
INFO 2950: Intro to Data Science

Lecture 25
2023-11-27

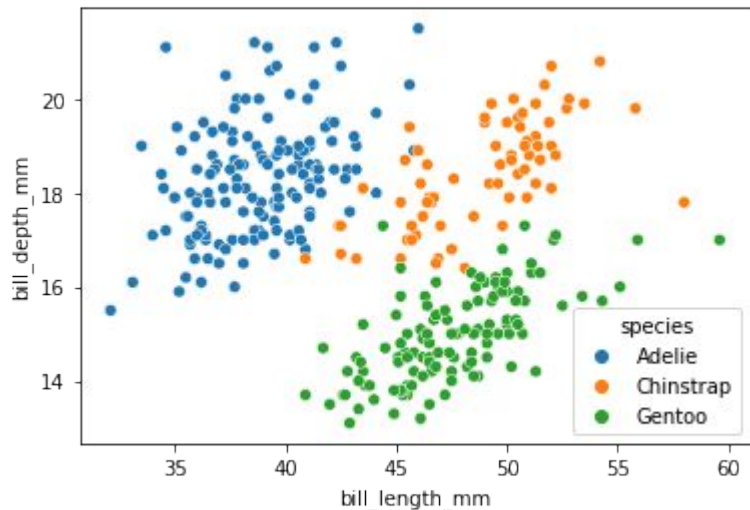
Agenda

1. Classification
 - a. Linear separability
2. Continuous output prediction
 - a. Linear Regression
 - b. Perceptrons
 - c. Neural Networks

Classifying

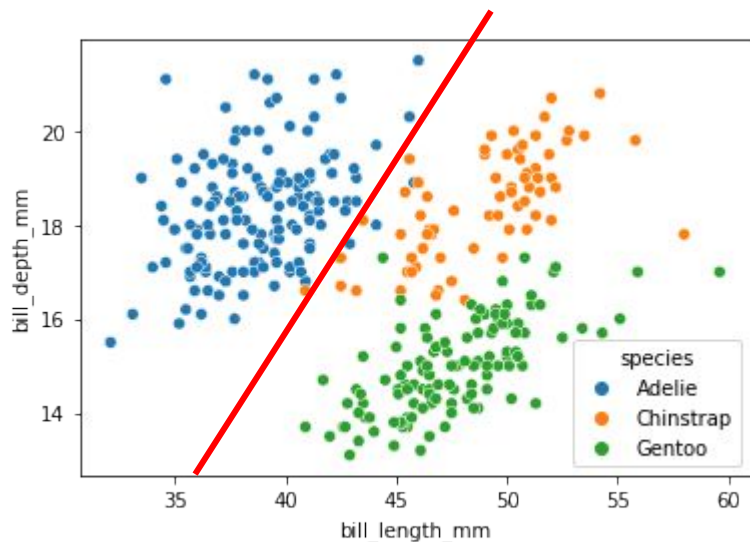
- To classify data (e.g., to predict whether binary output y is 0 or 1), we've learned how to use:
 - (Logistic) regression
 - Naive Bayes
 - K-means clustering

Linear separability



Given bill length and depth, can you draw a line that reliably separates **Adelie** penguins from **Chinstrap** and **Gentoo** penguins?

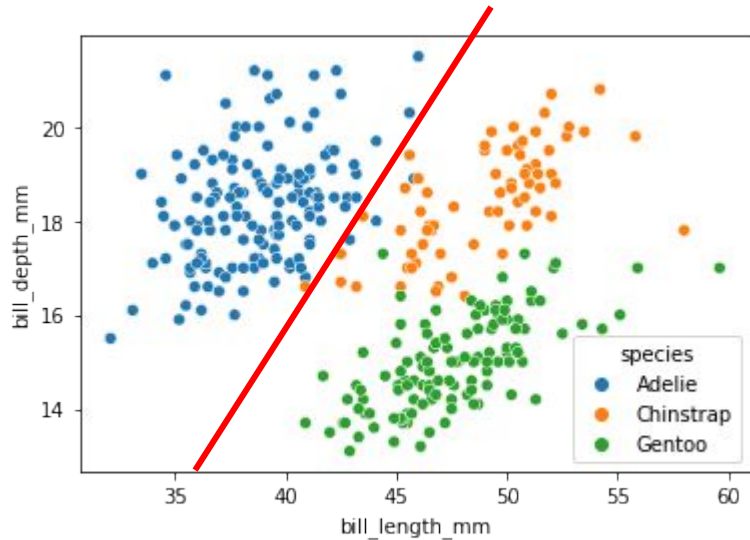
Linear separability



Given bill length and depth, can you draw a line that reliably separates **Adelie** penguins from **Chinstrap** and **Gentoo** penguins?

Yes! It's not perfect, but most of the blue dots (**Adelie**) are on one side, and most of the non-blue dots are on the other

Regression = Decision boundary



Math version:

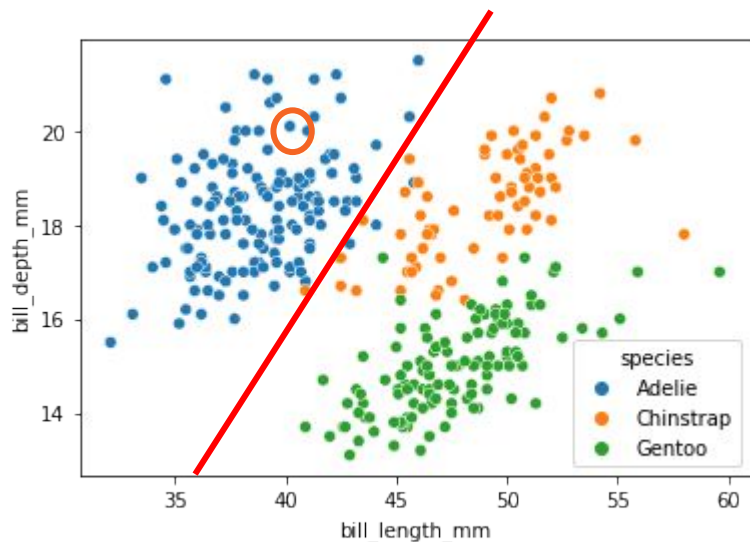
$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

Classification:

$y = 1$ if Adelie, 0 if Chinstrap, 0 if Gentoo

Calculate a linear output



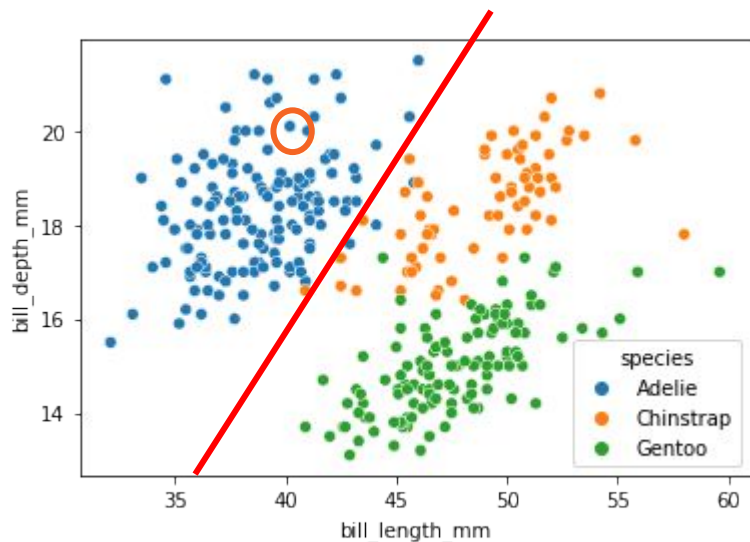
Math version:

$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

Is the value for length=40, depth=20
positive or negative?

Calculate a linear output



Math version:

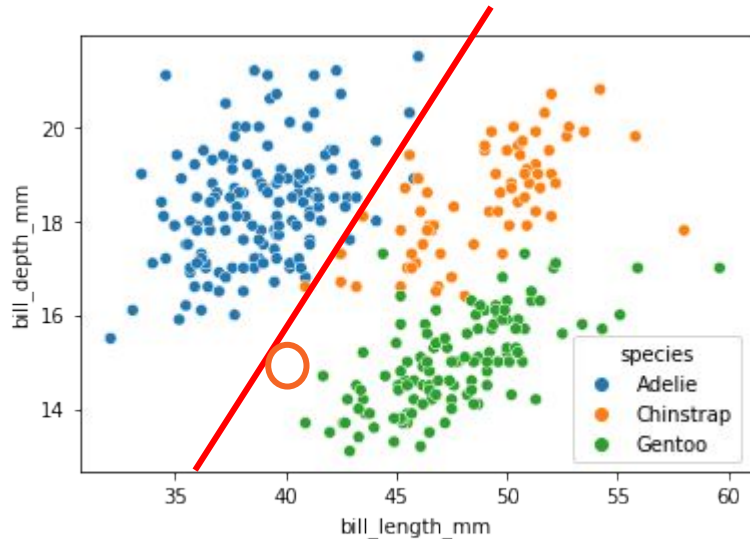
$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

Is the value for length=40, depth=20
positive or negative?

$$25 + (-60) + 40 = 5 \text{ *positive*}$$

Calculate a linear output



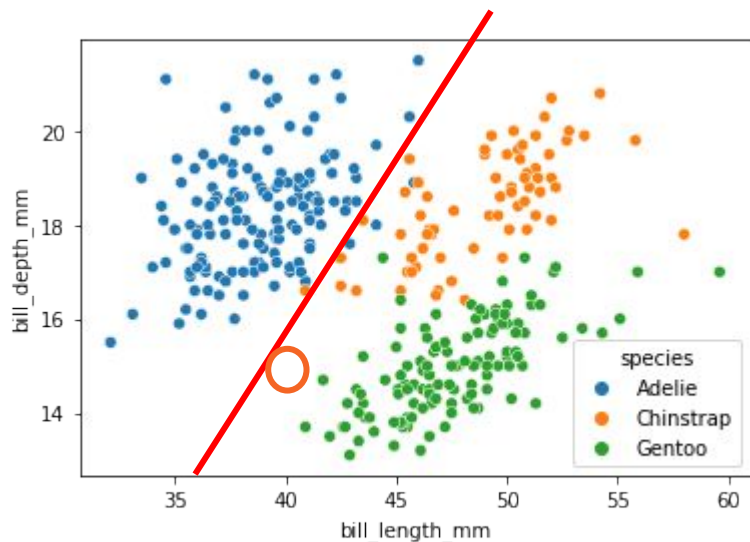
Math version:

$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

**Is the value for length=40, depth=15
positive or negative?**

Calculate a linear output



Math version:

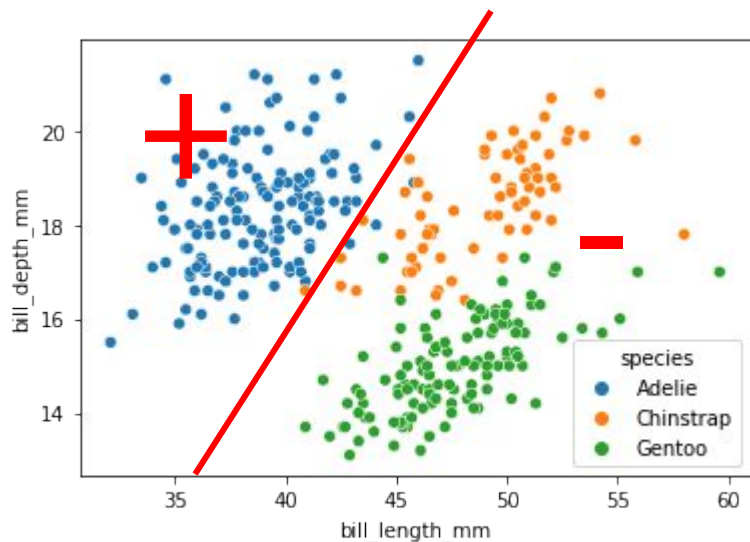
$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

Is the value for length=40, depth=15
positive or negative?

$$25 + (-60) + 30 = -5 \text{ *negative*}$$

Regions of positive / negative



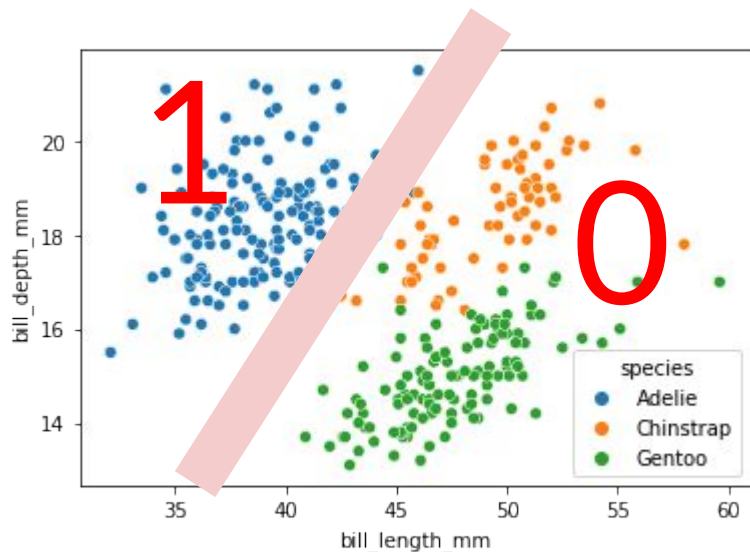
Math version:

$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

Applying these weights for any point left of the line gives us a positive value, and any point right of the line negative

logistic regression



Math version:

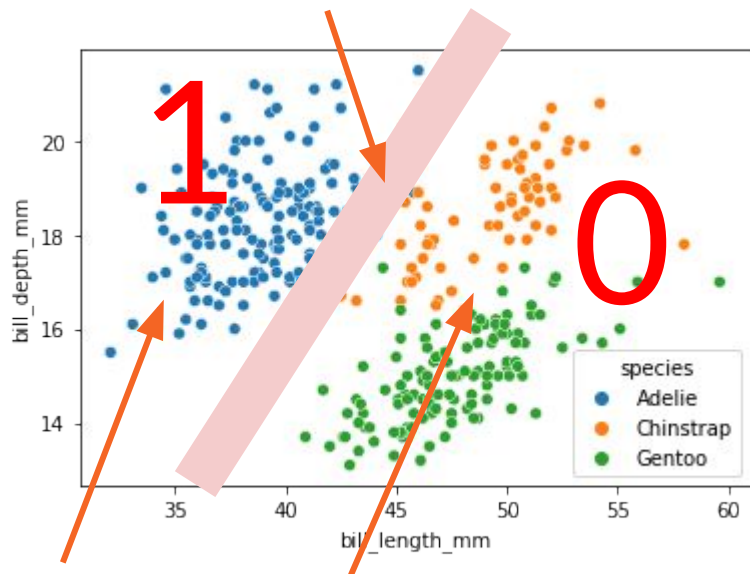
$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

The logistic function squashes negative towards 0 and positive towards 1

logistic regression

Cliff of uncertainty



Plateaus of confidence

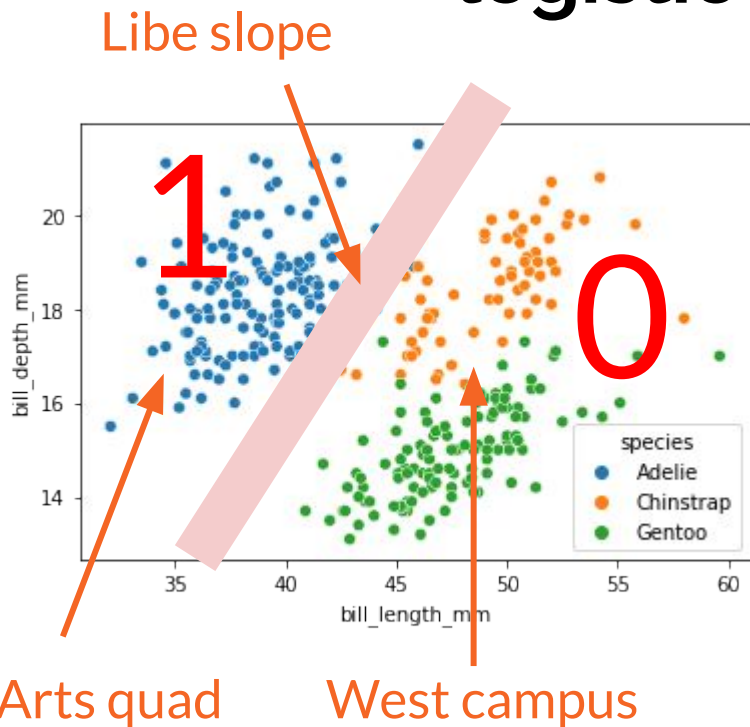
Math version:

$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

The logistic function squashes negative towards 0 and positive towards 1

logistic regression



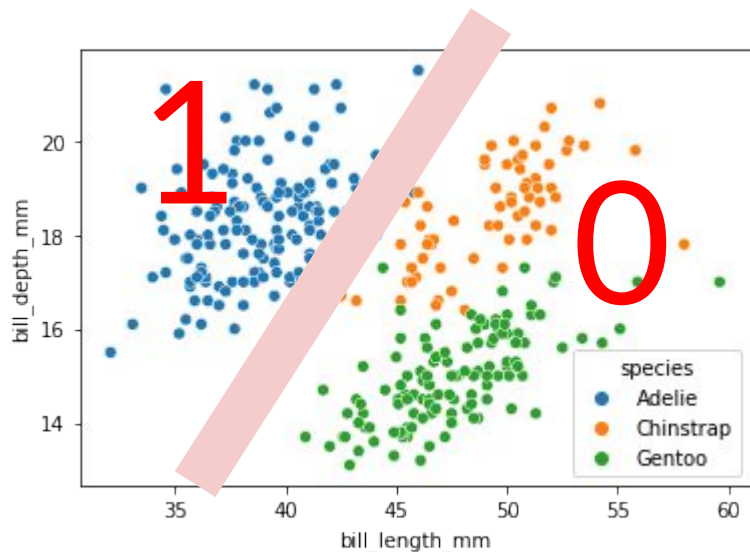
Math version:

$$\alpha = 25$$

$$\beta_{\text{length}} = -1.5, \beta_{\text{depth}} = 2.0$$

The logistic function squashes negative towards 0 and positive towards 1

Three views of Log. Reg.



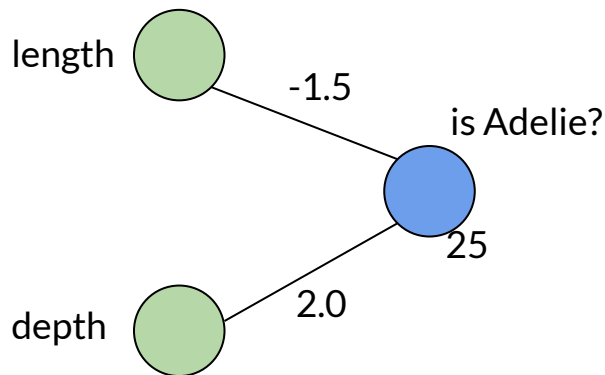
Mathematical expression

$$y = \sigma(\alpha + \beta_{\text{length}} X_{\text{length}} + \beta_{\text{depth}} X_{\text{depth}})$$

Python object

```
adelie_model.coef_, adelite_model.intercept_  
(array([[ -1.37,   1.98]]), array([24.29]))
```

A visual language for regression

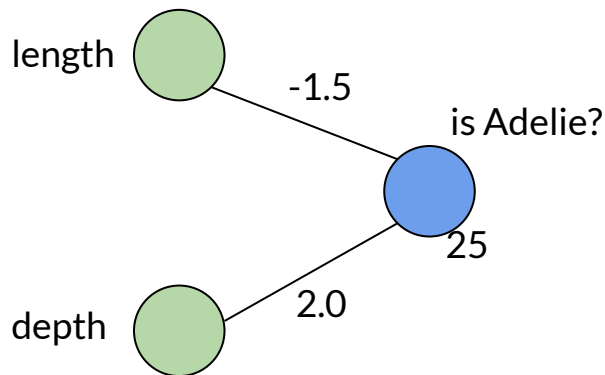


The *inputs* (length and depth) feed into the output (is the penguin Adelie?)

