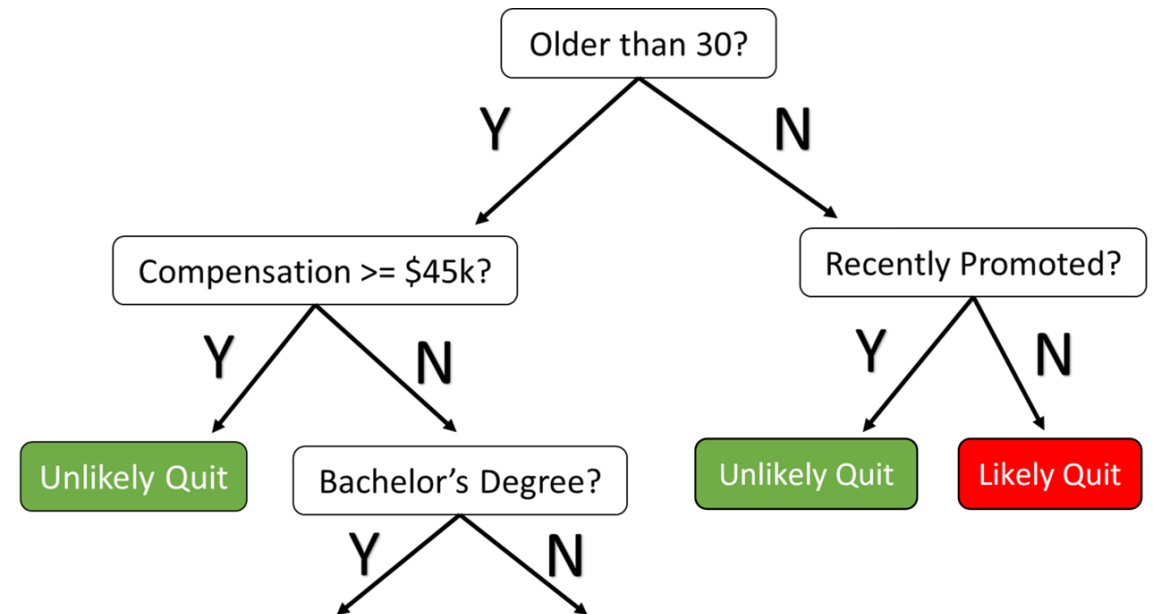
Pay	Age	Promoted?	College	Quit?
30	60	No	No	No
34	64	No	No	No

.....

Pay	Age	Promoted?	College	Quit?
46	33	No	Yes	No
80	40	No	Yes	No

Pay	Age	Promoted?	College	Quit?
42	21	Yes	Yes	No
160	18	Yes	No	No

Pay	Age	Promoted?	College	Quit?
44	23	No	No	Yes
44	27	No	No	Yes
60	19	No	Yes	Yes
120	25	No	Yes	Yes



Meet the algorithms: Tree-based learning

Pay	Age	Promoted?	College	Quit?
32	35	Yes	Yes	Yes
40	50	No	Yes	Yes

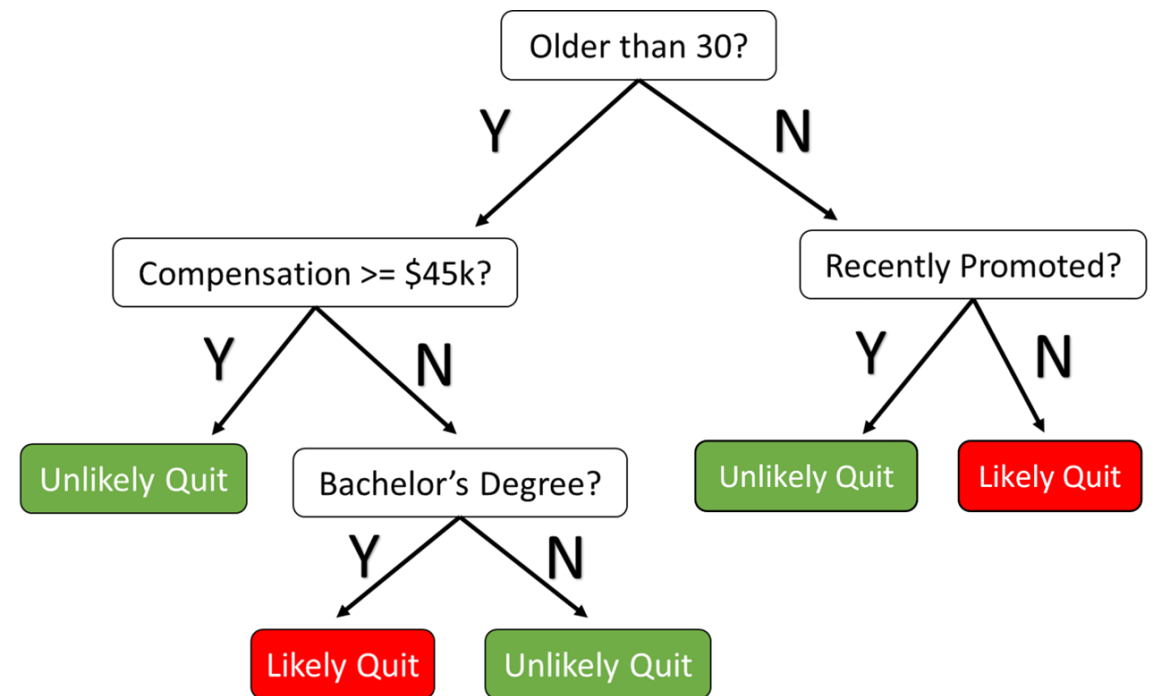
Pay	Age	Promoted?	College	Quit?
30	60	No	No	No
34	64	No	No	No

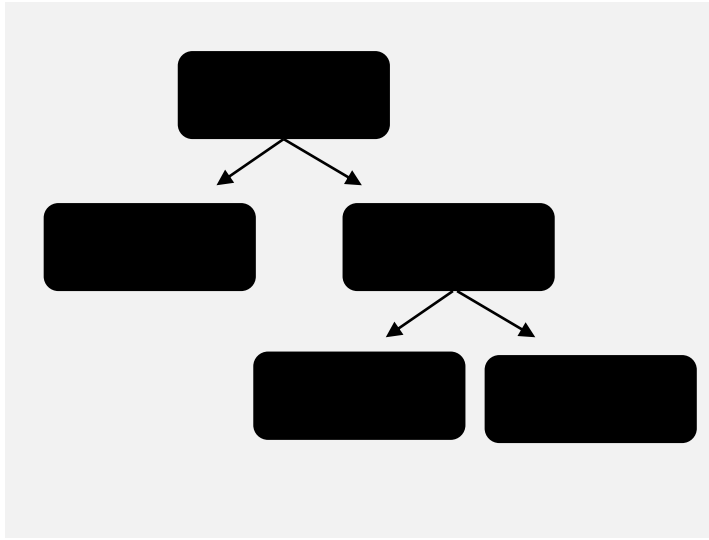
.....