Each input has a *weight*

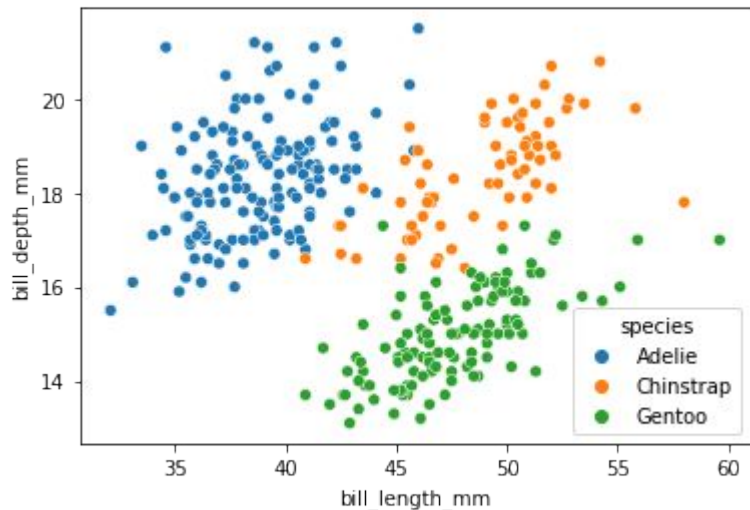
The output has a *bias* or *intercept*

A visual language for regression



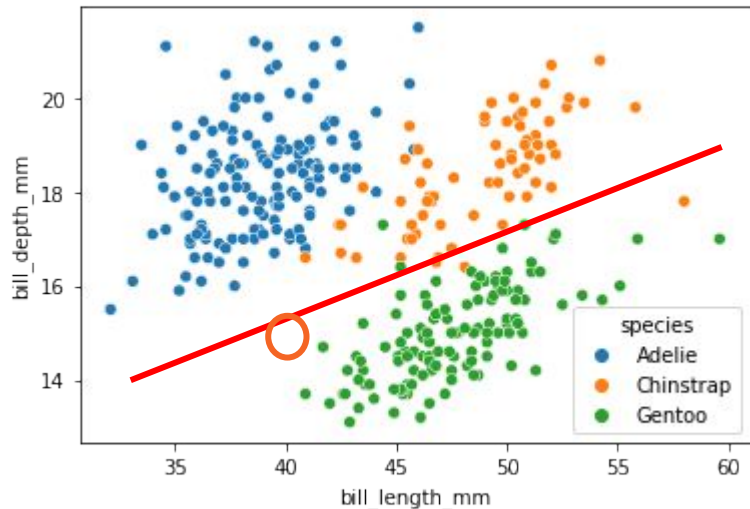
"To guess if a penguin is an Adelie, take -1.5 times the bill length plus 2.0 times the bill depth, and add 25. If the result is positive, guess yes (1), otherwise guess no (0)"

Linear separability



Given bill length and depth, can you draw a line that reliably separates **Gentoo** penguins from **Chinstrap** and **Adelie** penguins?

Linear separability



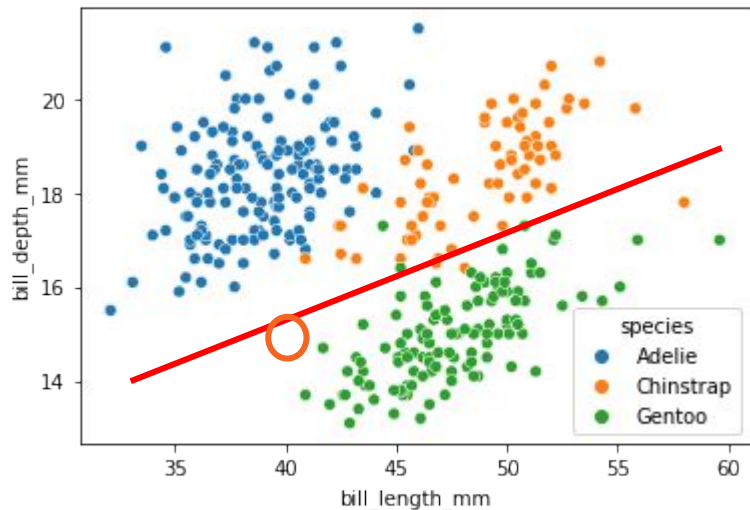
Math version:

$$\alpha = 28$$

$$\beta_{\text{length}} = 0.5, \beta_{\text{depth}} = -3.0$$

Is the value for length=40, depth=15
positive or negative?

Linear separability



Math version:

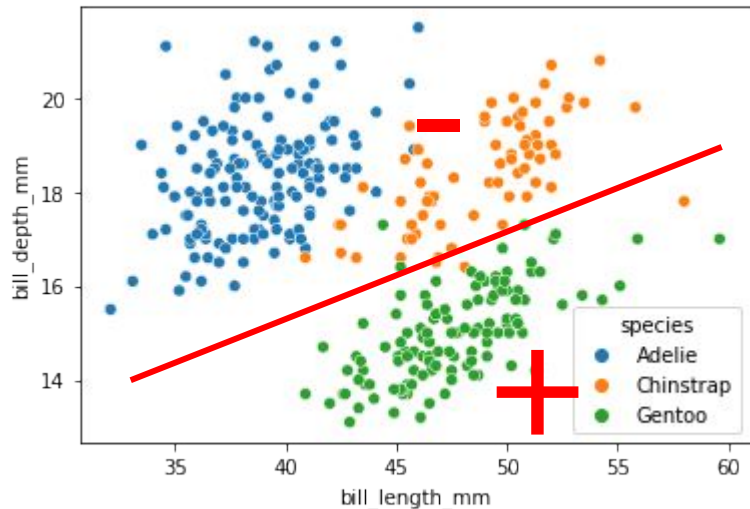
$$\alpha = 28$$

$$\beta_{\text{length}} = 0.5, \beta_{\text{depth}} = -3.0$$

Is the value for length=40, depth=15
positive or negative?

$$28 + 20 - 45 = 3 \text{ *positive*}$$

Regions of positive and negative



Math version:

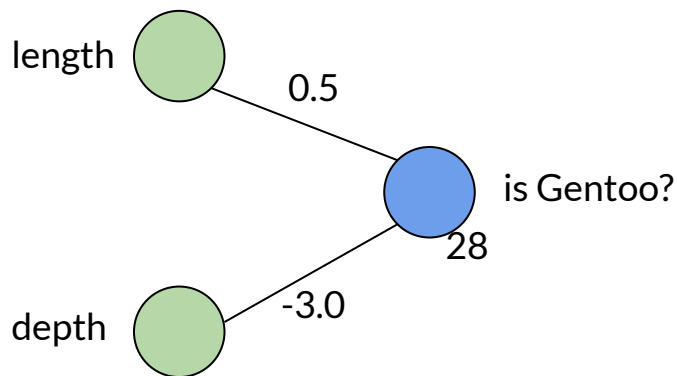
$$\alpha = 28$$

$$\beta_{\text{length}} = 0.5, \beta_{\text{depth}} = -3.0$$

Is the value for length=40, depth=15
positive or negative?

$$28 + 20 - 45 = 3 \text{ positive}$$

A visual language for regression

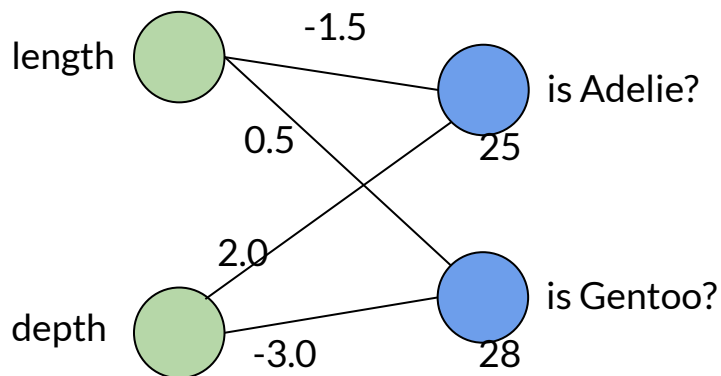


The *inputs* (length and depth) feed into the output (is the penguin Adelie?)

Each input has a *weight*

The output has a *bias* or *intercept*

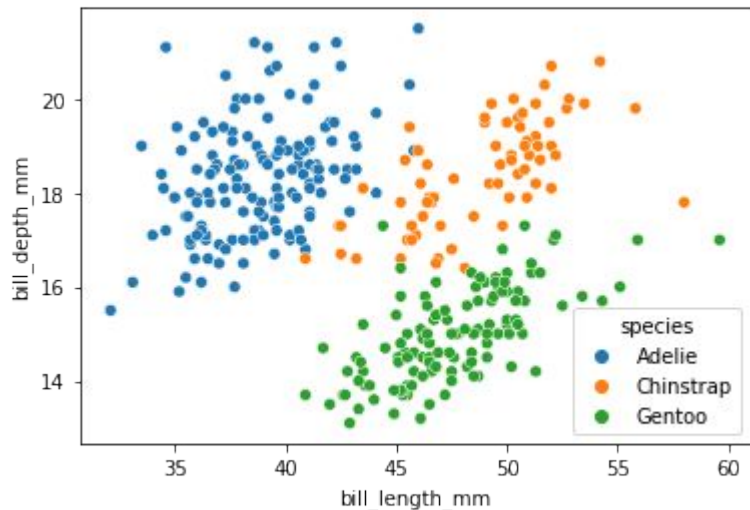
A visual language for regression



We can create more than one regression model from the same data!

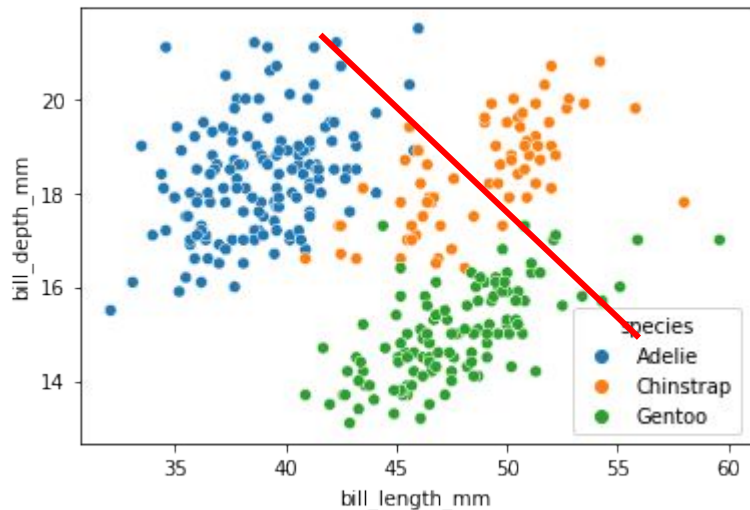
This is equivalent to drawing two decision boundaries

Linear separability



Given bill length and depth, can you draw a line that reliably separates **Chinstrap** penguins from **Gentoo** and **Adelie** penguins?

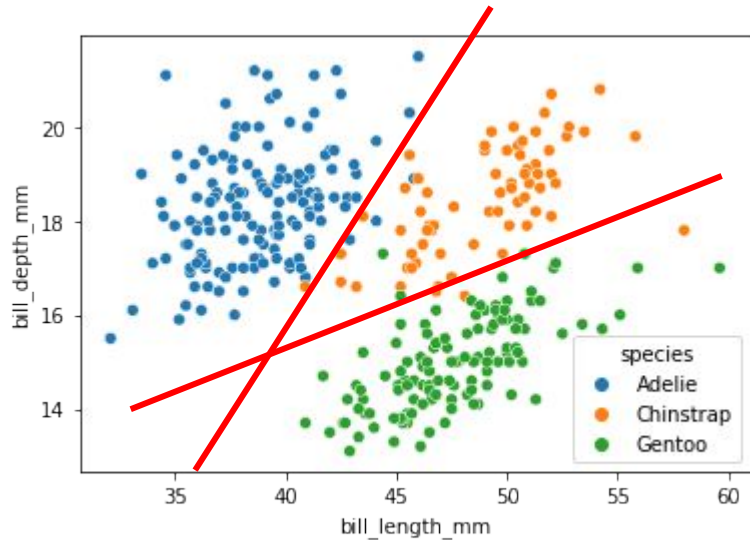
Linear separability



Given bill length and depth, can you draw a line that reliably separates **Chinstrap** penguins from **Gentoo** and **Adelie** penguins?

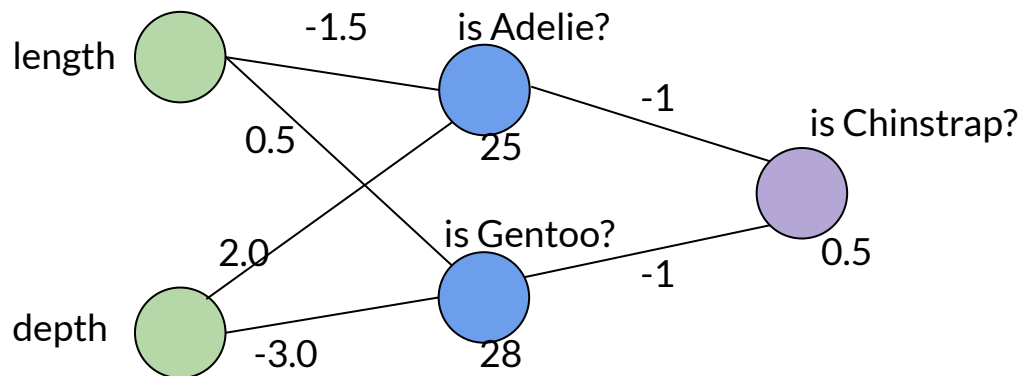
Not really, no

Outputs as inputs



If we combine the output of two linear classifiers, we can identify a non-linear region!

Multiple layers of regression



Outputs from layer 1
become inputs to layer 2

Classifying

- To classify data (e.g., to predict whether binary output y is 0 or 1), we've learned how to use:
 - (Logistic) regression
 - Naive Bayes
 - K-means clustering

Classifying

- To classify data (e.g., to predict whether binary output y is 0 or 1), we've learned how to use:
 - (Logistic) regression
 - Naive Bayes
 - K-means clustering
- What if we have *continuous* output y ?

Classifying

- To classify data (e.g., to predict whether binary output y is 0 or 1), we've learned how to use:
 - (Logistic) regression
 - Naive Bayes
 - K-means clustering
- What if we have *continuous* output y ? We can use **linear regression!**

Predicting continuous data

- We know we can do the following things with linear regressions:
 - **Predict** (continuous) outputs
 - **Summarize** relationships between input and output variables
 - Describe outliers/oddities

Predicting continuous data

- We know we can do the following things with linear regressions:

What if we can get more
accurate predictions



- **Predict** (continuous) outputs
- **Summarize** relationships between input and output variables
- Describe outliers/oddities

Predicting continuous data

- We know we can do the following things with linear regressions:

What if we can get more
accurate predictions



- **Predict** (continuous) outputs

In exchange for less
interpretable results?



- **Summarize** relationships between input and output variables
- Describe outliers/oddities

Neural networks!

- Neural nets are **models**:
 - Still have multiple inputs x_i and make a single prediction \hat{y}
 - Still train on *train_set* and test on *test_set*

Neural networks!

- Neural nets are **models**:

Note: this prediction
can be classification
OR continuous

- Still have multiple inputs x_i and make a single prediction \hat{y}
- Still train on *train_set* and test on *test_set*

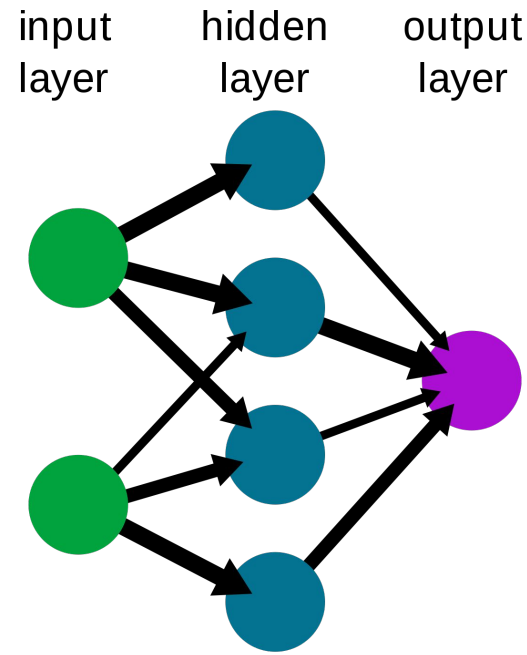
Neural networks!

- Neural nets are **models**:
 - Still have multiple inputs x_i and make a single prediction \hat{y}
 - Still train on *train_set* and test on *test_set*
- But, neural nets often outperform linear regressions on **big data** because they can account for **nonlinearities**

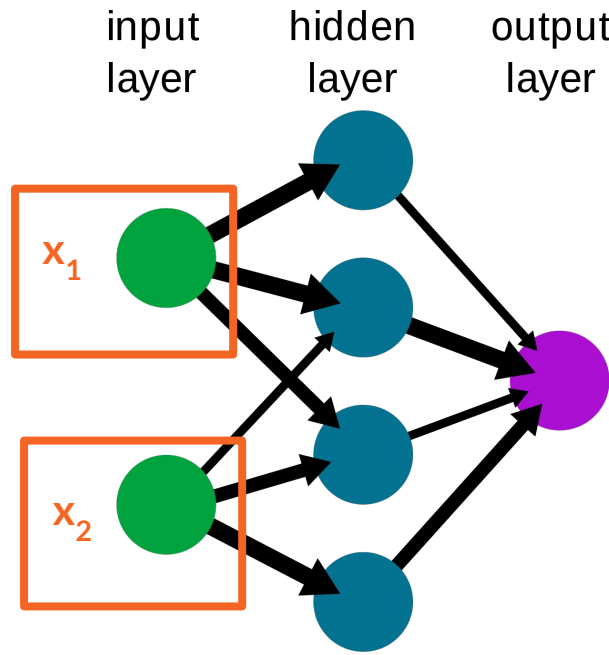
Neural nets work
pretty well, but
they aren't
perfect!



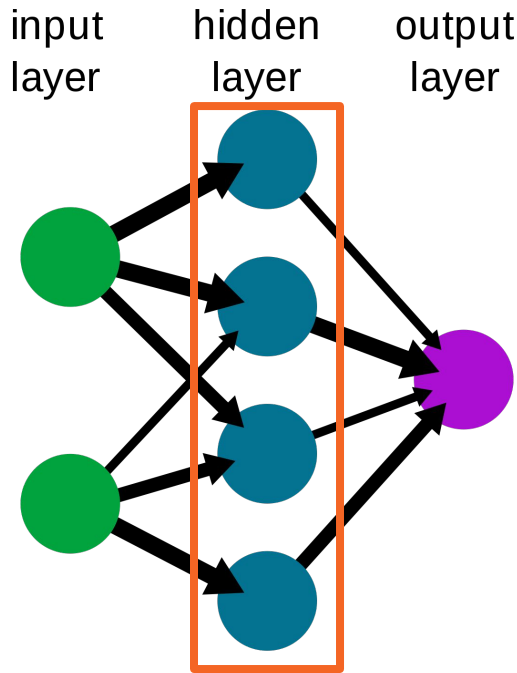
A simple neural net (NN)



A simple neural net (NN)

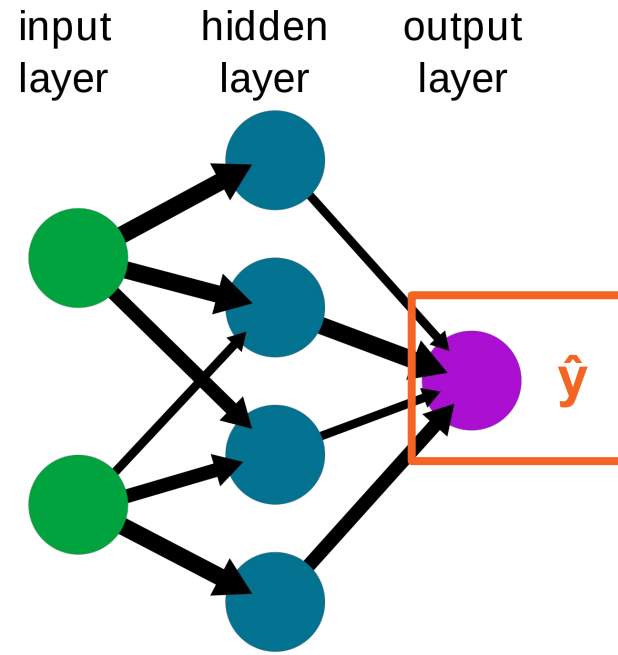


A simple neural net (NN)

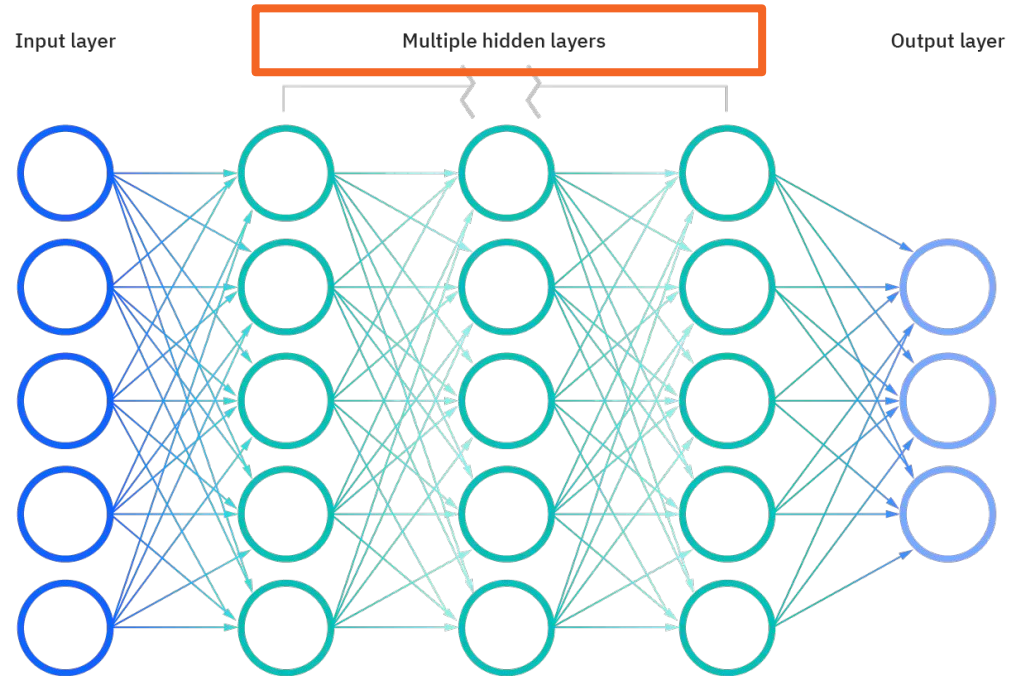


Each circle is a “node” with an input and an output

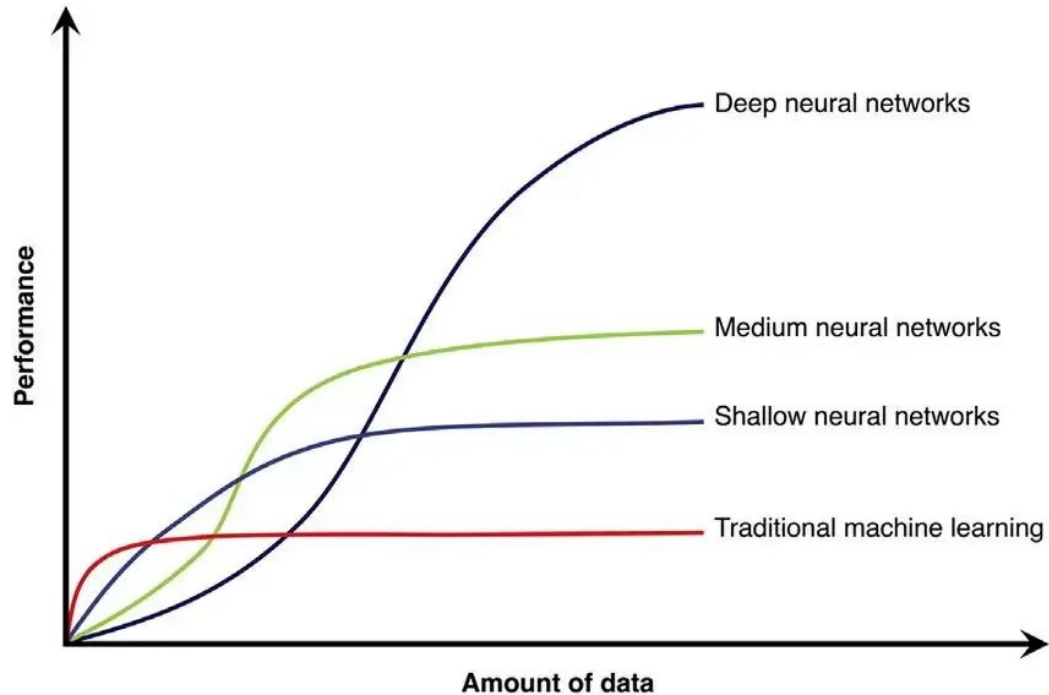
A simple neural net (NN)



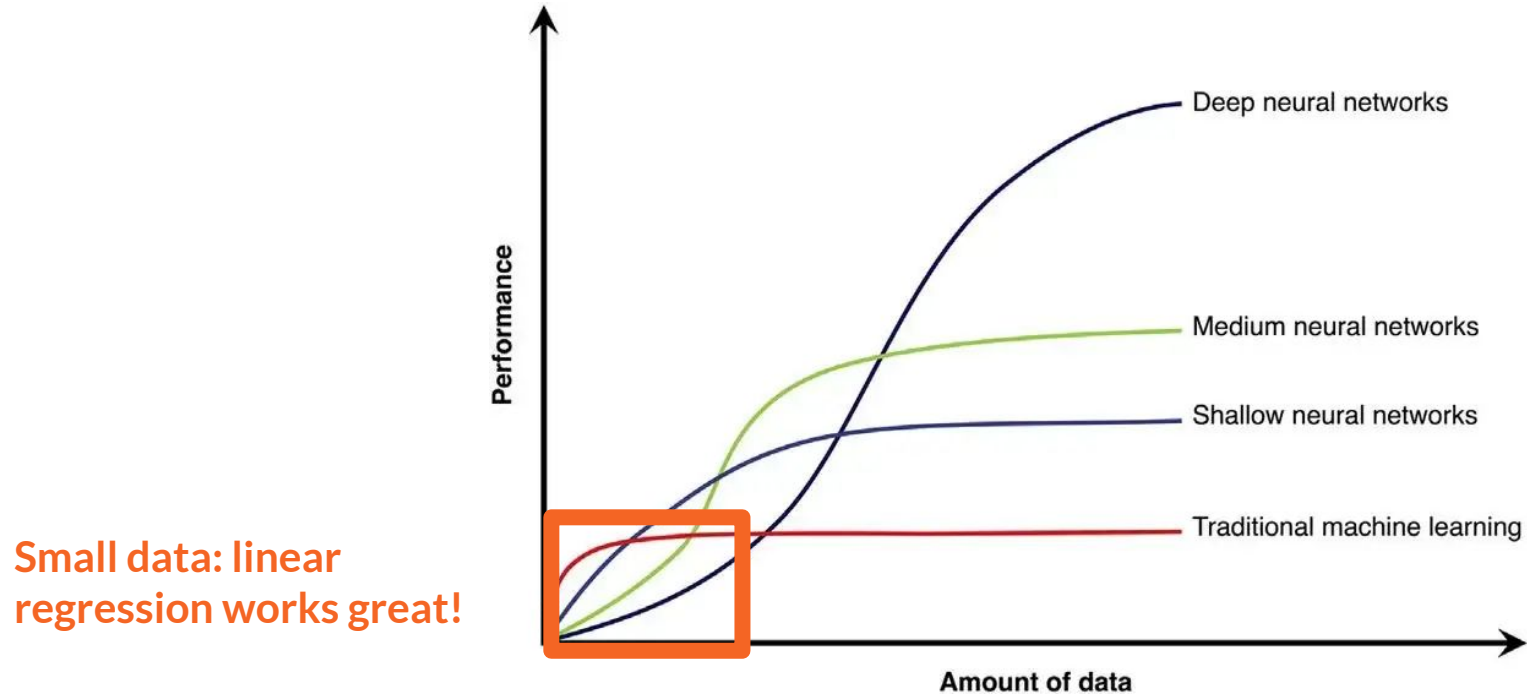
A **deep** neural net (NN)



Neural networks for big data

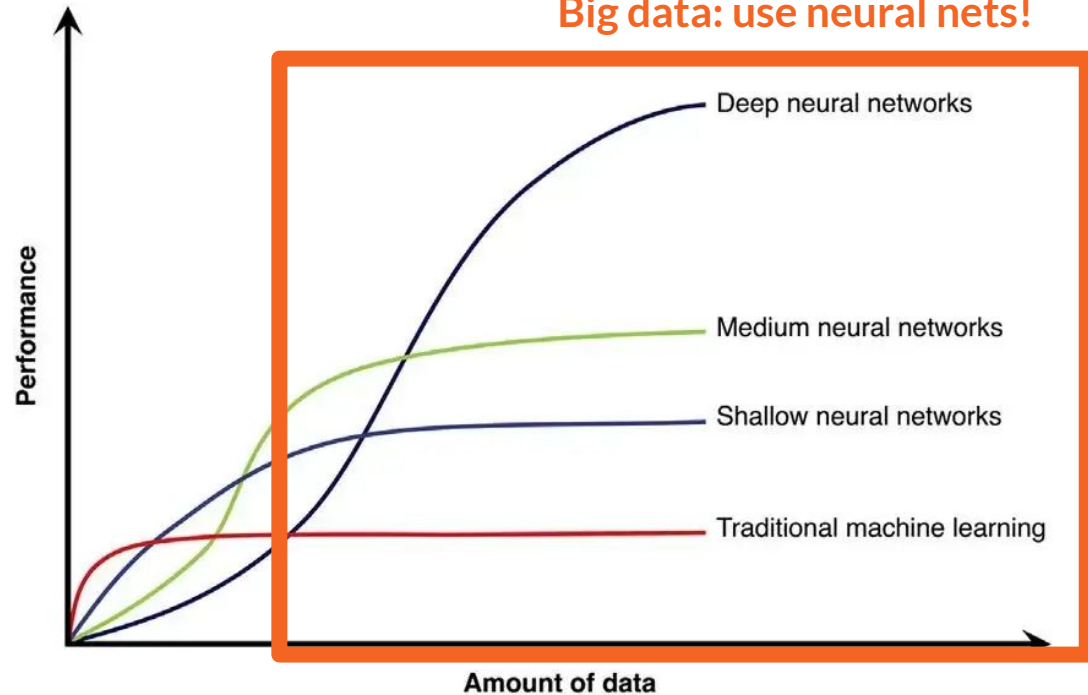


Neural networks for big data

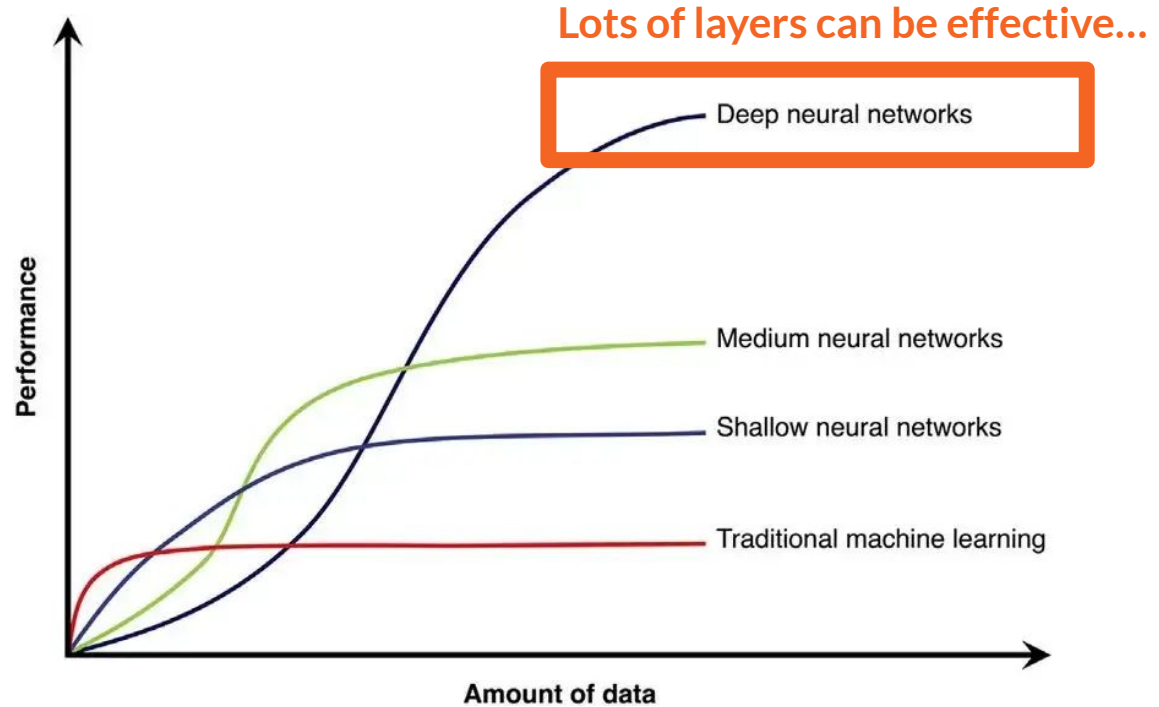


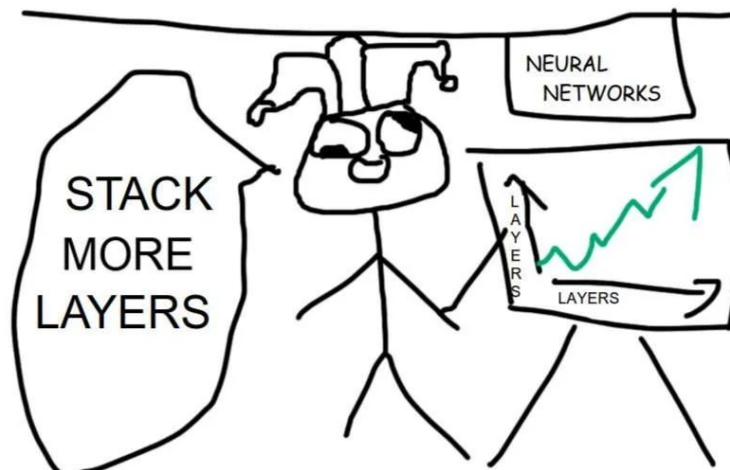
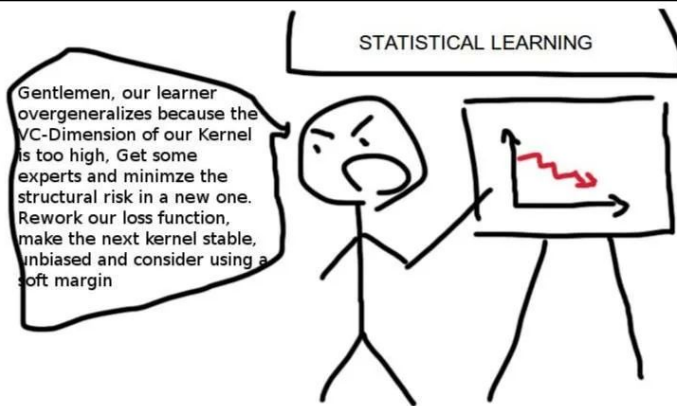
Neural networks for big data

Big data: use neural nets!



Neural networks for big data





Neural net drawbacks

- Neural nets are difficult for humans to interpret
- Neural nets are prone to overfitting on small data (needs big data to benefit from nonlinearities)
- Neural nets require a lot of compute
 - GPUs, CPUs, distributed computing (work in parallel)

Compute & sustainability

On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?

Emily M. Bender*
ebender@uw.edu
University of Washington
Seattle, WA, USA

Angelina McMillan-Major
aymm@uw.edu
University of Washington
Seattle, WA, USA

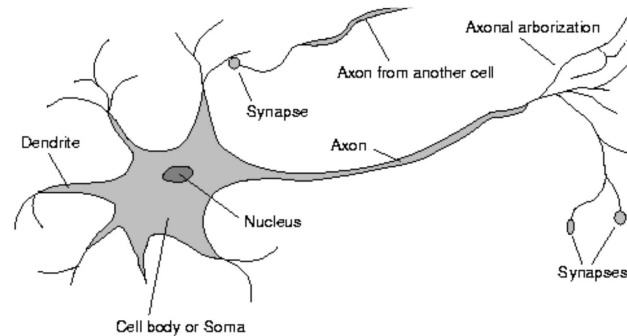
Timnit Gebru*
timnit@blackinai.org
Black in AI
Palo Alto, CA, USA

Shmargaret Shmitchell
shmargaret.shmitchell@gmail.com
The Aether

- “Training a single BERT base model (without hyperparameter tuning) on GPUs was estimated to require **as much energy as a trans-American flight.**”

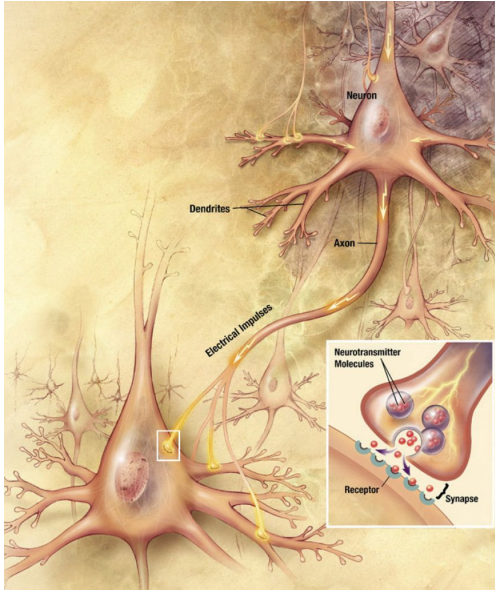
What is a neural net?

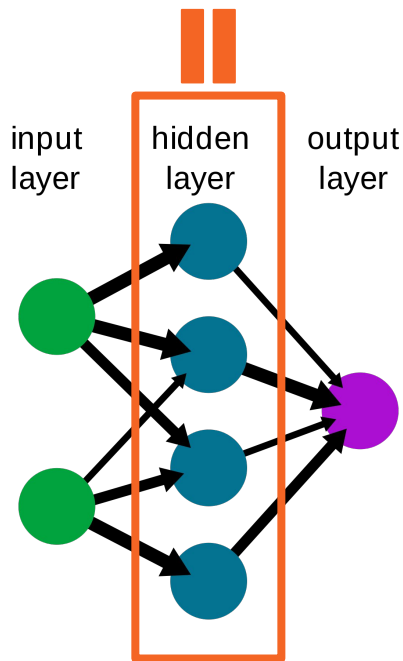
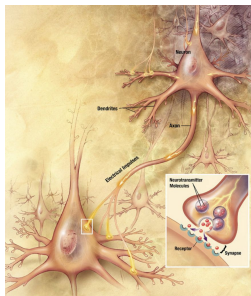
- Inspired by the human brain:
 - Different parts of the brain have different functions, which occur by firing neurons (10^{11})



What is a neural net?

- Inspired by the human brain:
 - Different parts of the brain have different functions, which occur by firing neurons (10^{11})
 - We don't exactly know how different brain functions are assigned (hard to interpret!), but we know neurons connect via synapses (10^{14})





What is a neural net?

- Inspired by the human brain:
 - Different parts of the brain have different functions, which occur by firing neurons (10^{11})
 - We don't exactly know how different brain functions are assigned (hard to interpret!), but we know neurons connect via synapses (10^{14})

What is a neural net?

- Inspired by the human brain:
 - Different parts of the brain have different functions, which occur by firing neurons (10^{11})
 - We don't exactly know how different brain functions are assigned (hard to interpret!), but we know neurons connect via synapses (10^{14})
 - The output is “consciousness”

What is a neural net?

- Inspired by the human brain:
 - Different parts of the brain have different functions, which occur by firing neurons (10^{11})
 - We don't exactly know how different brain functions are assigned (hard to interpret!), but we know neurons connect via synapses (10^{14})
 - The output is “consciousness”

Computers: “only”
 10^9 gates, 10^9 bits
RAM. Brain wins!

Example: predicting house prices

- What are some input x 's that are useful to predict output $y = \text{house price}$?

Example: predicting house prices

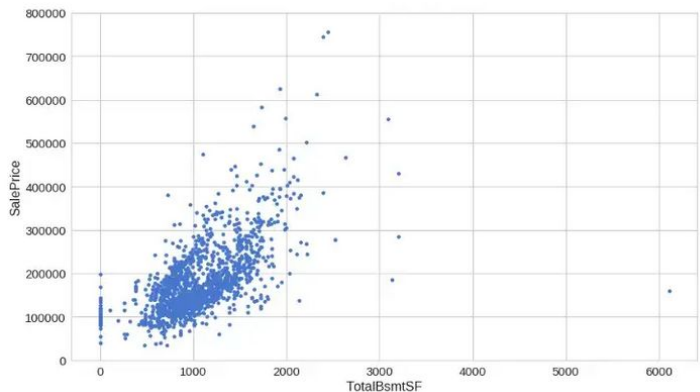
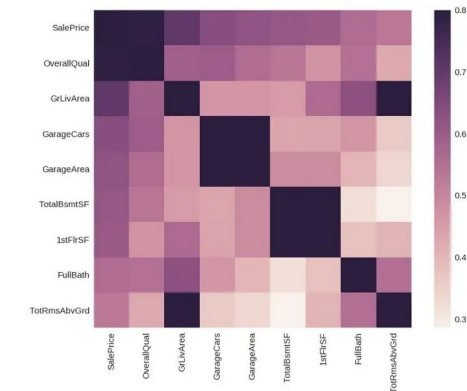
- What are some input x 's that are useful to predict output y = house price?
 - Square footage
 - # bed, # bath
 - Location
 - Neighborhood house \$'s
 - ...