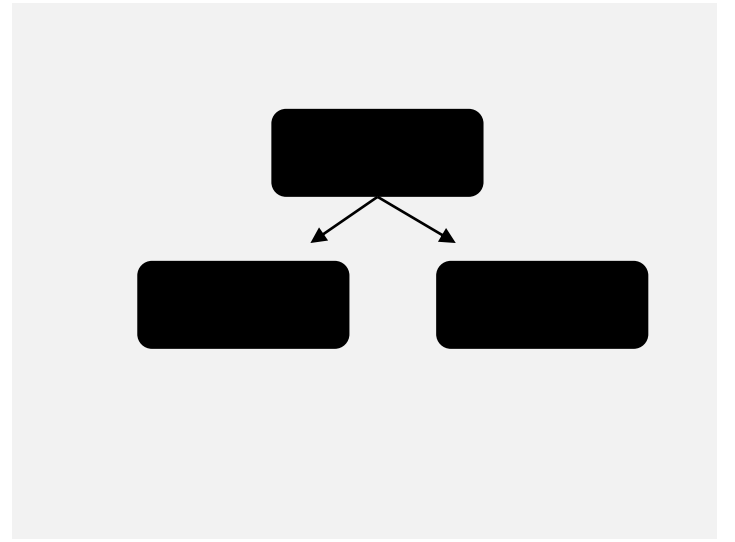
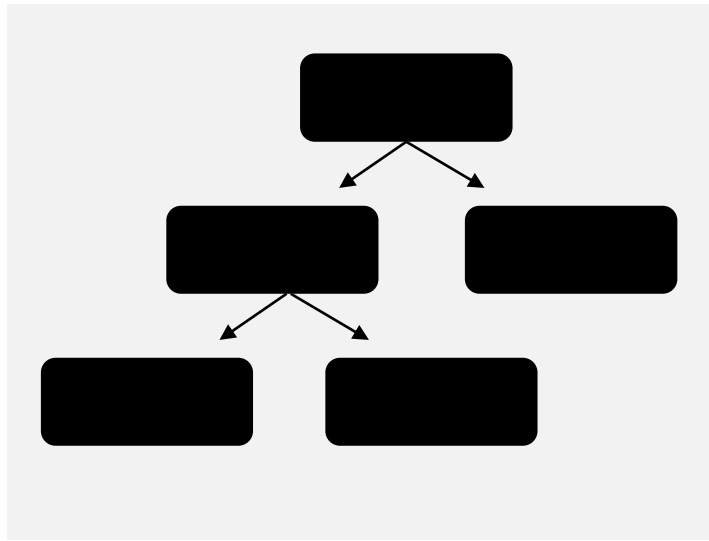
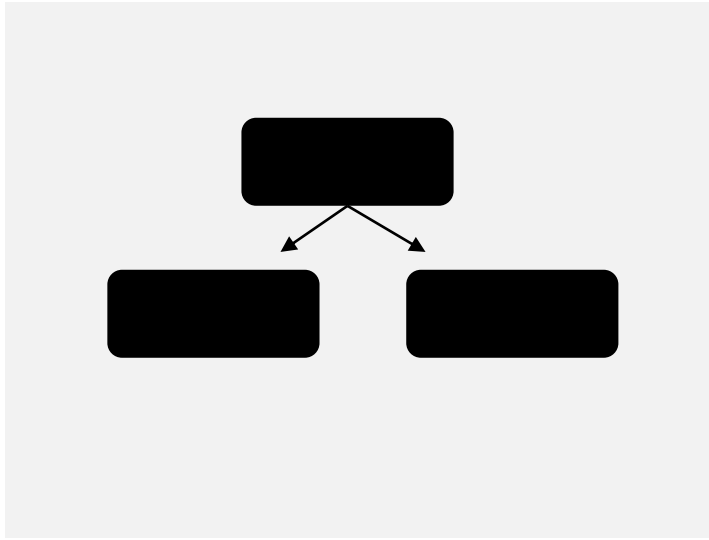
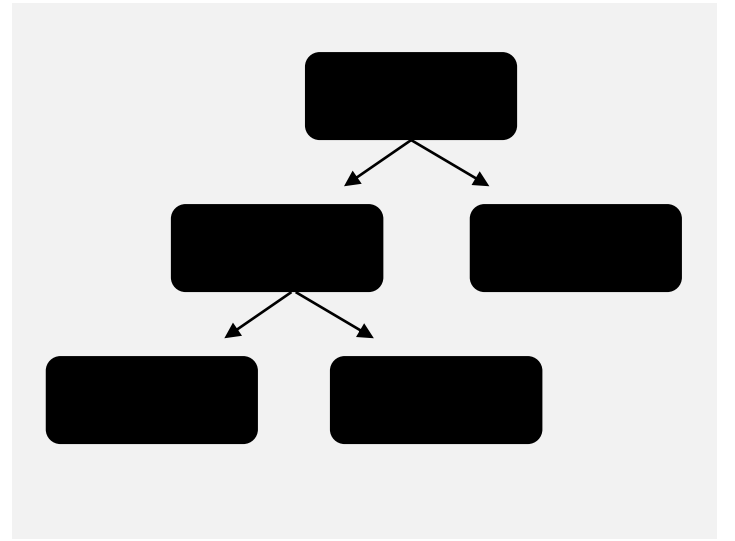
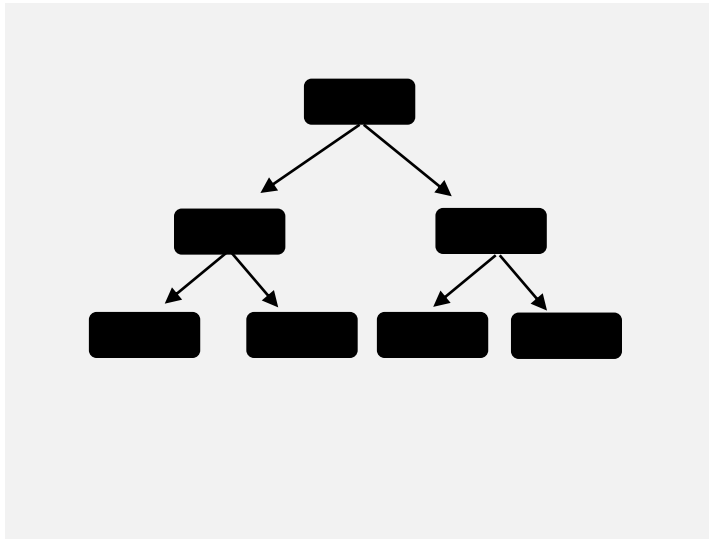
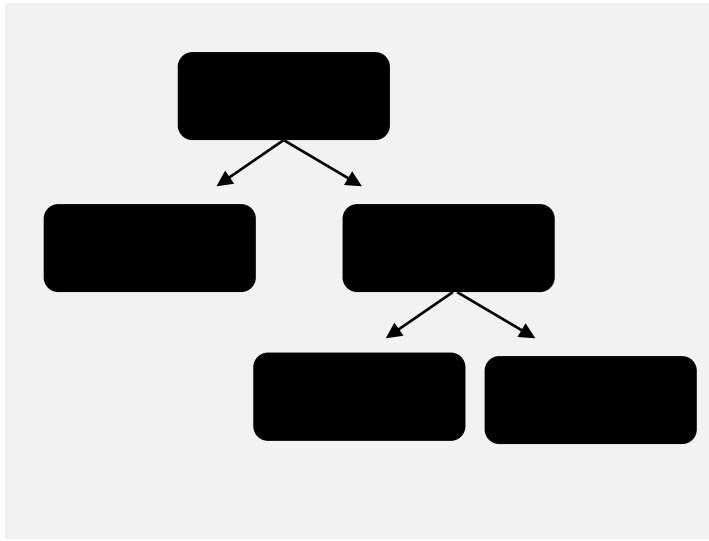
Pay	Age	Promoted?	College	Quit?
46	33	No	Yes	No
80	40	No	Yes	No