Refresher: linear regression

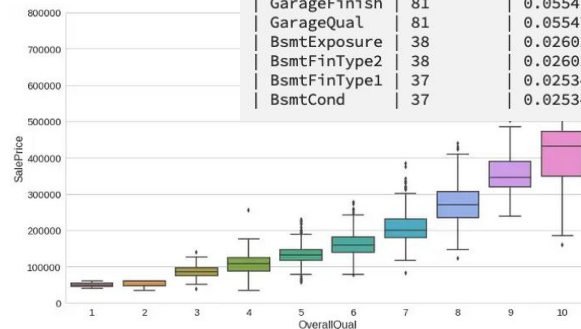
- Step 1: decide on your inputs and outputs
 - Plot/summarize...
 - Outliers
 - Collinearity
 - Missing data

Refresher: linear regression

- Step 1: decide on your inputs and outputs
 - Plot/summarize...
 - Outliers
 - Collinearity
 - Missing data



	Row count	Percentage
PoolQC	1453	0.995205
MiscFeature	1406	0.963014
Alley	1369	0.937671
Fence	1179	0.807534
FireplaceQu	690	0.472603
LotFrontage	259	0.177397
GarageCond	81	0.055479
GarageType	81	0.055479
GarageYrBlt	81	0.055479
GarageFinish	81	0.055479
GarageQual	81	0.055479
BsmtExposure	38	0.026027
BsmtFinType2	38	0.026027
BsmtFinType1	37	0.025342
BsmtCond	37	0.025342

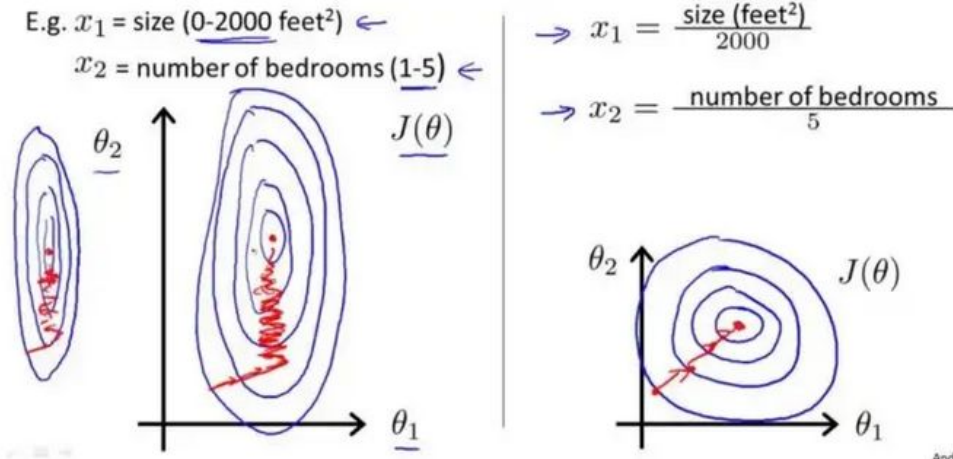


Refresher: linear regression

- Step 2: perform data preprocessing
 - Missing values (imputation)
 - Transformations, normalization

Refresher: linear regression

- Step 2: perform data preprocessing
 - Missing values (imputation)
 - Transformations, normalization

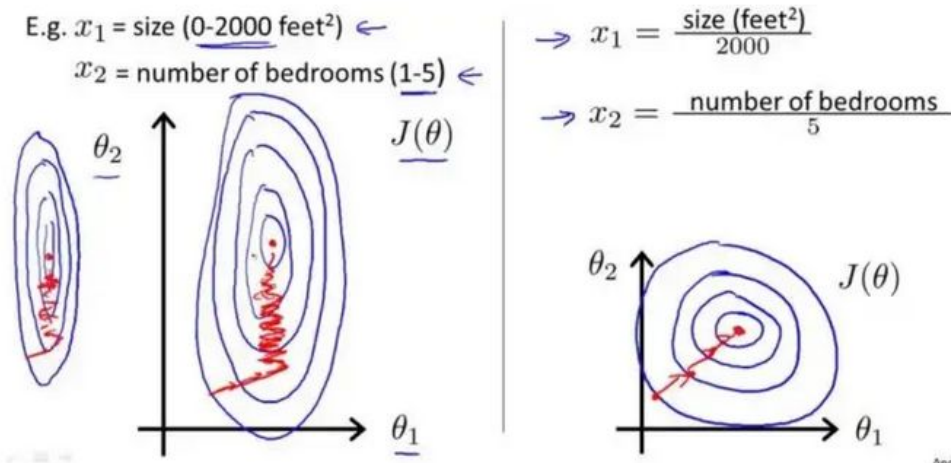


Andrew Ng

Refresher: linear regression

- Step 2: perform data preprocessing
 - Missing values (imputation)
 - Transformations, normalization

In addition to having more easily comparable coefficients, our algorithms (SGD) converge faster when data is scaled!



Andrew Ng

Refresher: linear regression

- Step 3: decide on an evaluation metric
 - Choose one based on data type of output y
 - For housing prices: _____

Refresher: linear regression

- Step 3: decide on an evaluation metric
 - Choose one based on data type of output y
 - For housing prices: **RMSE, MSE, MAE, MAPE, ... any evaluation metric that is reasonable for continuous output y**

Refresher: linear regression

- Step 3: decide on an evaluation metric
 - Choose one based on data type of output y
 - For housing prices, maybe RMSE
 - `def rmse(predict, actual) :`
`return np.sqrt(np.mean(np.square(predict - actual)))`

Refresher: linear regression

- Step 3: decide on an evaluation metric
 - Choose one based on data type of output y
 - For housing prices, maybe RMSE
- Step 4: split your data
 - Train / val / test sets
 - `X_train, X_test, y_train, y_test`
`=train_test_split(X, y, test_size=0.4)`

Refresher: linear regression

- Step 5: run a regression

- $y \sim x_1 + x_2 + x_3$

$$Y(x_1, x_2, x_3) = w_1x_1 + w_2x_2 + w_3x_3 + w_0$$

Different notation: instead of β 's for coefficients, we now use w to represent “weights”

Refresher: linear regression

- Step 5: run a regression

- $y \sim x_1 + x_2 + x_3 + \dots$

```
model=LinearRegression().fit(X_train,y_train)
predictions_test=model.predict(X_test)
rmse(predictions_test,y_test)
```

Refresher: linear regression

- Step 5: run a regression

- $y \sim x_1 + x_2 + x_3 + \dots$

```
model=LinearRegression().fit(X_train,y_train)
predictions_test=model.predict(X_test)
rmse(predictions_test,y_test)
```

- Step 6: interpret results!
 - Predict, summarize, outliers/oddities

Model Recap

- Step 1: decide on inputs / outputs
- Step 2: data preprocessing
- Step 3: decide on an evaluation metric
- Step 4: split your data
- Step 5: run the model
- Step 6: interpret results

But what about neural nets?!

- A lot of these steps are the same!

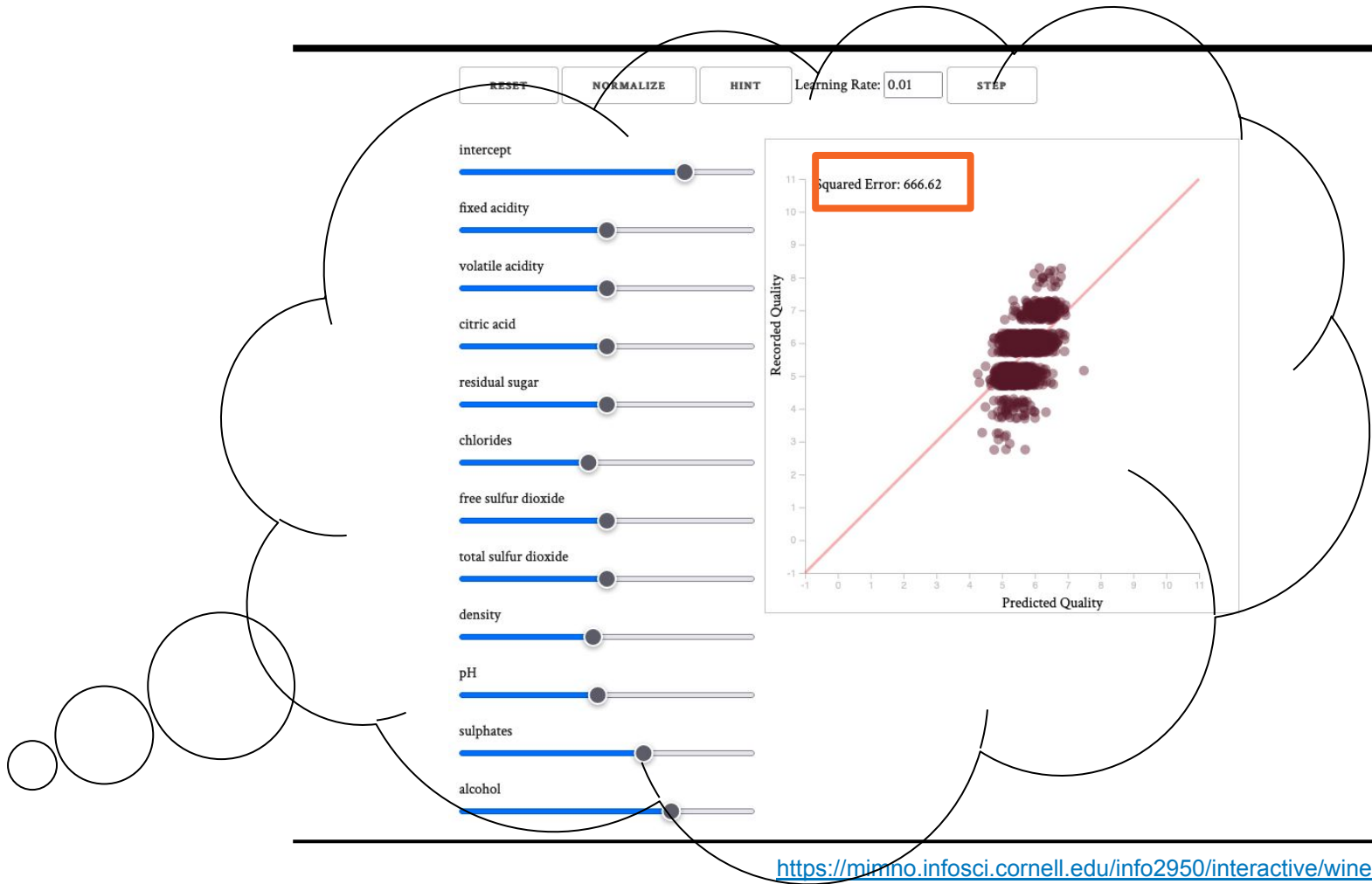
We do this before even
touching linear regression
/ neural nets!

- Step 1: decide on inputs / outputs
- Step 2: data preprocessing
- Step 3: decide on an evaluation metric
- Step 4: split your data
- Step 5: run the model
- Step 6: interpret results

But what about neural nets?!

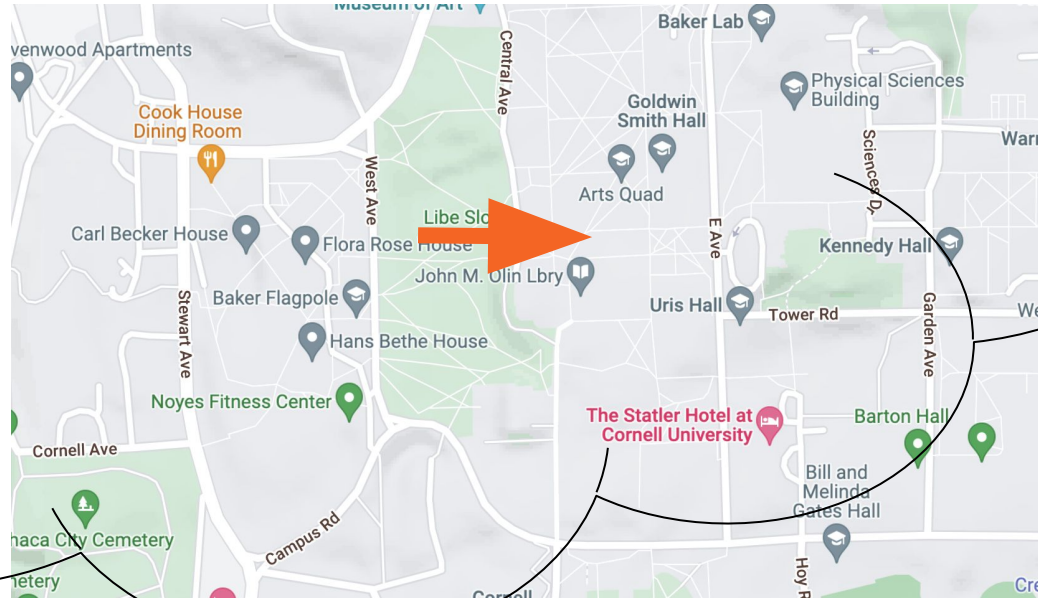
- A lot of these steps are the same!
 - Step 1: decide on inputs / outputs
 - Step 2: data preprocessing
 - Step 3: decide on an evaluation metric
 - Step 4: split your data
 - Step 5: run the model
 - Step 6: interpret results

Caveat: more explanation
needed here



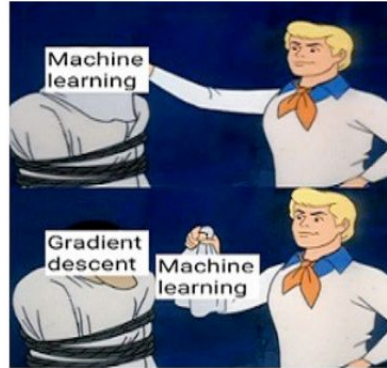
Gradient is a hint in multiple directions

The slope is steep in the
East-West direction,
but flat in the
North-South direction



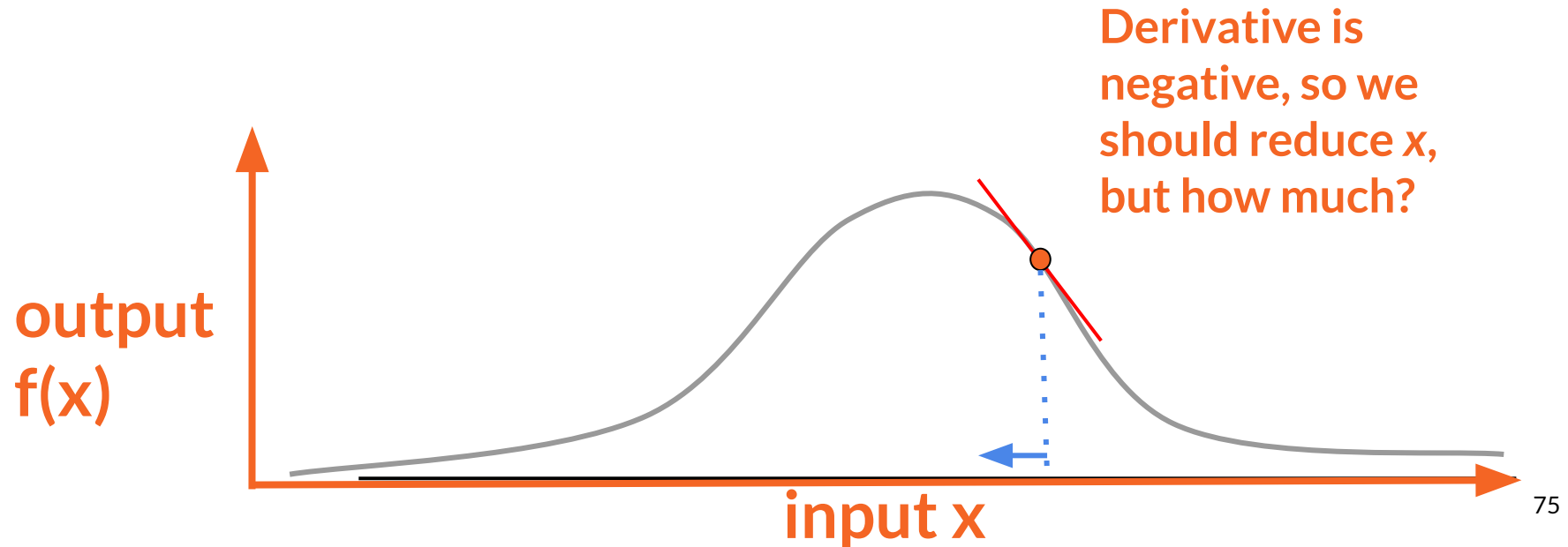
Takeaways on gradients

- “Stochastic gradient descent” is used to find minima / maxima for complicated models (e.g. multivariable regression)
 - choose a **learning** rate to do this efficiently
 - this is the core of modern machine **learning**!

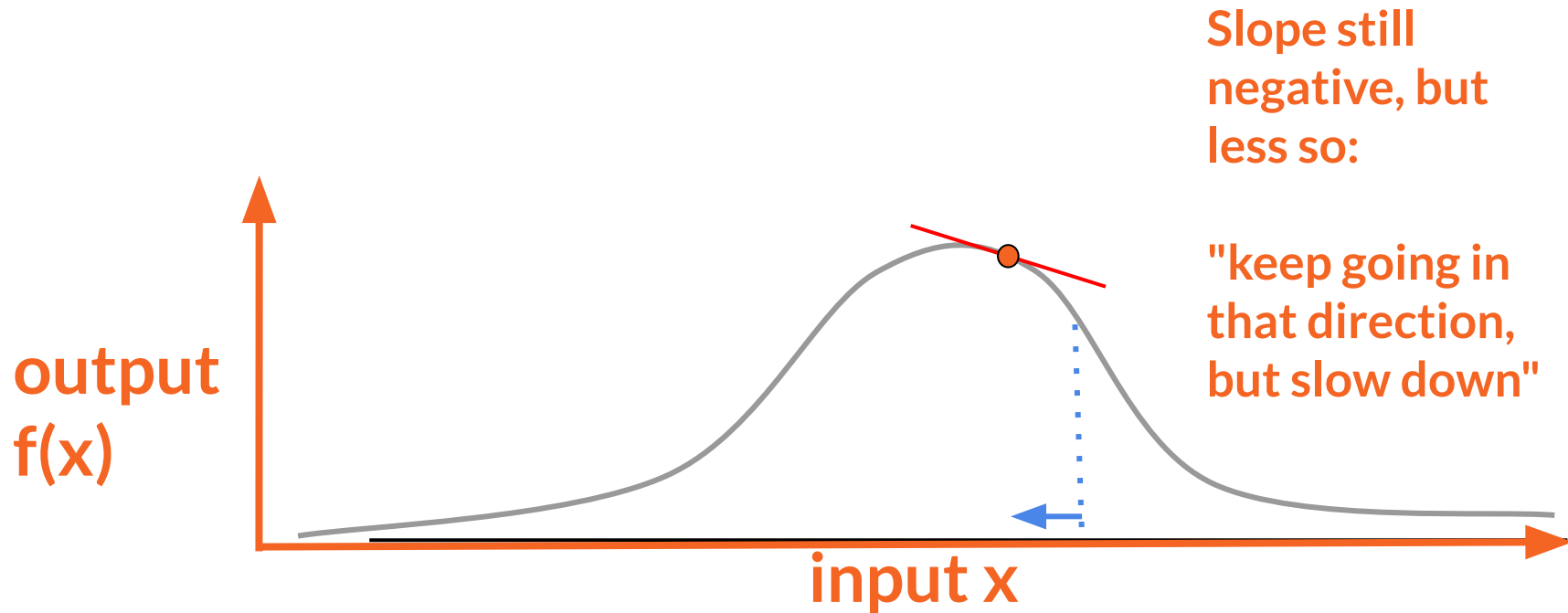


<https://towardsdatascience.com/batch-mini-batch-stochastic-gradient-descent-7a62ecba642a>

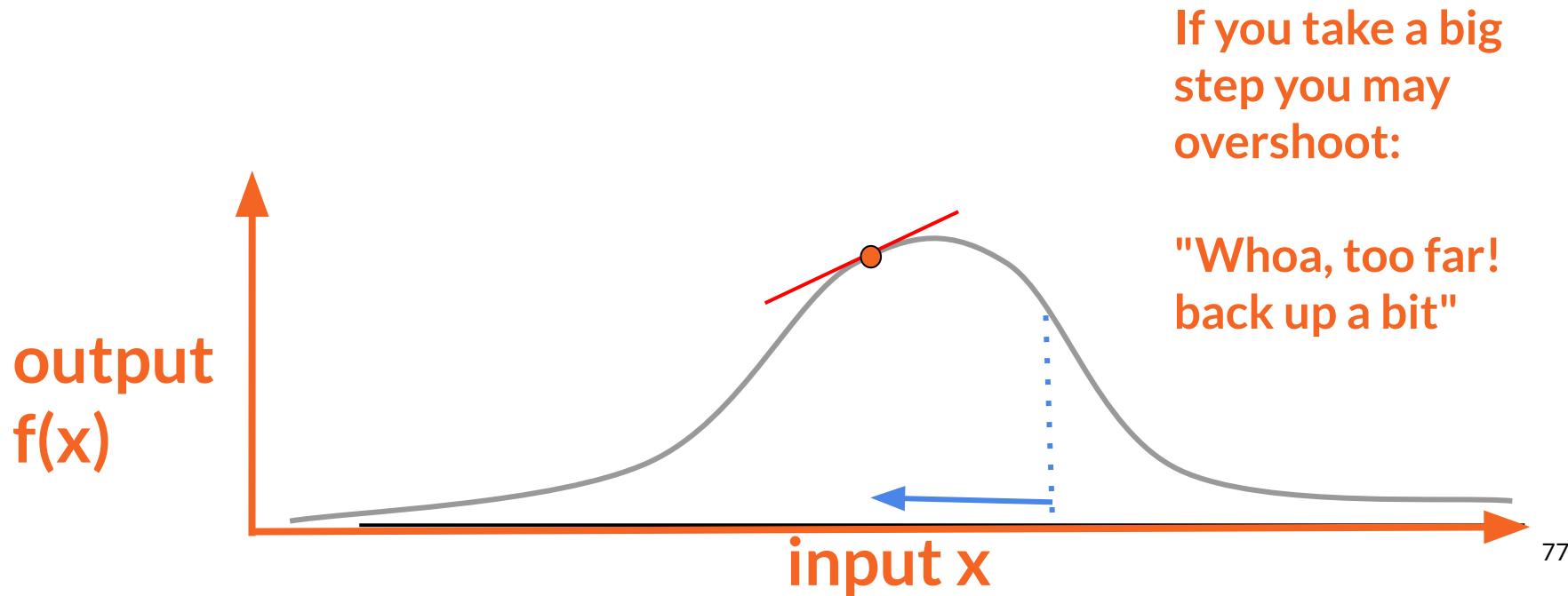
Stochastic Gradient Descent (SGD)



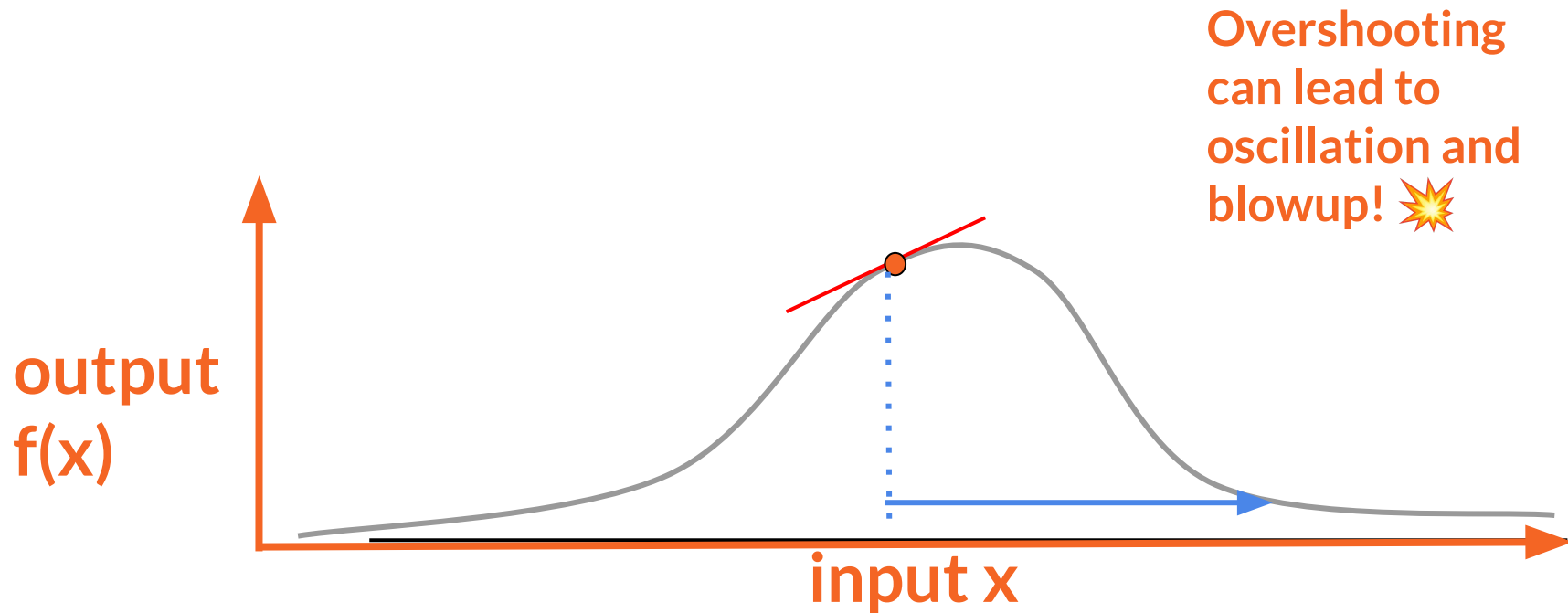
SGD: move β in a direction specified by the gradient



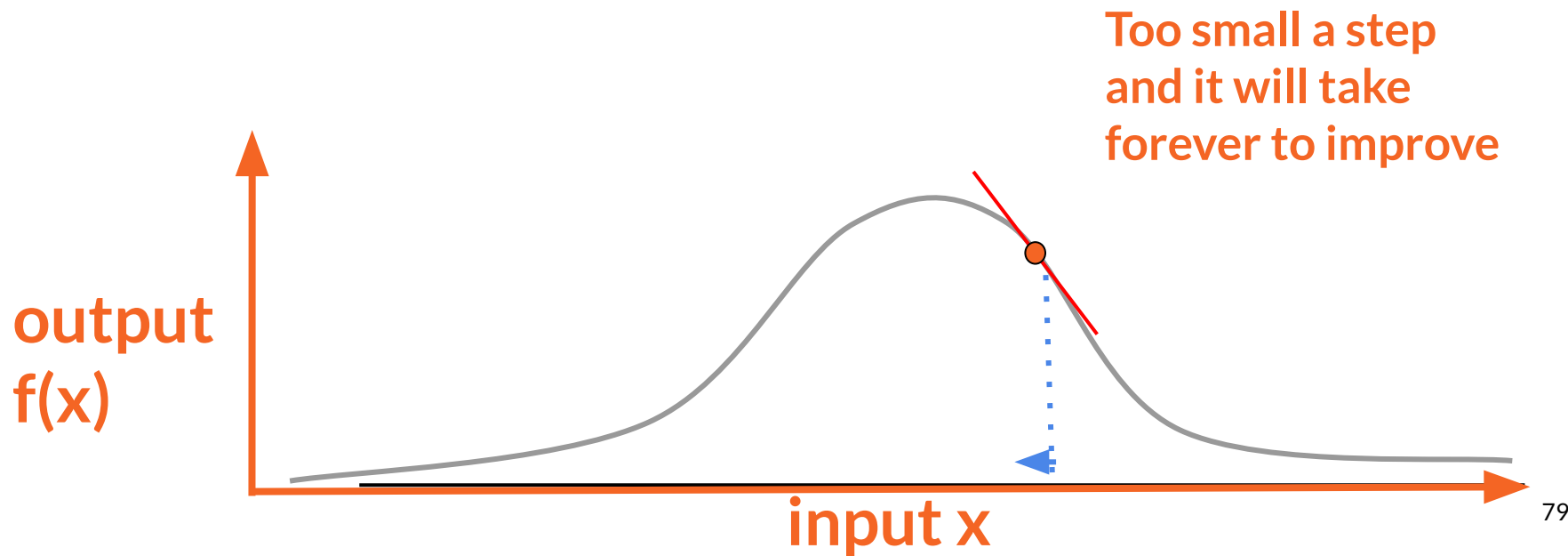
SGD: "step size" or "learning rate" is important



SGD: "step size" or "learning rate" is important



SGD: "step size" or "learning rate" is important

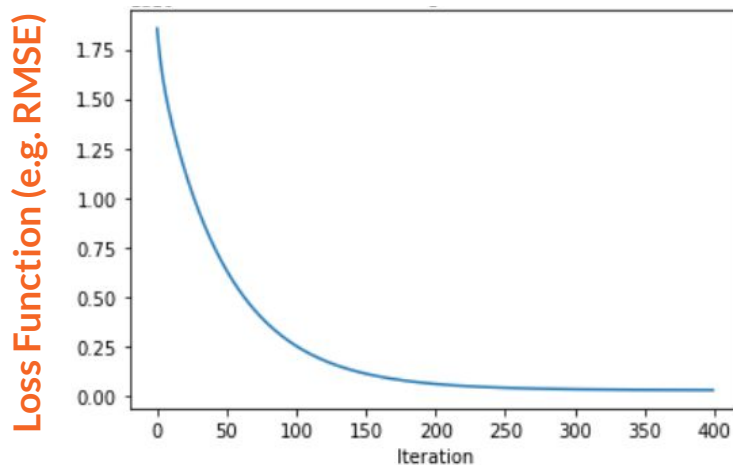


Step 3: decide on an evaluation metric

- In machine learning, we use an algorithm (e.g. SGD) that steps towards convergence and “learns” when we’ve reached a local minimum
 - “Loss function” = how you evaluate when to stop
 - RMSE is one choice of loss function

Step 3: decide on an evaluation metric

- In machine learning, we use an algorithm (e.g. SGD) that steps towards convergence and “learns” when we’ve reached a local minimum



Step 3: decide on an evaluation metric

- In Step 5, we've used sklearn to run:

```
model=LinearRegression().fit(X_train,y_train)
```

- But, we can re-define “fit” to include things like...


```
def fit(self, X, y, n_iter=100000, lr=0.01):  
    # [gradient descent code that runs for n_iter iterations and  
    # uses lr learning rate step sizes]  
    return self
```

Step 3: decide on an evaluation metric

- In Step 5, we've used sklearn to run:

```
model=LinearRegression().fit(X_train,y_train)
```

- But, we can re-define “fit” to include things like...

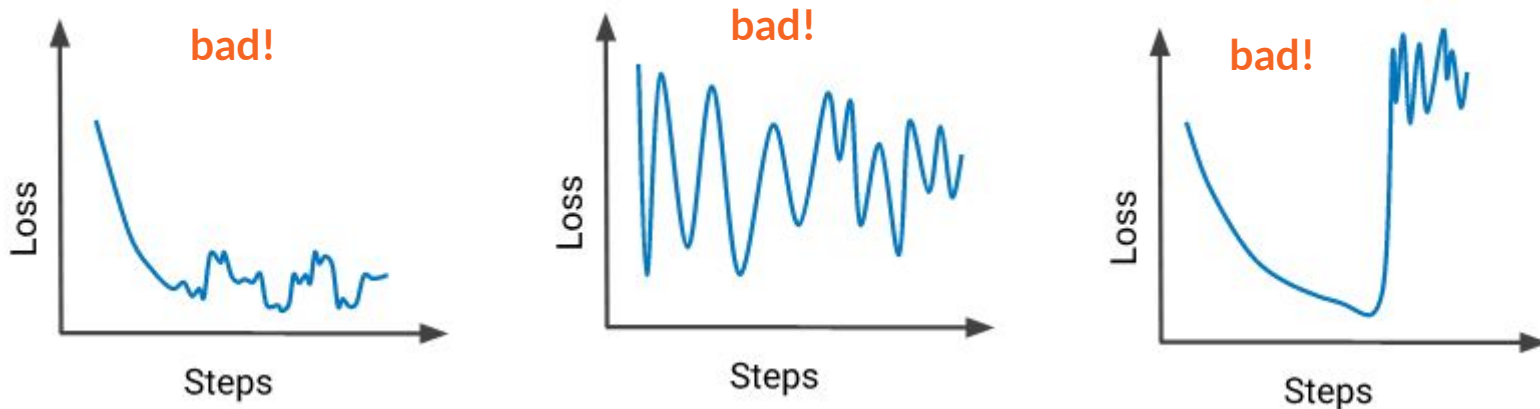


```
def fit(self, X, y, n_iter=100000, lr=0.01):  
    # [gradient descent code that runs for n_iter iterations and  
    # uses lr learning rate step sizes]  
    return self
```

```
model=LinearRegression().fit(X_train,y_train, 2000, 0.01)
```

Step 3: decide on an evaluation metric

- Your loss should be smoothly converging to a local min
- Debug by checking for NaNs or repeated input data, changing learning rate, etc. if you see plots like:

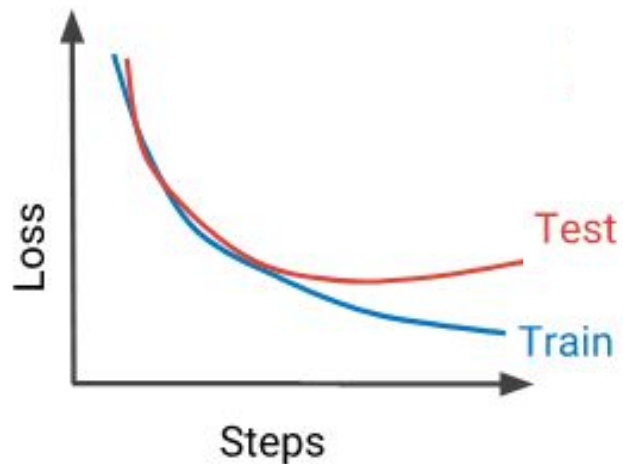


Step 3: decide on an evaluation metric

- It's important to check the loss for each of your train / val / test sets separately, to ensure no overfitting is happening!

Step 3: decide on an evaluation metric

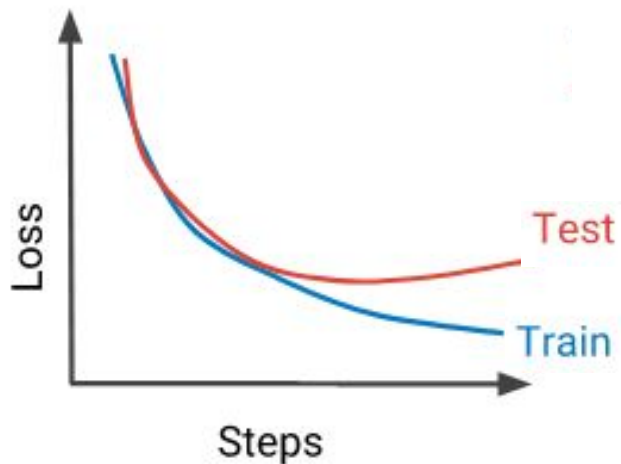
- It's important to check the loss for each of your train / val / test sets separately, to ensure no overfitting is happening!



Is this indicative
of overfitting?

Step 3: decide on an evaluation metric

- It's important to check the loss for each of your train / val / test sets separately, to ensure no overfitting is happening!



**This is overfit!
We're performing
way better on train
than test! (We want
to *minimize* loss)**

But what about neural nets?!

- A lot of these steps are the same!
 - Step 1: decide on inputs / outputs
 - Step 2: data preprocessing
 - Step 3: decide on an evaluation metric
 - Step 4: split your data
 - Step 5: run the model
 - Step 6: interpret results

Specifically, make sure
you have code to plot
loss function



But what about neural nets?!

- A lot of these steps are the same!
 - Step 1: decide on inputs / outputs
 - Step 2: data preprocessing
 - Step 3: decide on an evaluation metric
 - Step 4: split your data
 - Step 5: run the model
 - Step 6: interpret results

How do we do this with
neural nets?



1 min break & attendance!



tinyurl.com/7haktr95

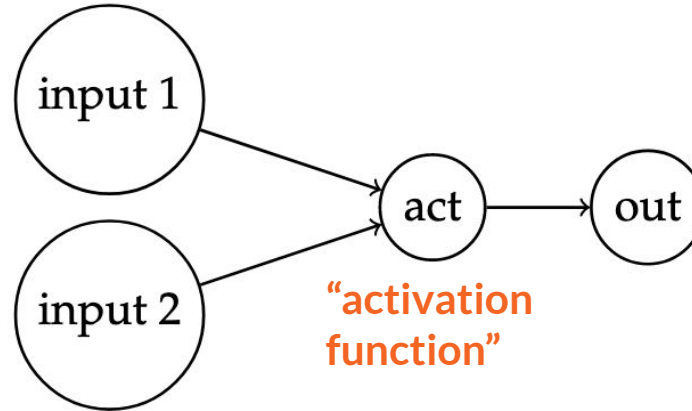
But what about neural nets?!

- A lot of these steps are the same!
 - Step 1: decide on inputs / outputs
 - Step 2: data preprocessing
 - Step 3: decide on an evaluation metric
 - Step 4: split your data
 - Step 5: run the model
 - Step 6: interpret results

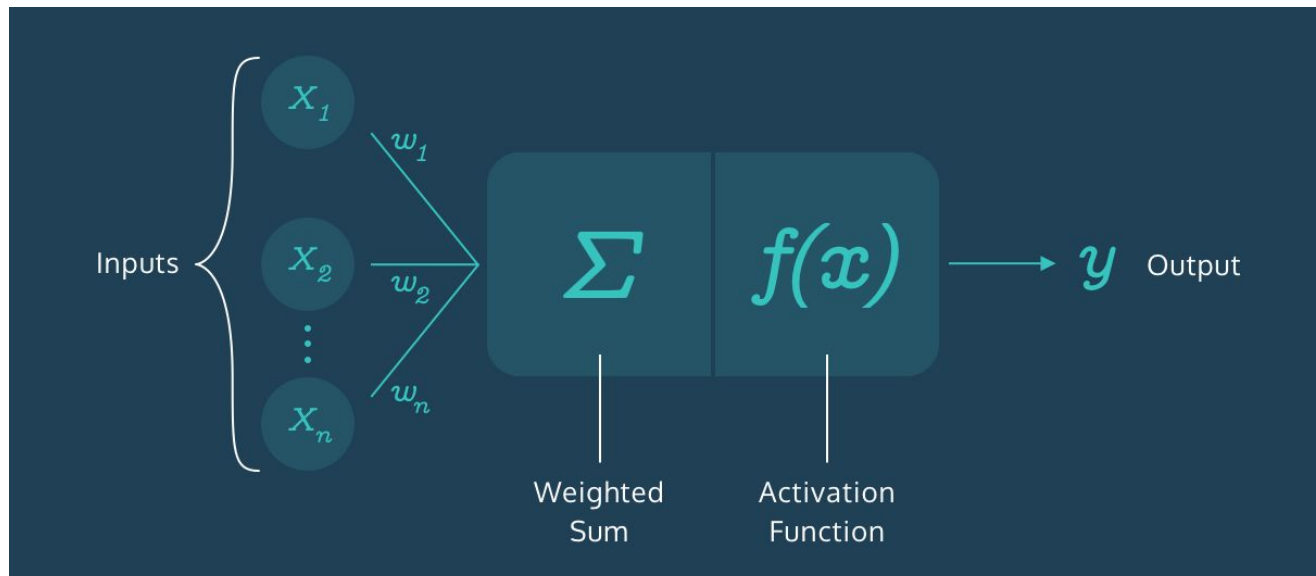
How do we do this with
neural nets?



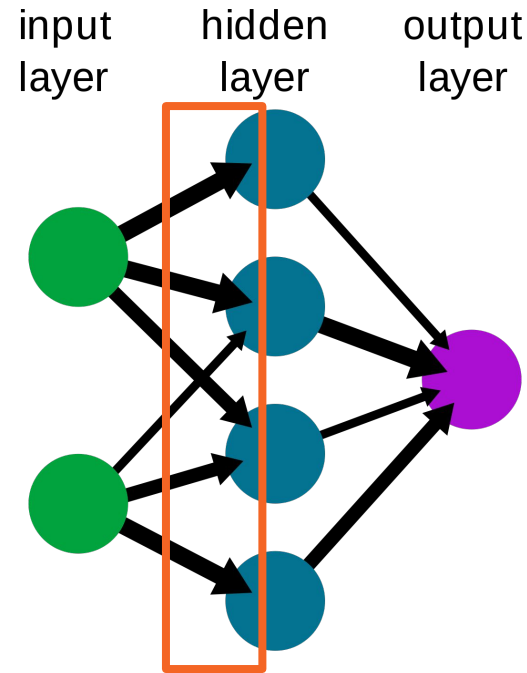
Perceptron



Perceptron (Minsky-Papert, 1969)

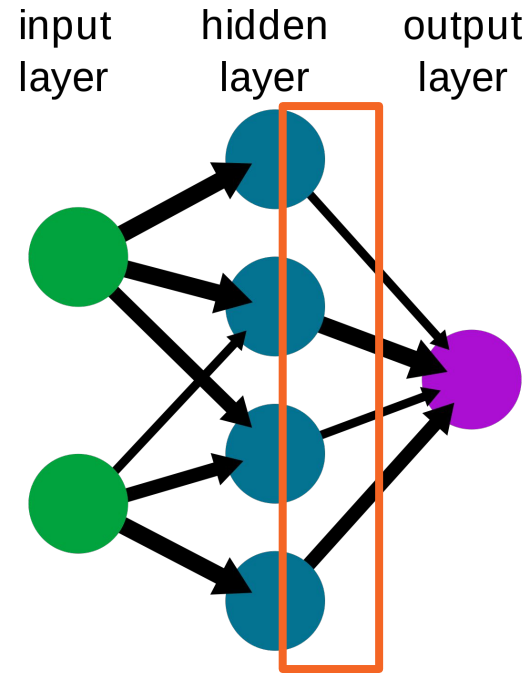


Single-layer perceptron



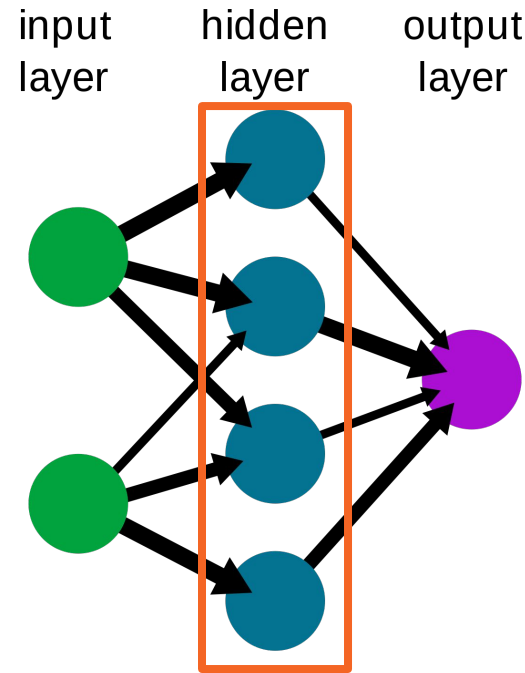
Each circle is a
“node” with an
input

Single-layer perceptron



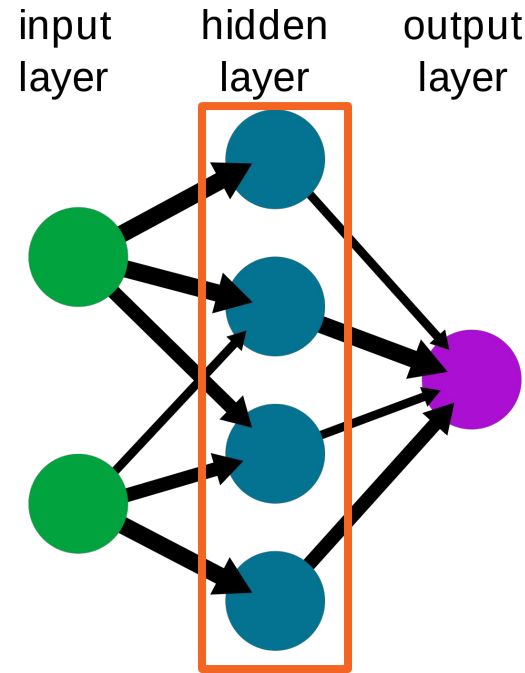
**Each circle is a
“node” with an
input and an
output**