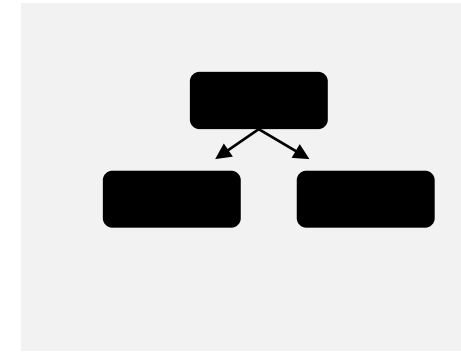
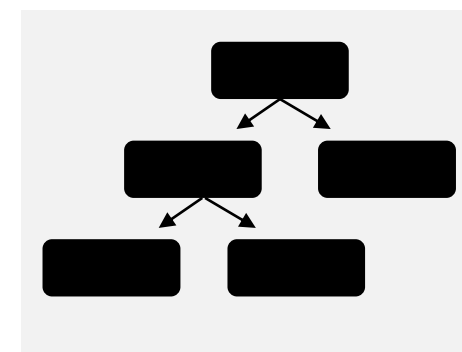
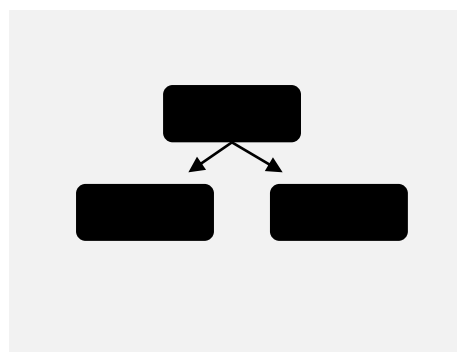
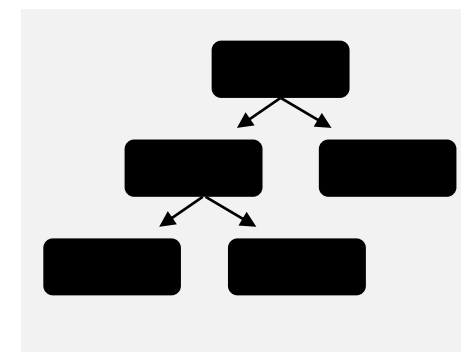
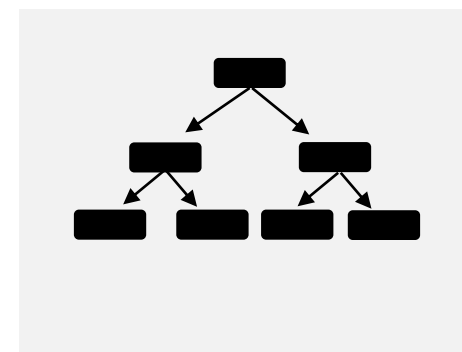
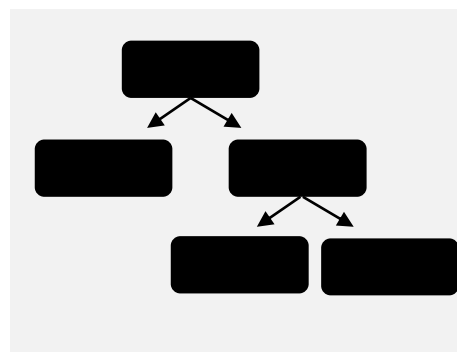
Pay	Age	Promoted?	College	Quit?
42	21	Yes	Yes	No
160	18	Yes	No	No