Single-layer perceptron



**Each node also
has a
corresponding
weight and *bias***

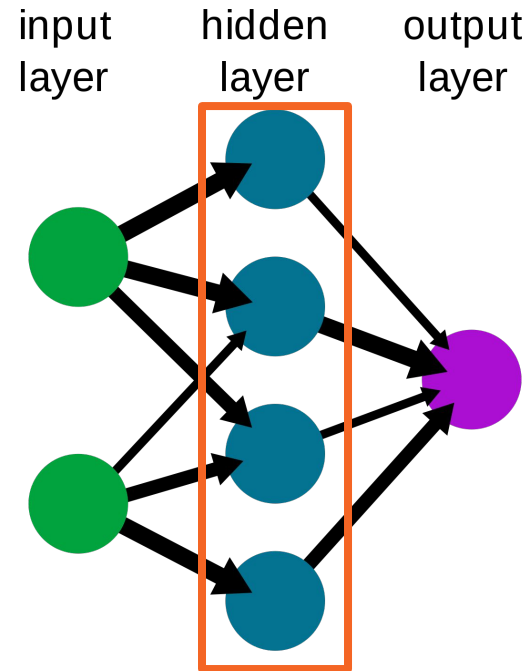
Single-layer perceptron



**Each node also
has a
corresponding
weight and *bias***

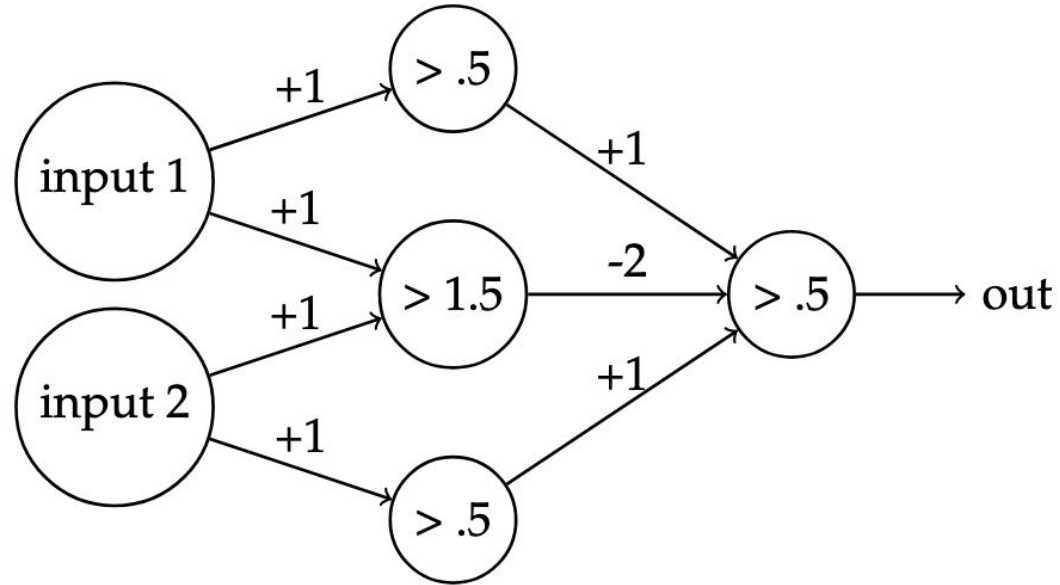
coefficient intercept

Single-layer perceptron

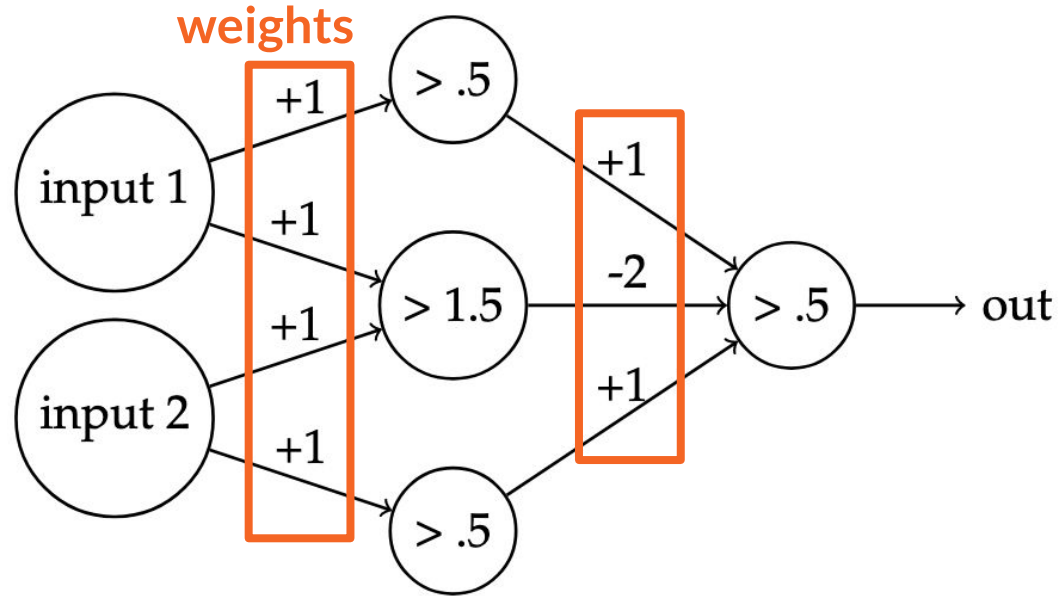


**Think of each
node as its own
linear regression**

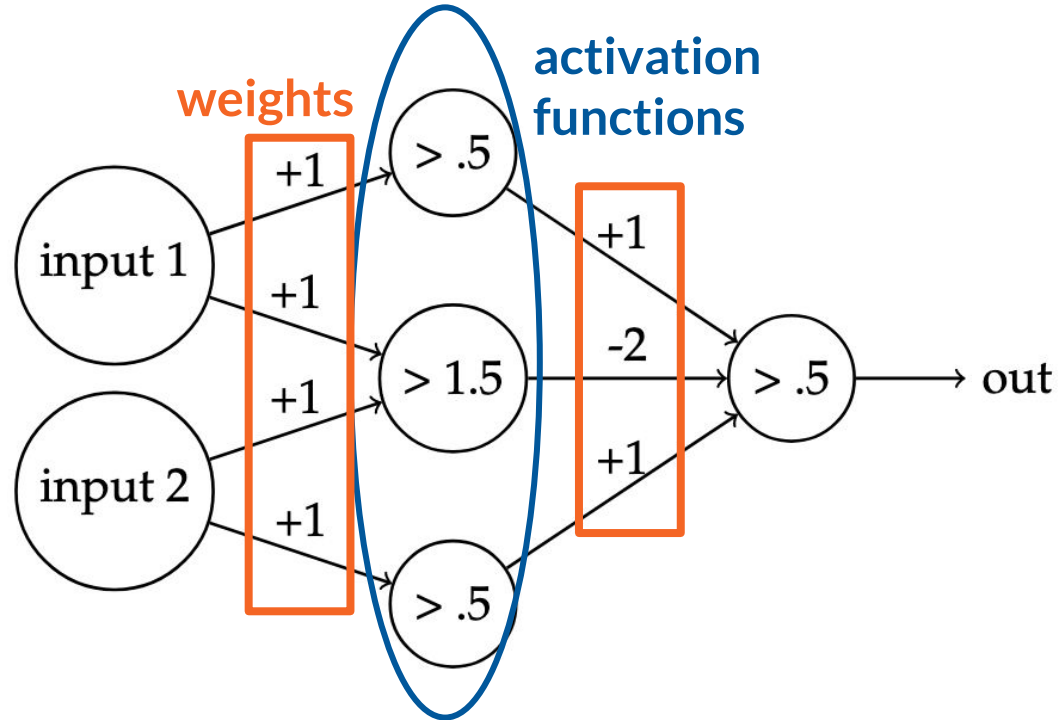
Multi-layer perceptron: XOR, 1985



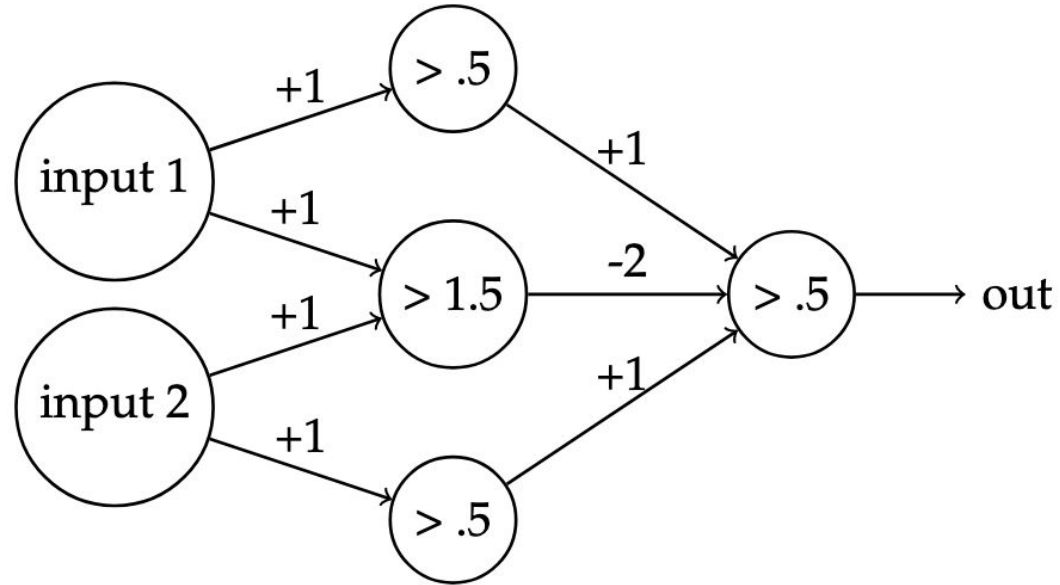
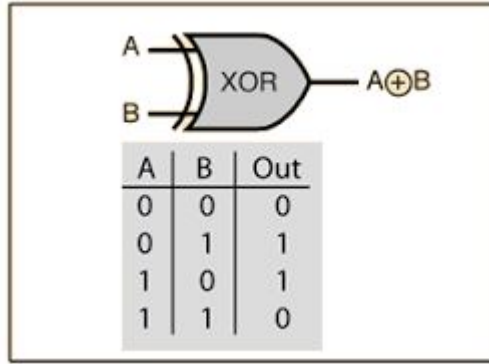
Multi-layer perceptron: XOR, 1985



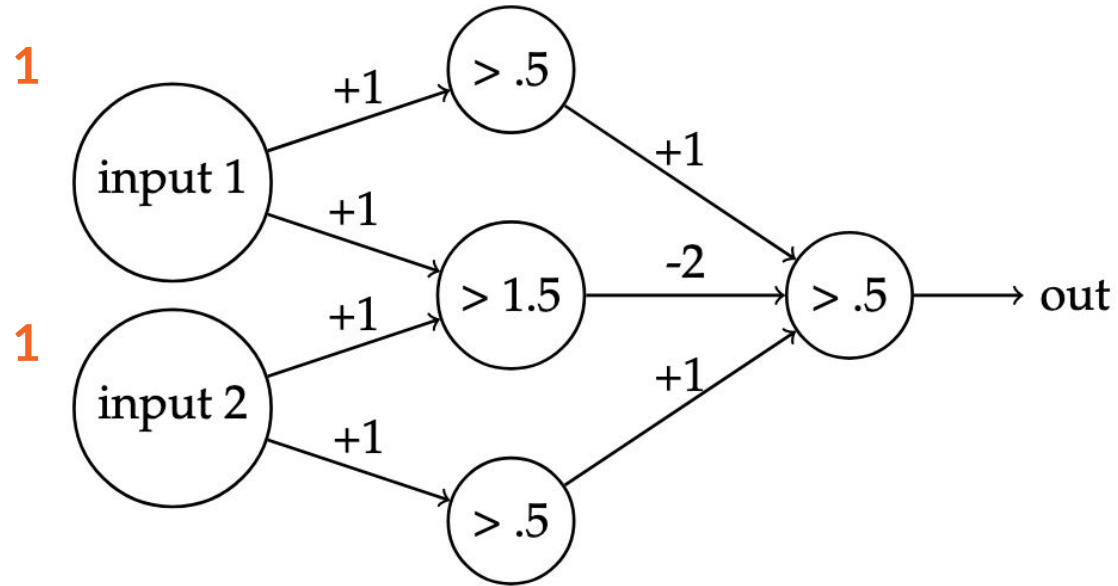
Multi-layer perceptron: XOR, 1985



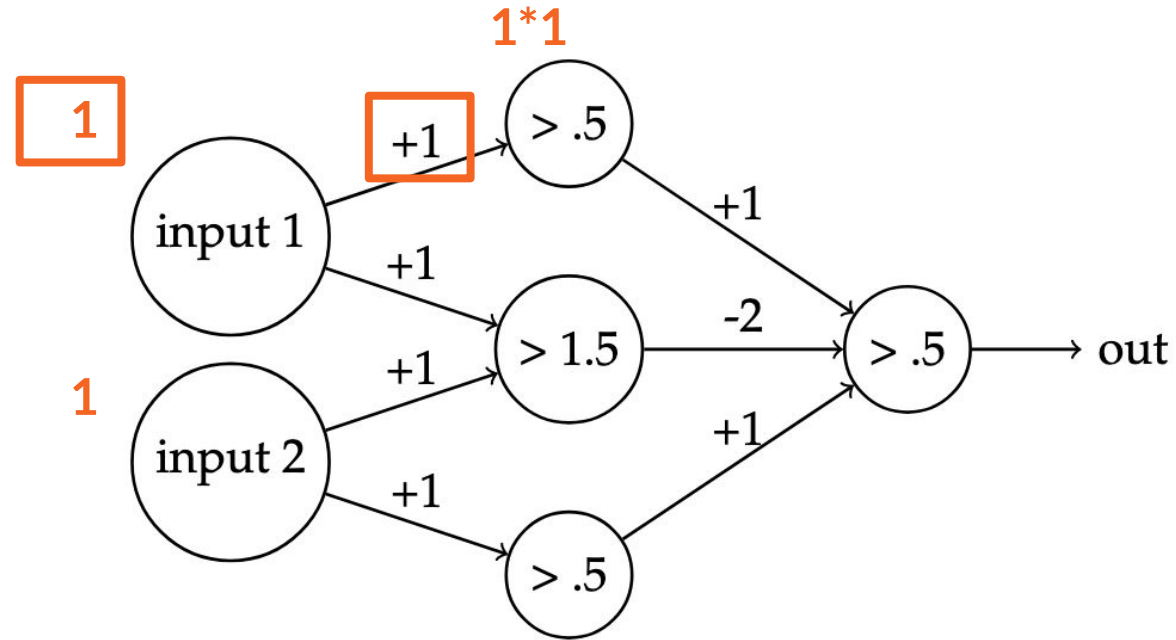
Multi-layer perceptron: XOR, 1985



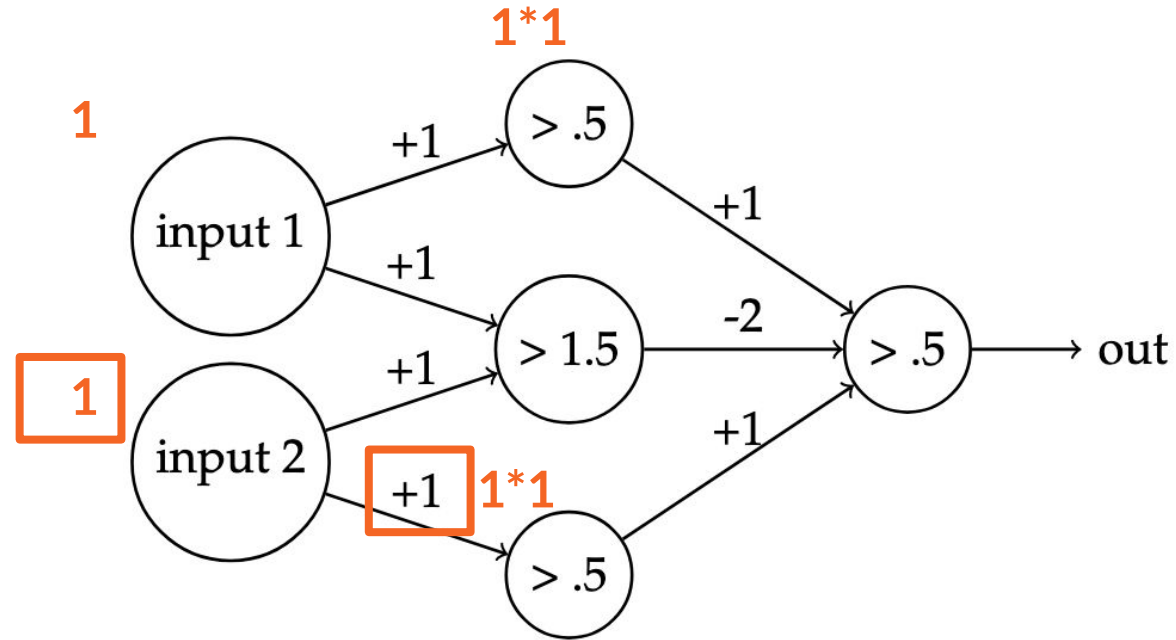
Multi-layer perceptron: XOR, 1985



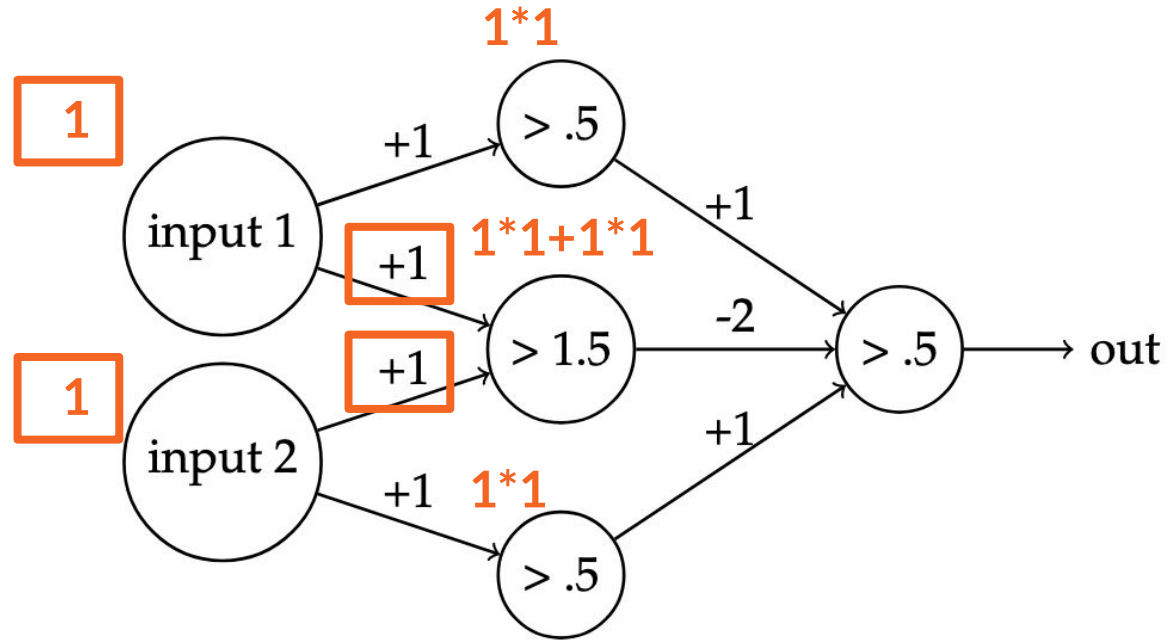
Multi-layer perceptron: XOR, 1985



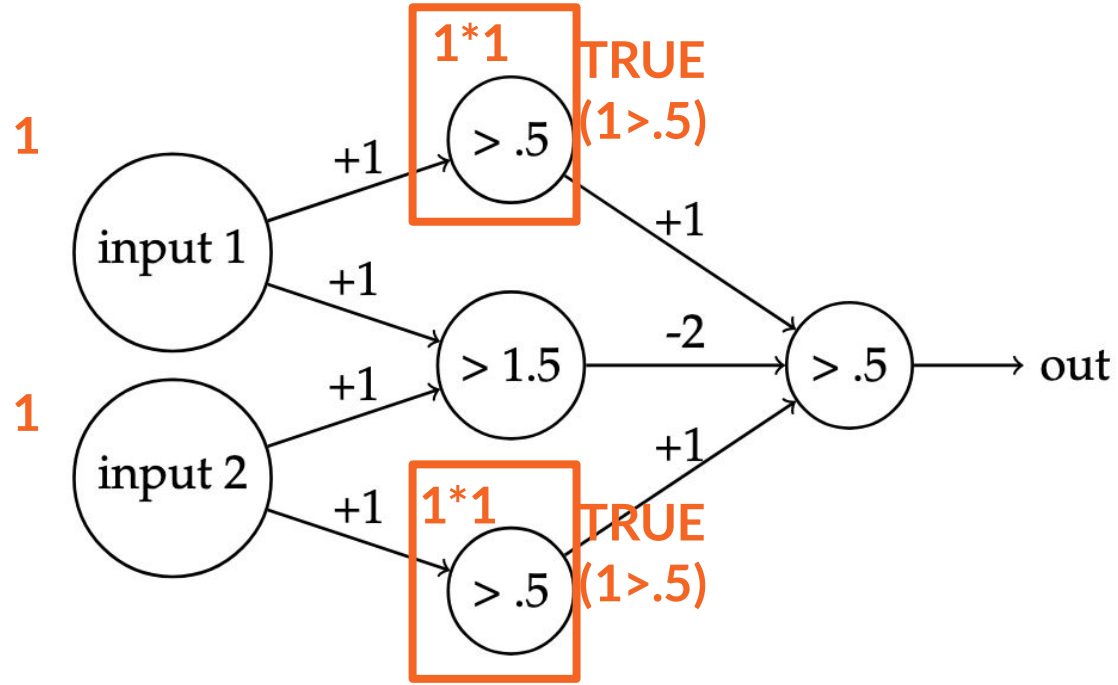
Multi-layer perceptron: XOR, 1985



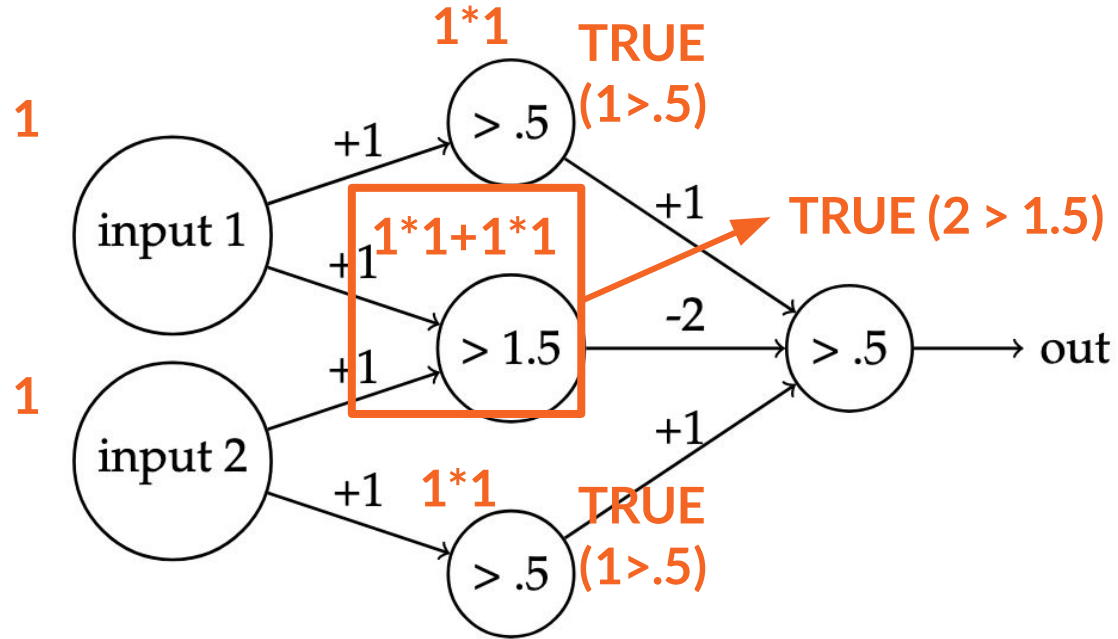
Multi-layer perceptron: XOR, 1985



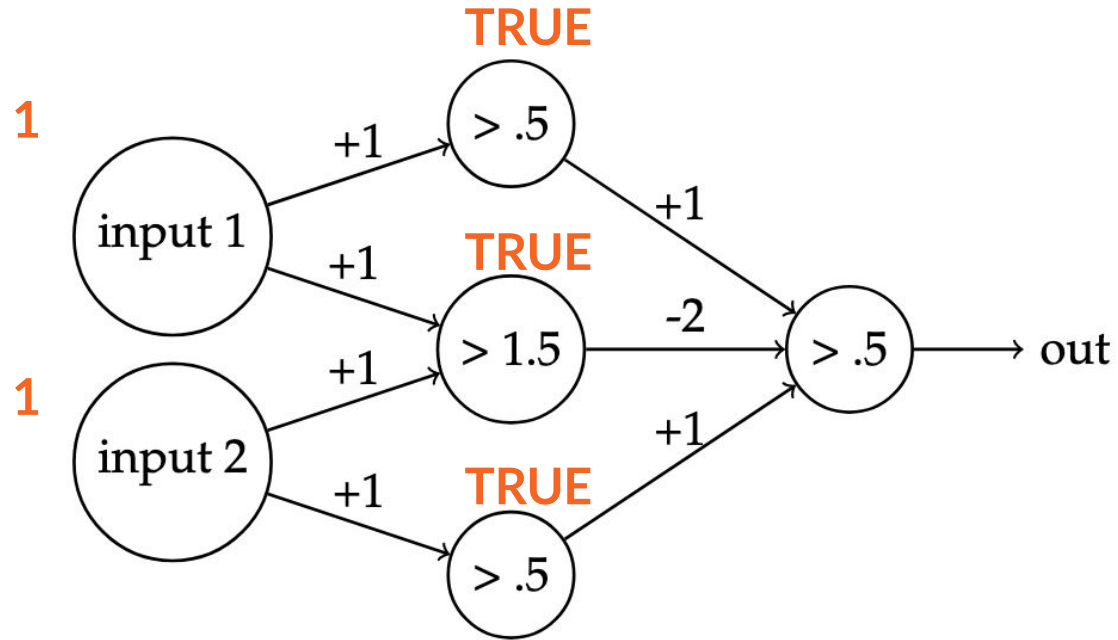
Multi-layer perceptron: XOR, 1985



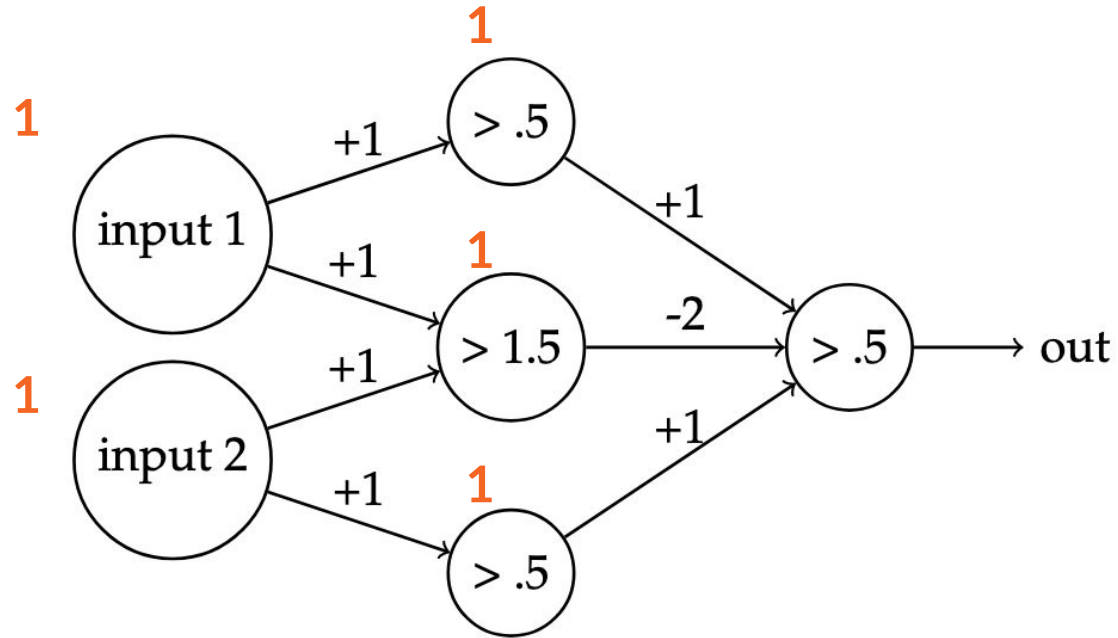
Multi-layer perceptron: XOR, 1985



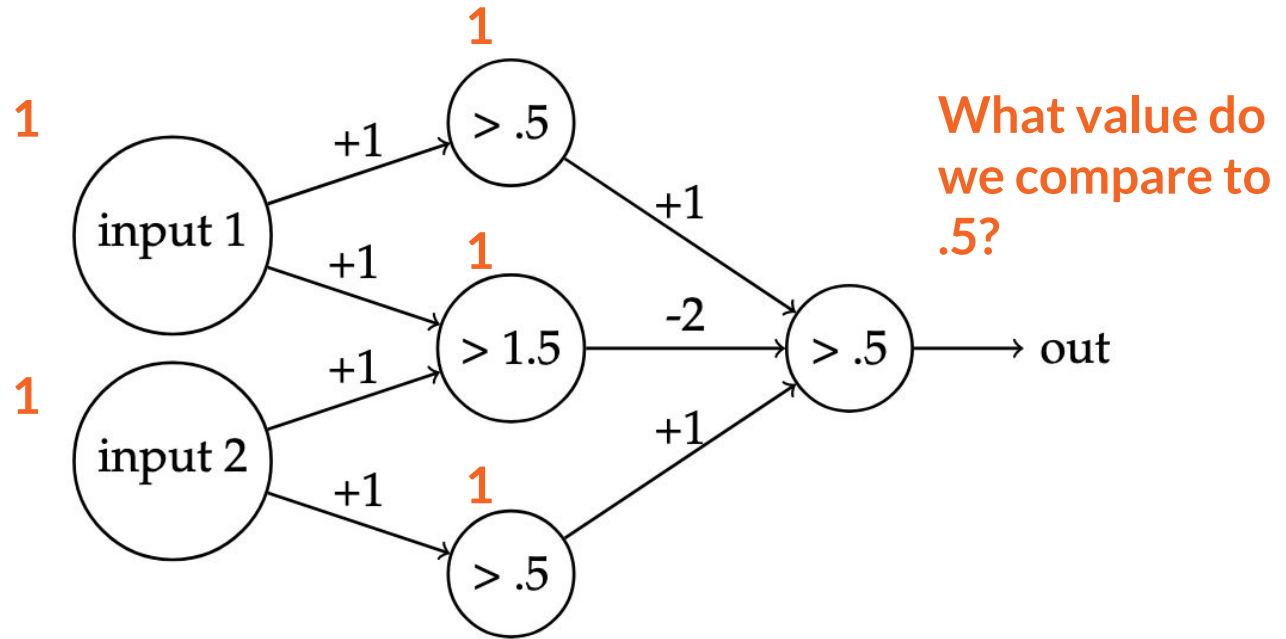
Multi-layer perceptron: XOR, 1985

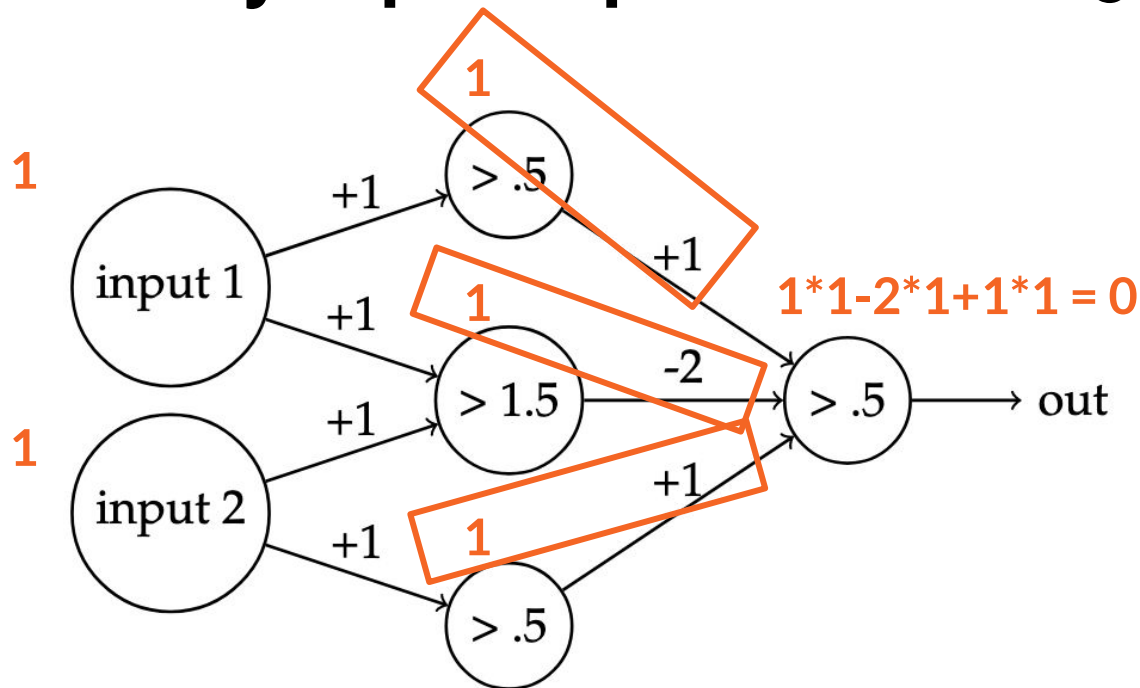


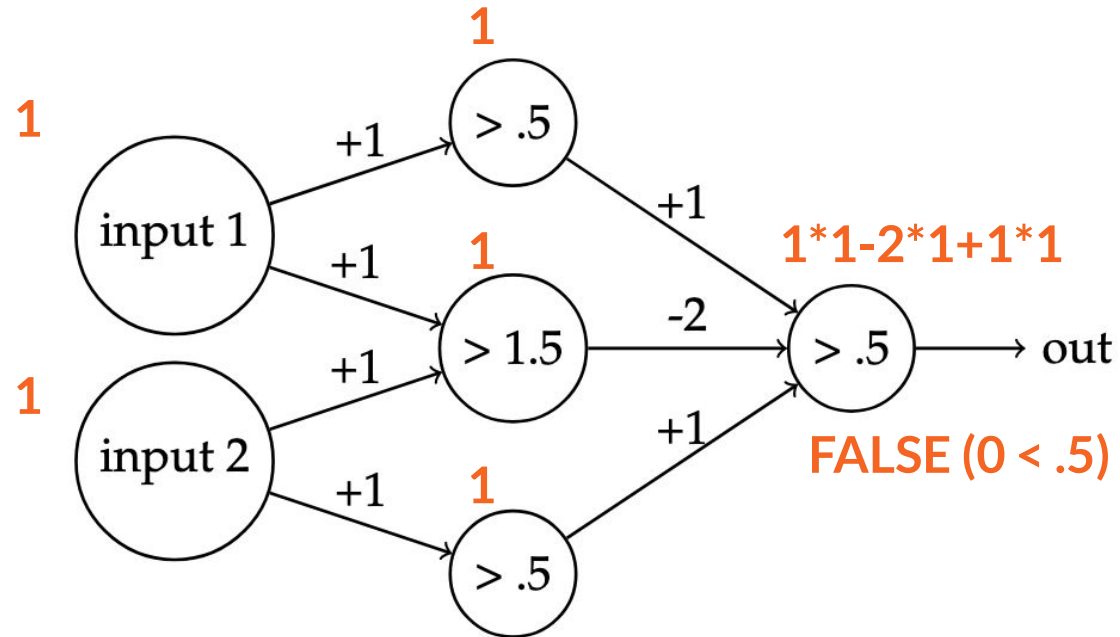
Multi-layer perceptron: XOR, 1985



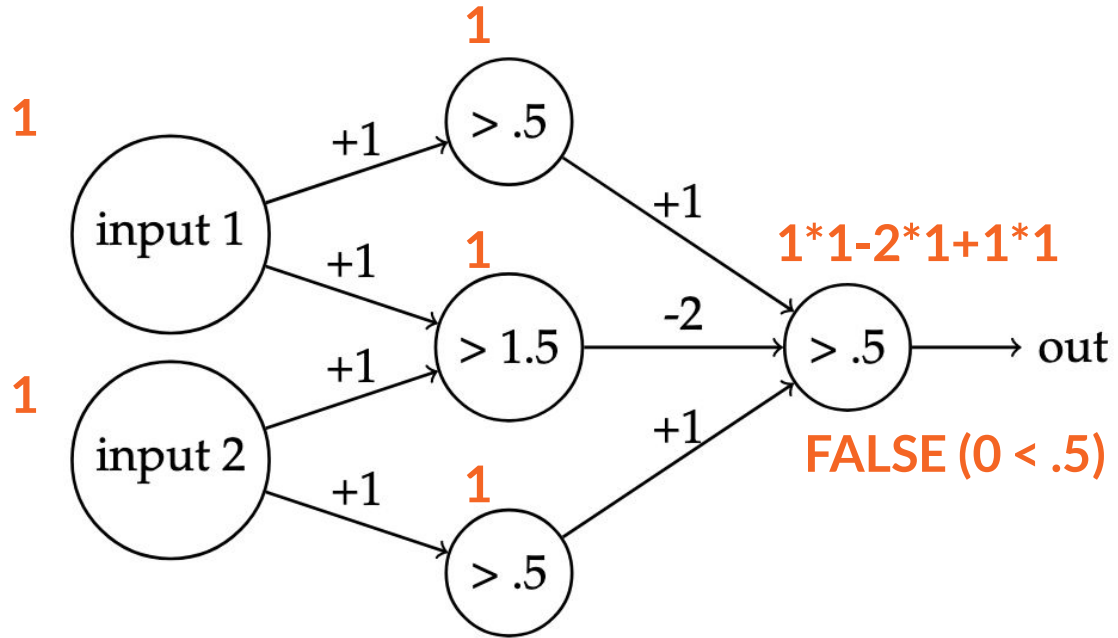
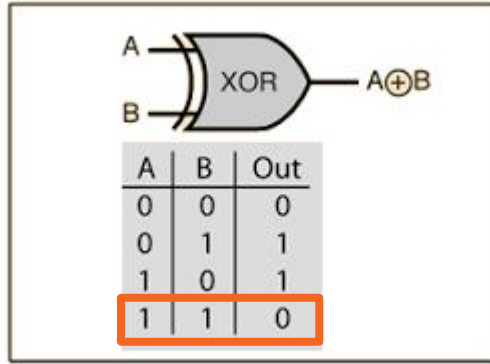
Multi-layer perceptron: XOR, 1985



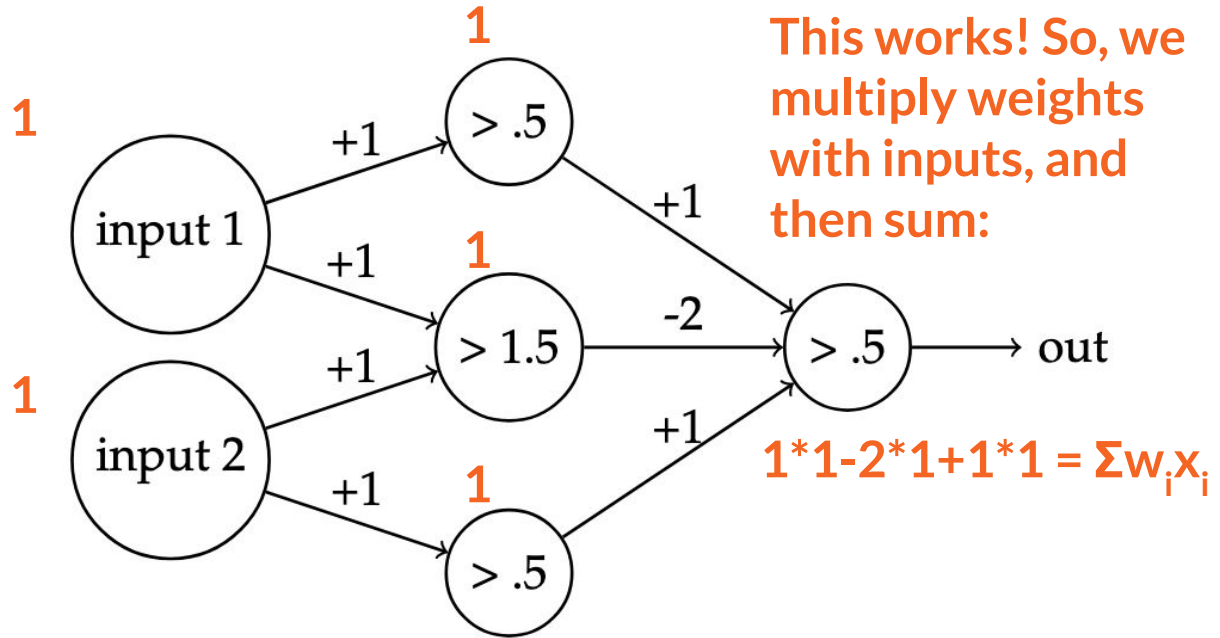
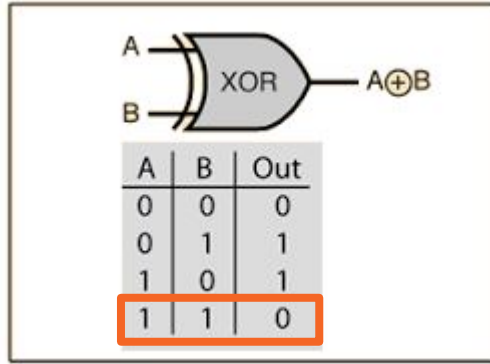