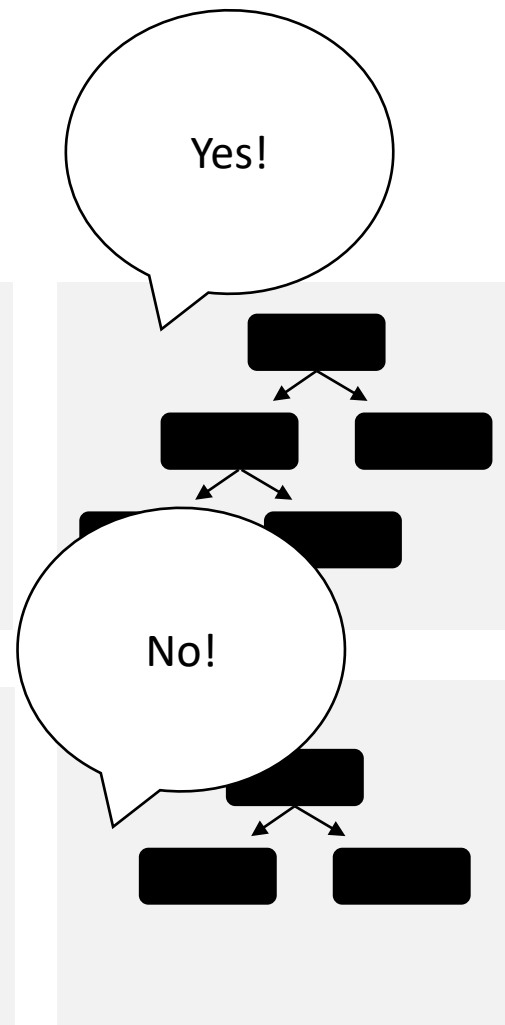
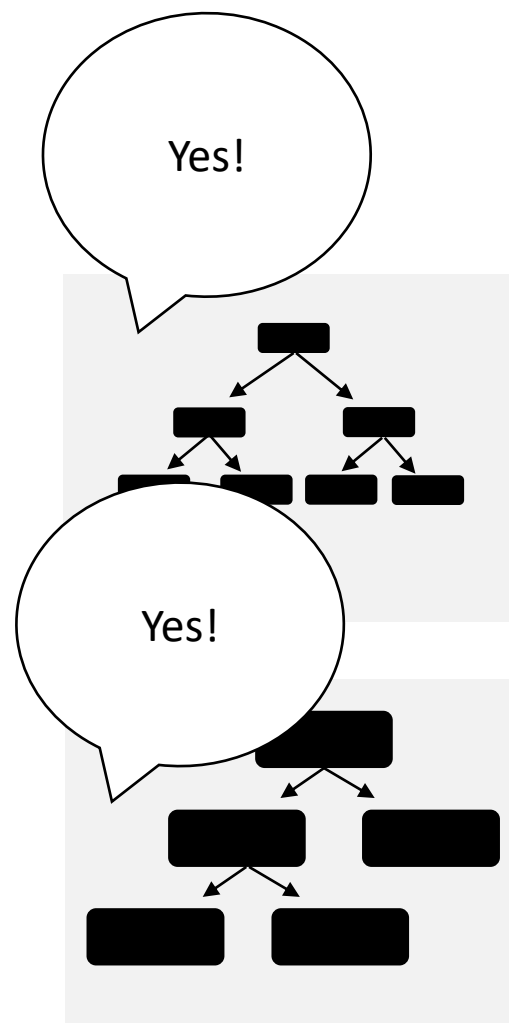
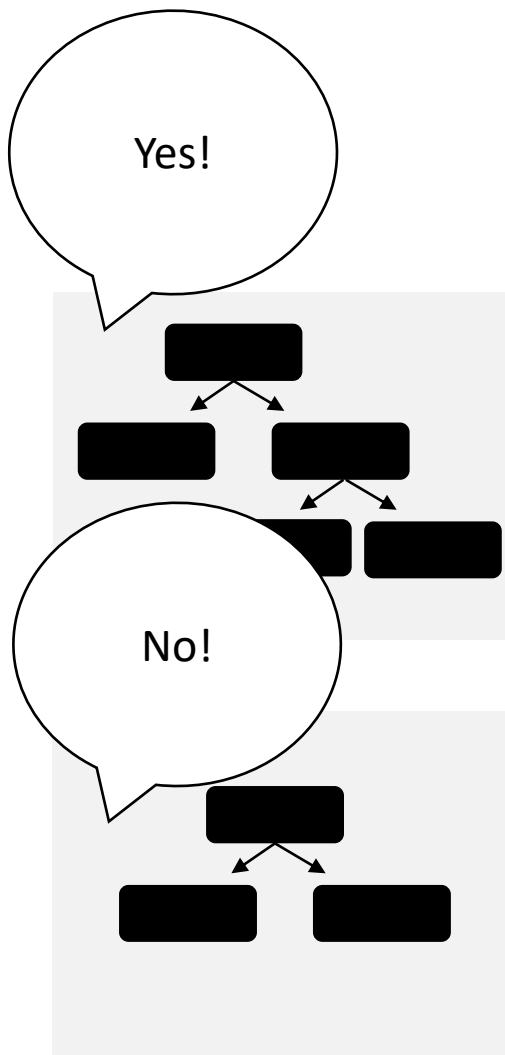
Pay	Age	Promoted?	College	Quit?
44	23	No	No	Yes
44	27	No	No	Yes
60	19	No	Yes	Yes
120	25	No	Yes	Yes










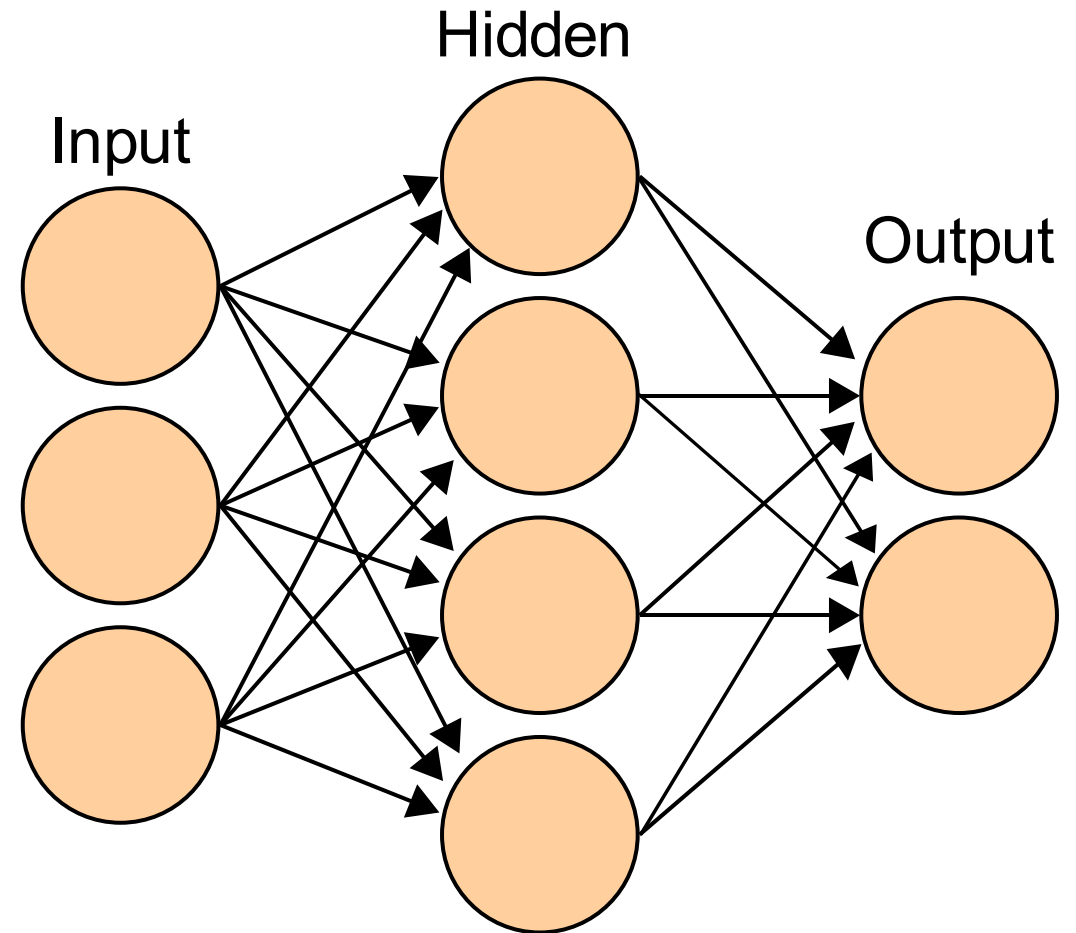




Motivating deep learning

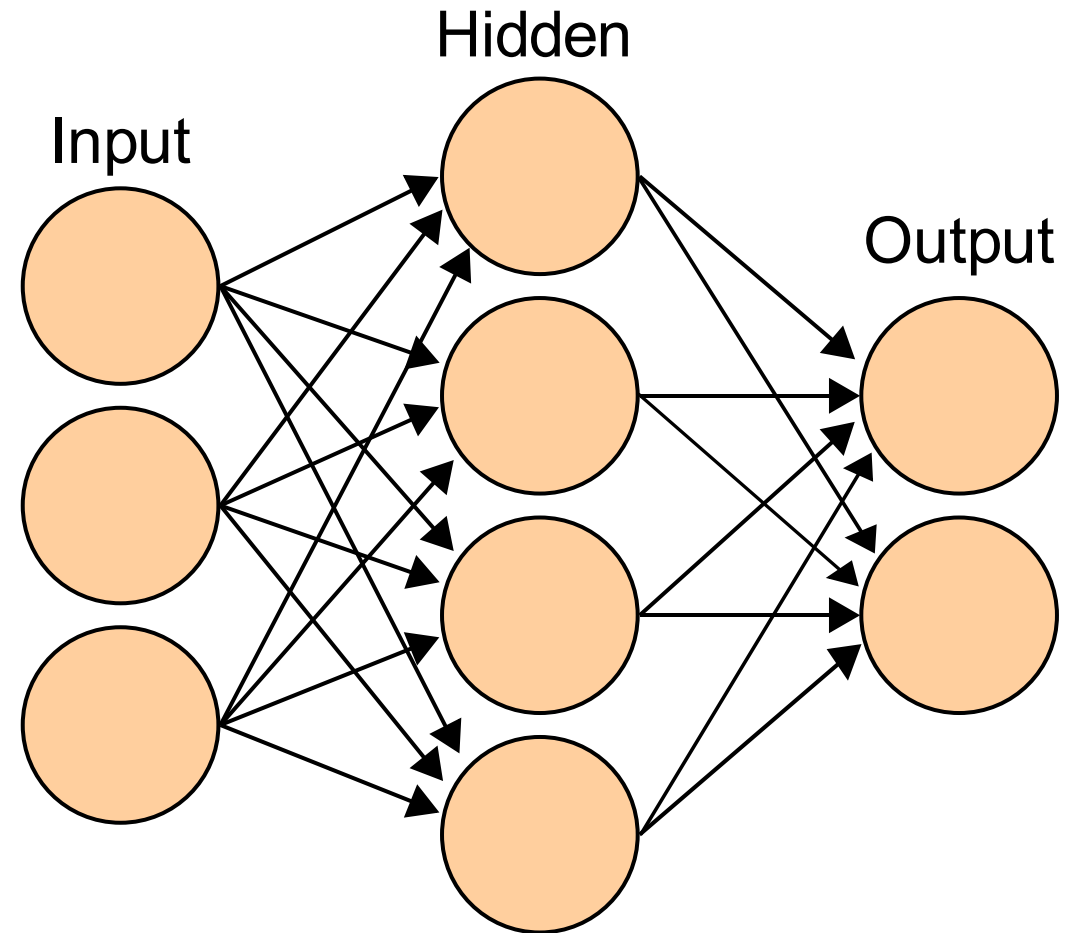
- In animal brains, **neurons** are connected to each other (at sites called **synapses**)
- When a neuron becomes **activated**, it typically sends signals to other neurons which in turn can activate them (in some cases it also inhibits them)
- So, when your brain receives sense data, a giant **network** of interconnected neurons fire so that you can make sense of and respond to what you're seeing
- The human brain is, in some ways, the thing we want to mimic with these algorithms. So why don't we try making an ML algorithm that copies these ideas?

Meet the algorithms: Deep learning



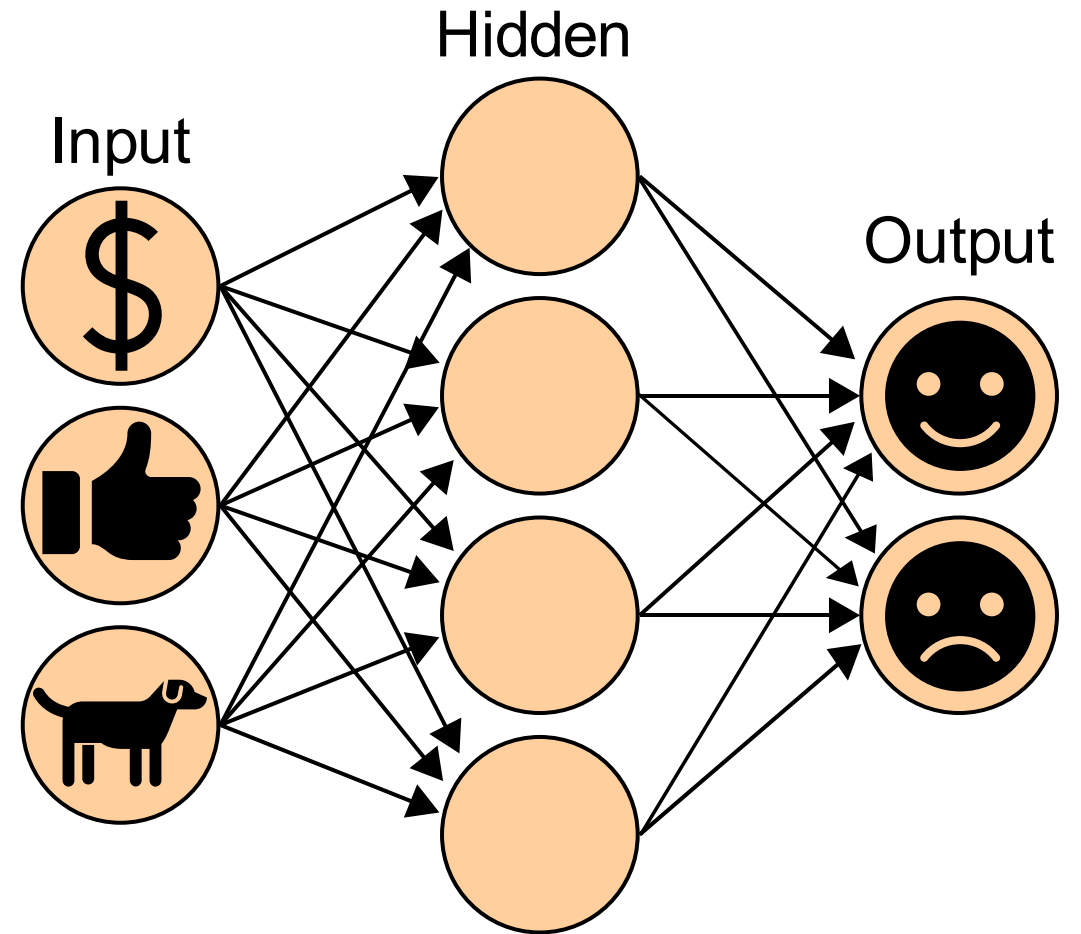
Meet the algorithms: Deep learning

Happy?	Rich?	Enjoys work?	Has dog?
Yes	No	Yes	Yes
No	No	No	No
Yes	Yes	No	Yes
No	Yes	Yes	No



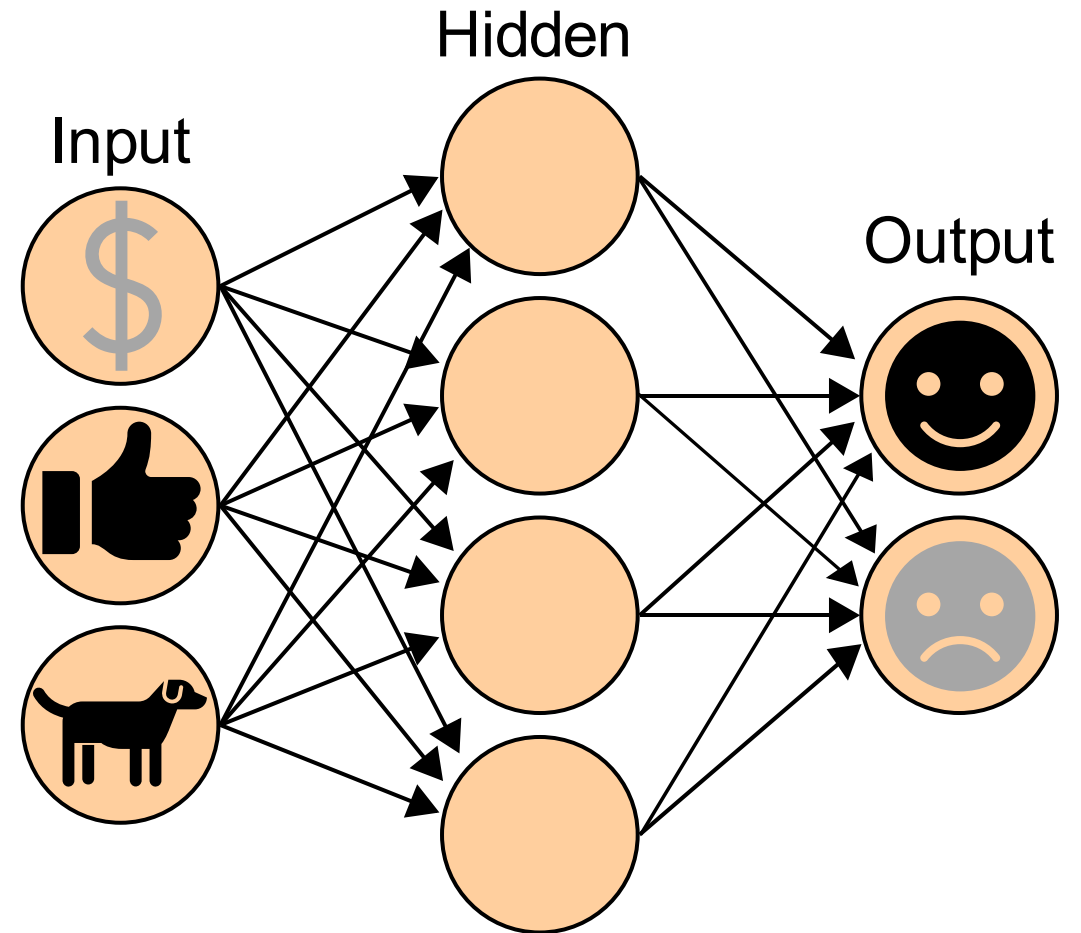
Meet the algorithms: Deep learning

Happy?	Rich?	Enjoys work?	Has dog?
Yes	No	Yes	Yes
No	No	No	No
Yes	Yes	No	Yes
No	Yes	Yes	No



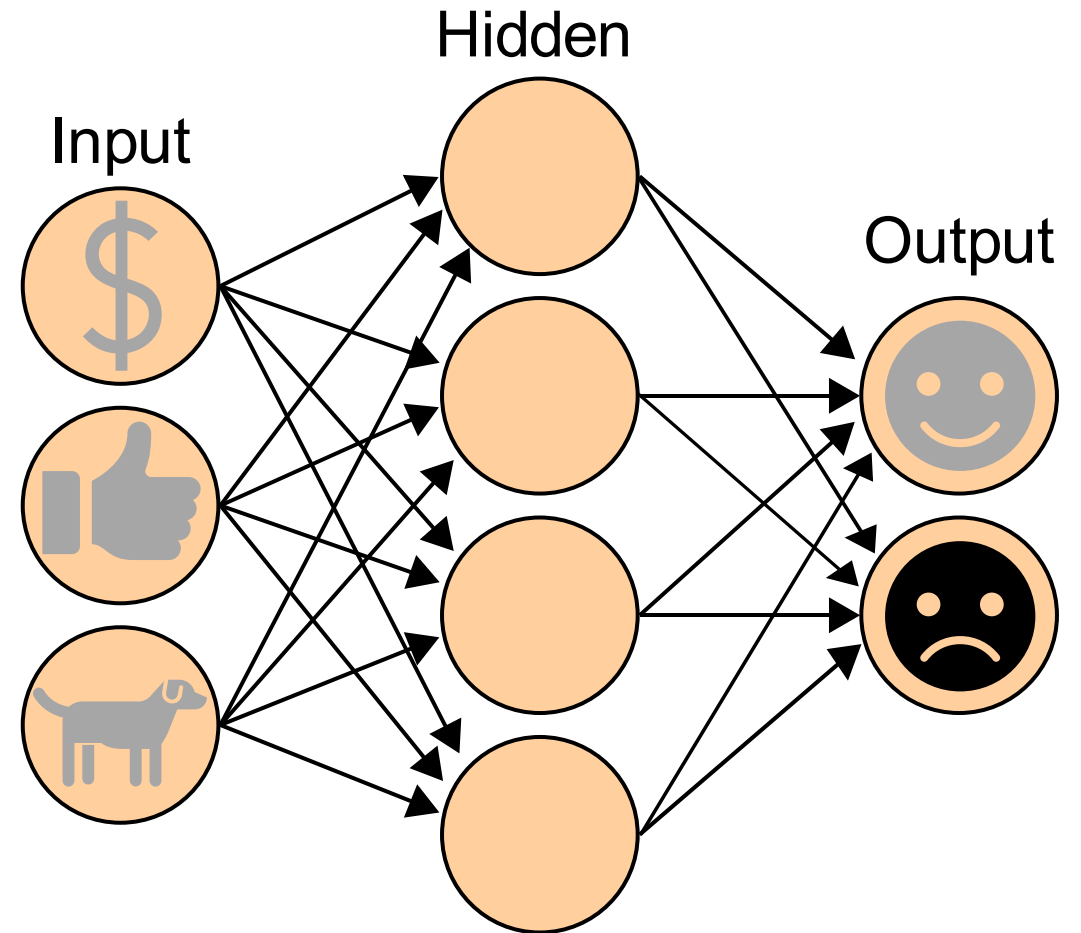
Meet the algorithms: Deep learning

Happy?	Rich?	Enjoys work?	Has dog?
Yes	No	Yes	Yes
No	No	No	No
Yes	Yes	No	Yes
No	Yes	Yes	No



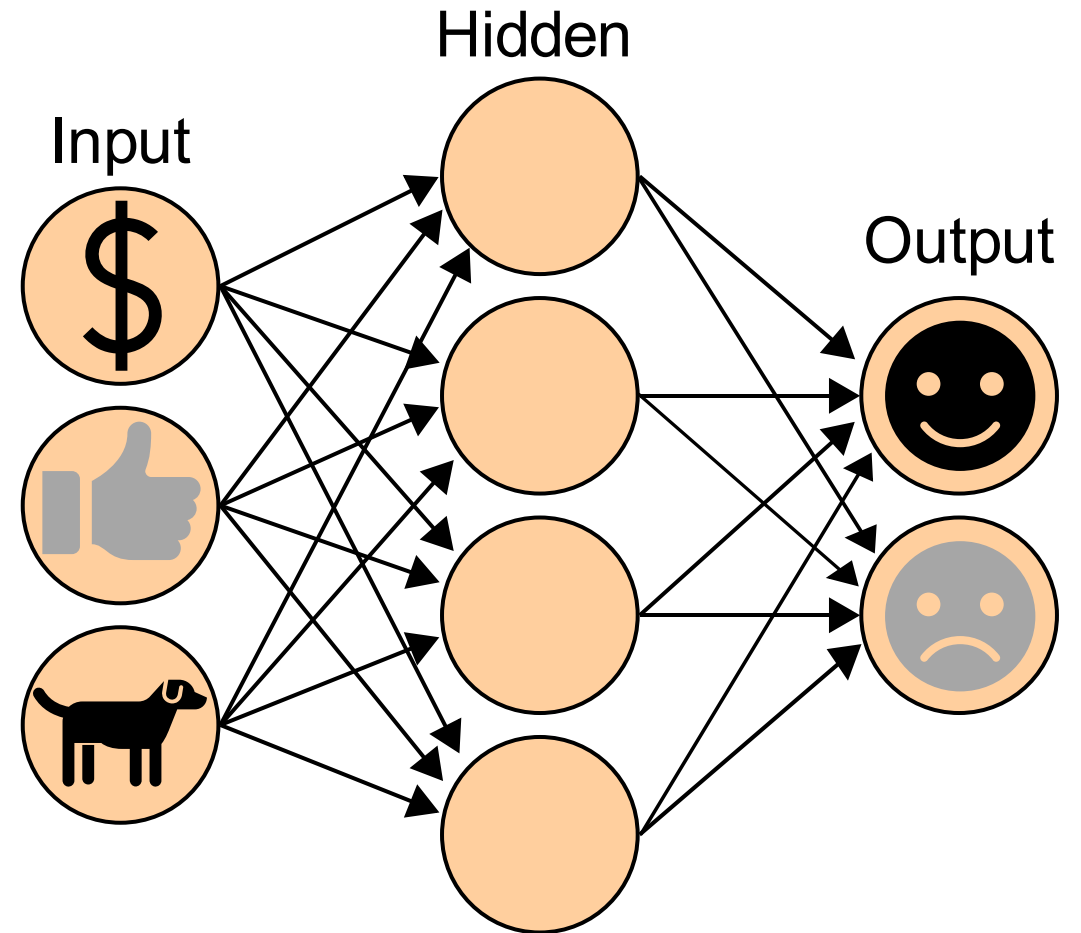
Meet the algorithms: Deep learning

Happy?	Rich?	Enjoys work?	Has dog?
Yes	No	Yes	Yes
No	No	No	No
Yes	Yes	No	Yes
No	Yes	Yes	No



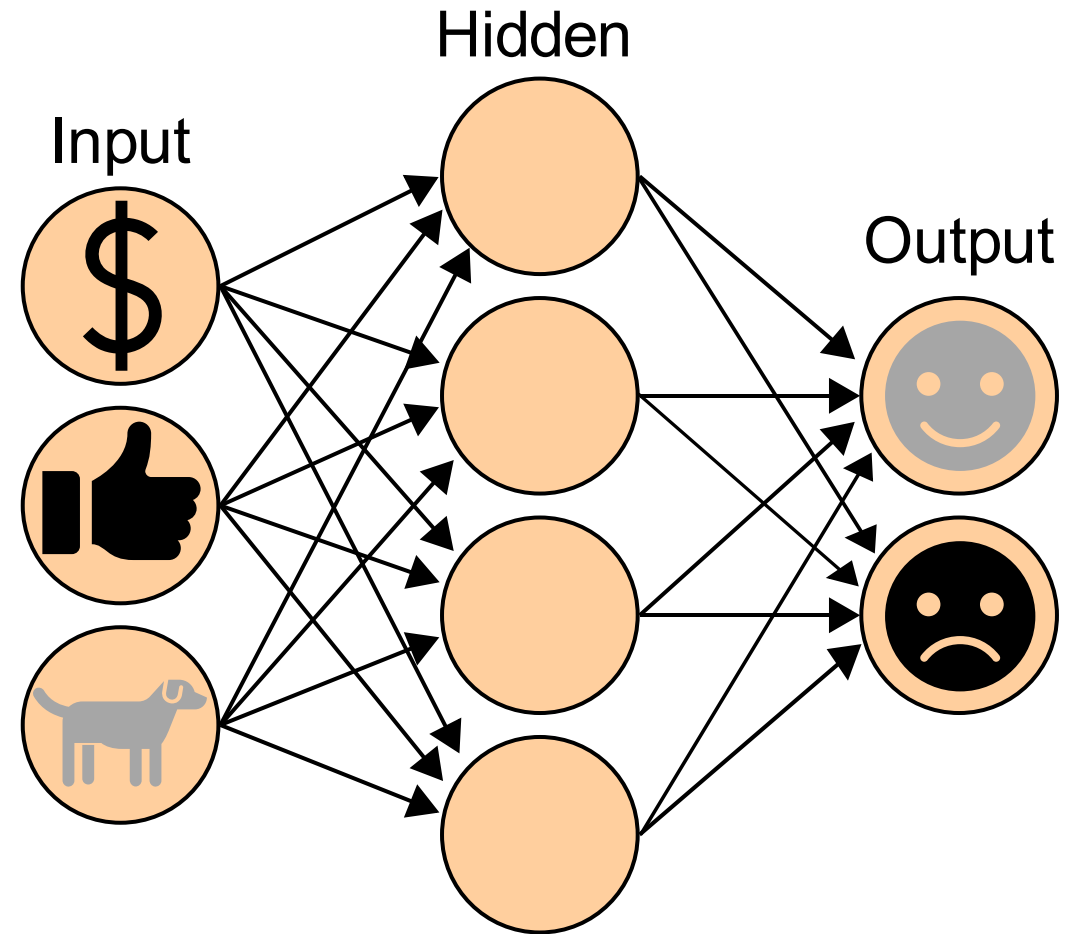
Meet the algorithms: Deep learning

Happy?	Rich?	Enjoys work?	Has dog?
Yes	No	Yes	Yes
No	No	No	No
Yes	Yes	No	Yes
No	Yes	Yes	No



Meet the algorithms: Deep learning

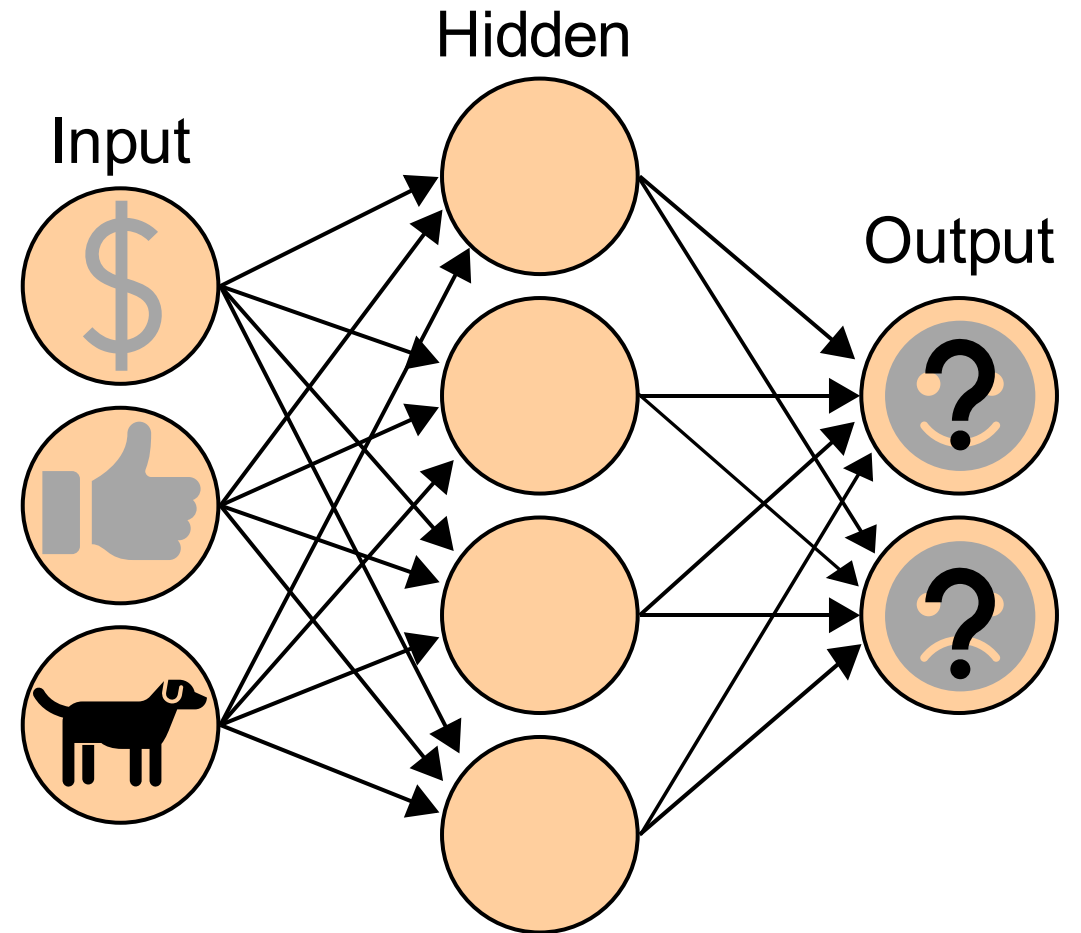
Happy?	Rich?	Enjoys work?	Has dog?
Yes	No	Yes	Yes
No	No	No	No
Yes	Yes	No	Yes
No	Yes	Yes	No



Meet the algorithms: Deep learning

Happy?	Rich?	Enjoys work?	Has dog?
Yes	No	Yes	Yes
No	No	No	No
Yes	Yes	No	Yes
No	Yes	Yes	No

?	No	No	Yes
---	----	----	-----



Other important topics

- Cross-validation
- Support vector machines
- Unsupervised ML (K-means, autoencoders)
- Reinforcement learning
- Adversarial learning

Reminders

- I've extended the deadline for discussion papers for the organizational theory readings... They are now due this Thursday (6/30) before class.
- Method Module 1 is available on Canvas now. If you want to complete it for credit on your final paper, finish it and turn it in on Canvas before class a week from today (7/5)
- Office hours are tomorrow in my office from 11 AM to noon.

See everyone Thursday!

If you'd like to take a quick look at the Methods Module, I'm happy to do so now