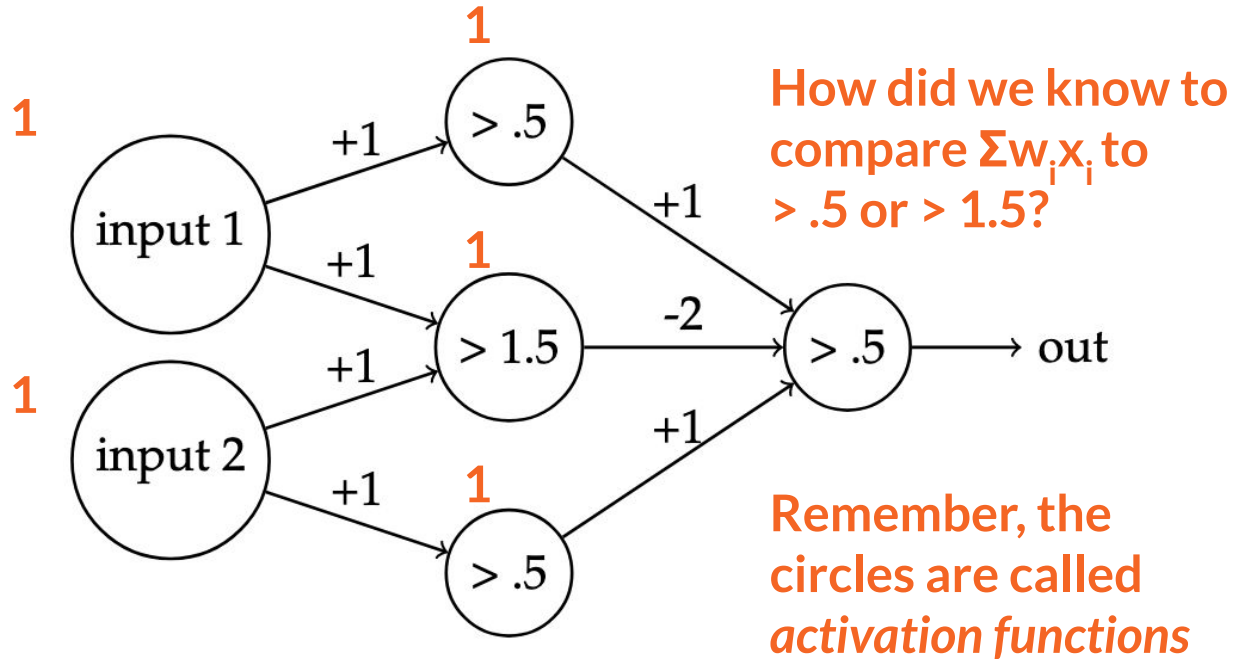
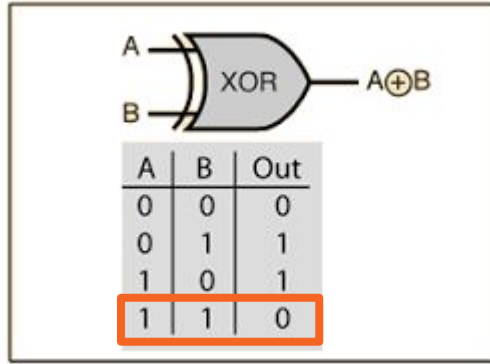
Multi-layer perceptron: XOR, 1985



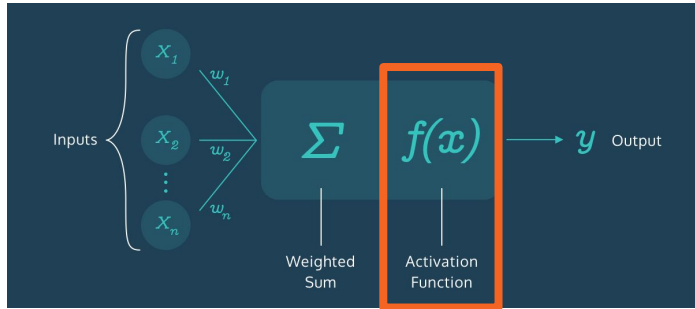
Multi-layer perceptron: XOR, 1985



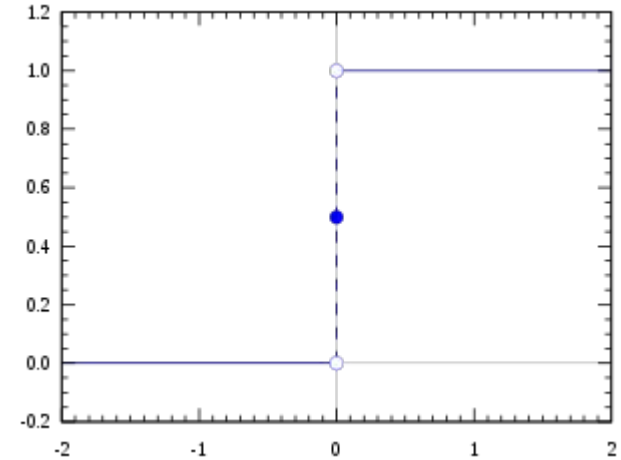
Multi-layer perceptron: XOR, 1985



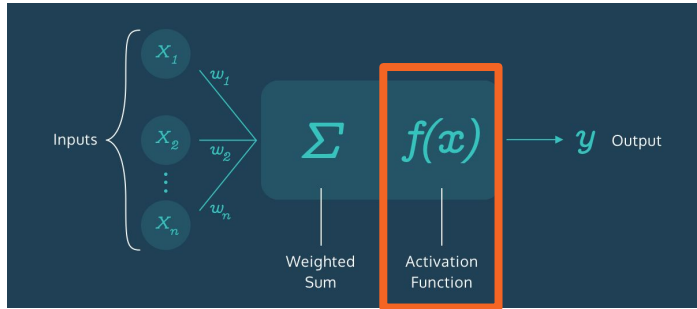
Common activation functions



$$\text{Step function: } f(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

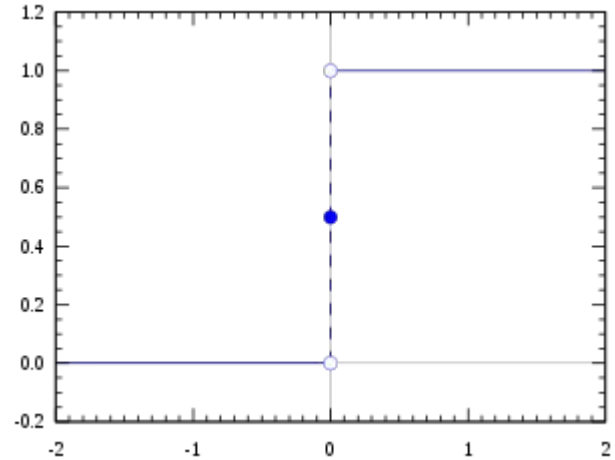


Common activation functions

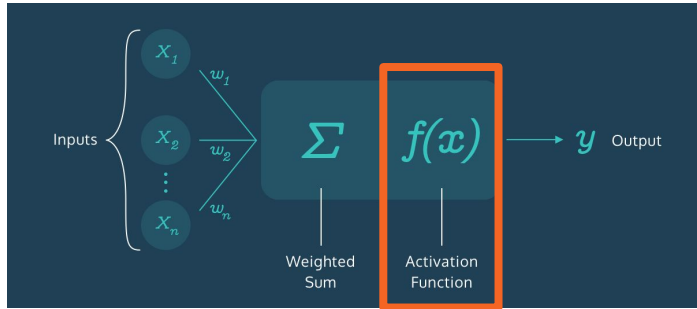


Step function: $f(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$

(We compared $\sum w_i x_i$ to being over/under .5 or 1.5 instead of over/under 0, which is just shifting a step function!)

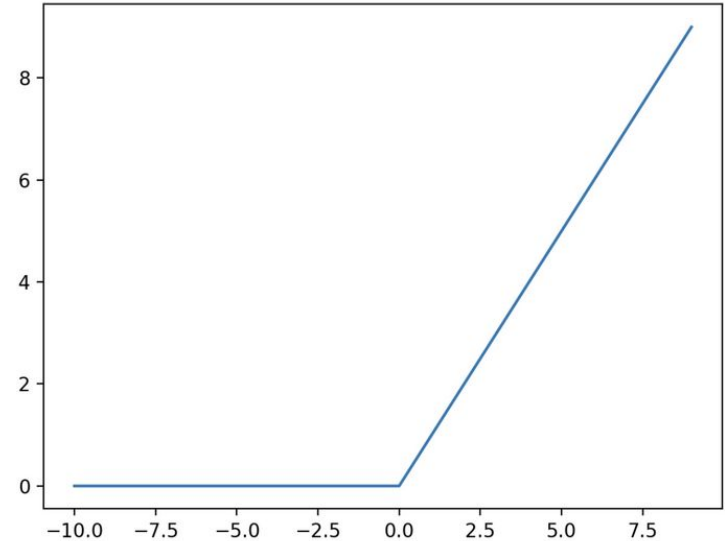
$$= \begin{cases} 1 & \text{if } \sum w_i x_i \geq 0 \\ 0 & \text{if } \sum w_i x_i < 0 \end{cases}$$


Common activation functions

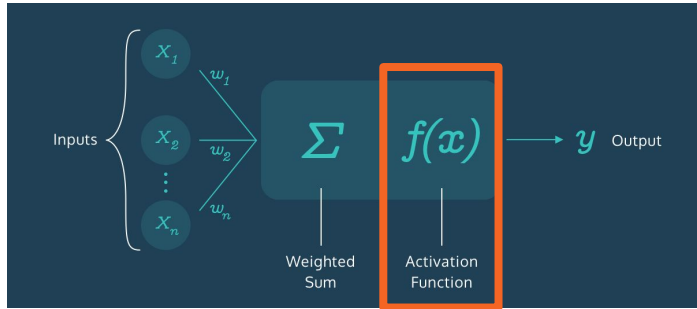


**Rectified Linear
Units (ReLUs):**

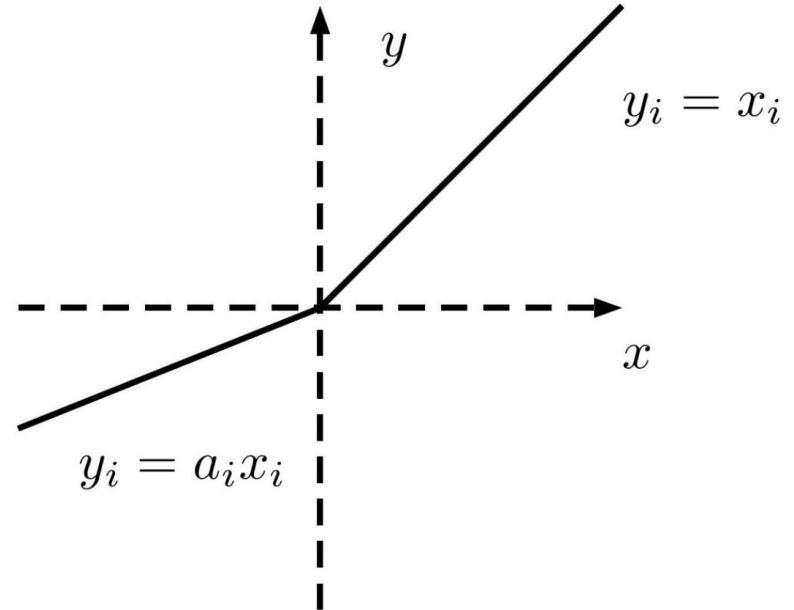
$$f(x) = \max(0, x)$$



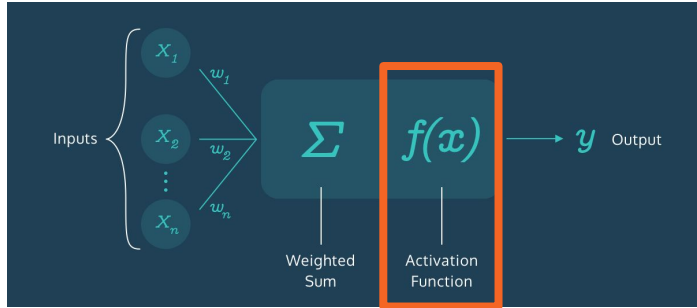
Common activation functions



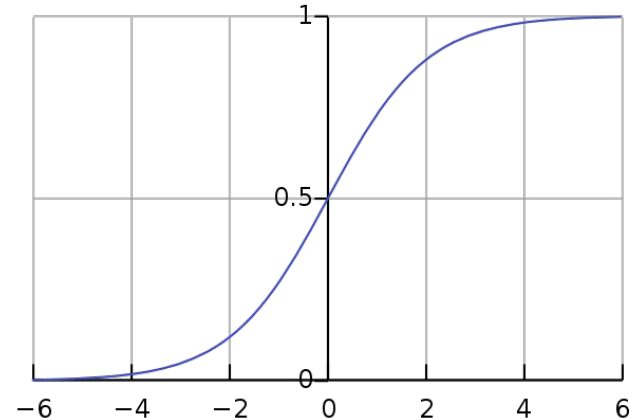
Leaky ReLU:



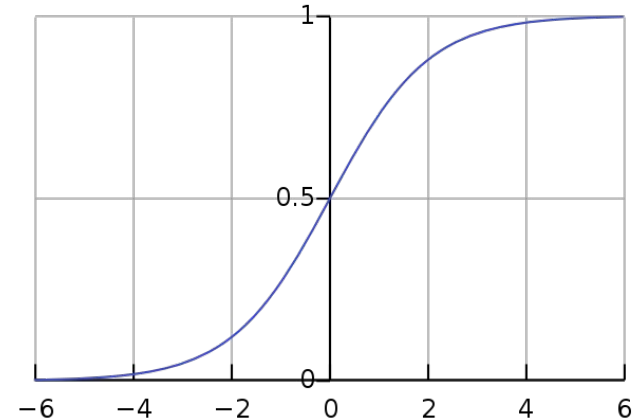
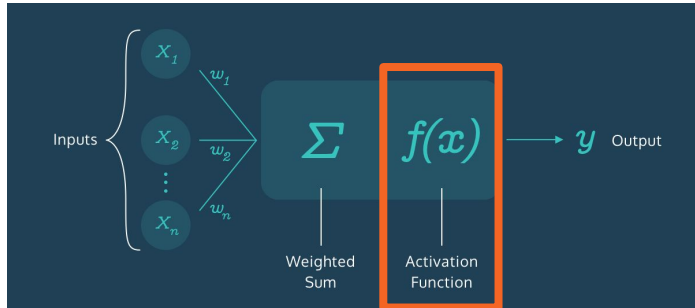
Common activation functions



**What's this
function called?**

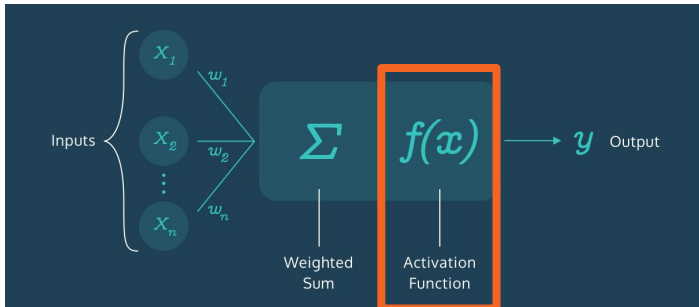


Common activation functions

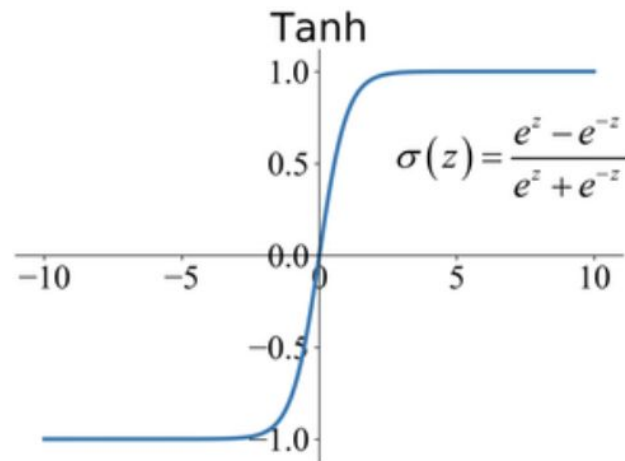


Sigmoid (logistic) activation function: $g(z) = \frac{1}{1 + e^{-z}}$

Common activation functions



Hyperbolic tangent (Tanh) looks like sigmoid, but goes from -1 to 1 instead of 0 to 1.



Neural nets

- NNs are just multi-layer perceptrons!

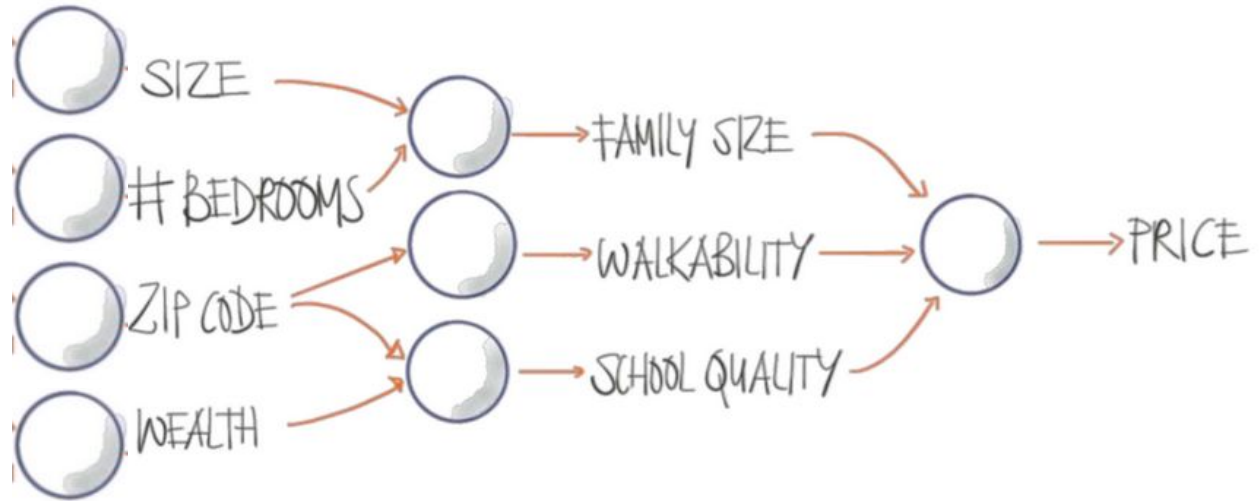
Neural nets

- NNs are just multi-layer perceptrons!
- How do you get weight and bias values?
 - Similar to regressions: “training the model” results in a weight (coefficient) and bias (intercept) for each node

Neural nets

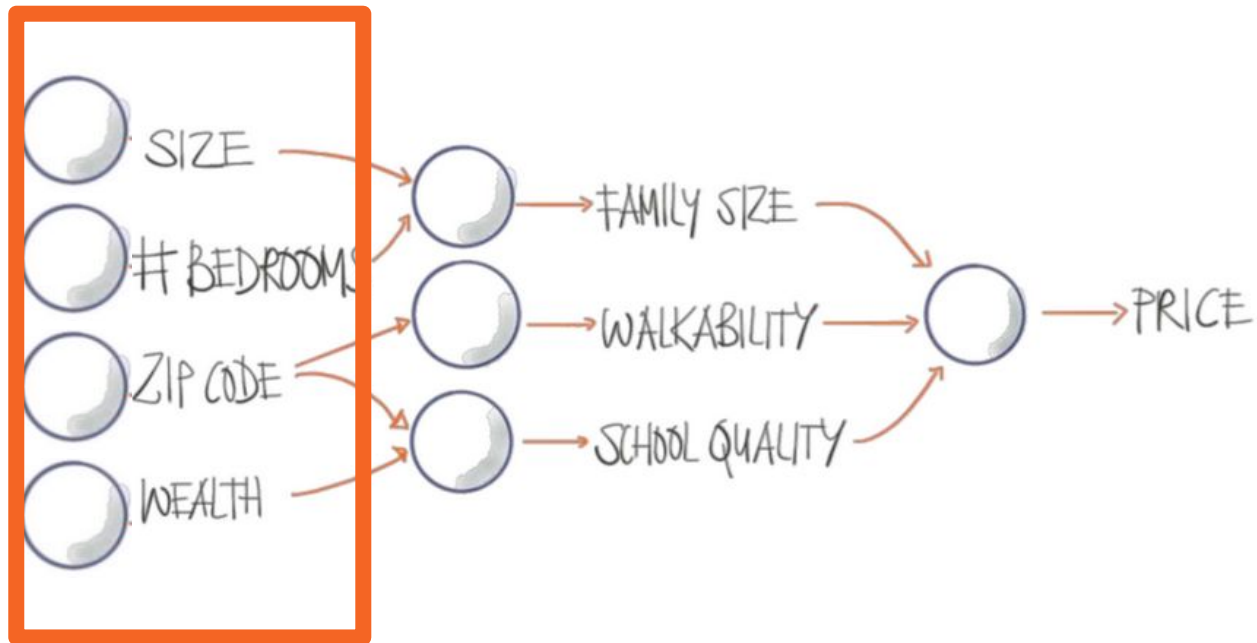
- NNs are just multi-layer perceptrons!
- How do you get weight and bias values?
 - Similar to regressions: “training the model” results in a weight (coefficient) and bias (intercept) for each node
 - Randomly initialize your weights. Iterate:
 - **Gradient descent:** finds where to minimize loss function
 - **Backpropagation:** updates weights based on gradient of loss

Interpreting neural nets

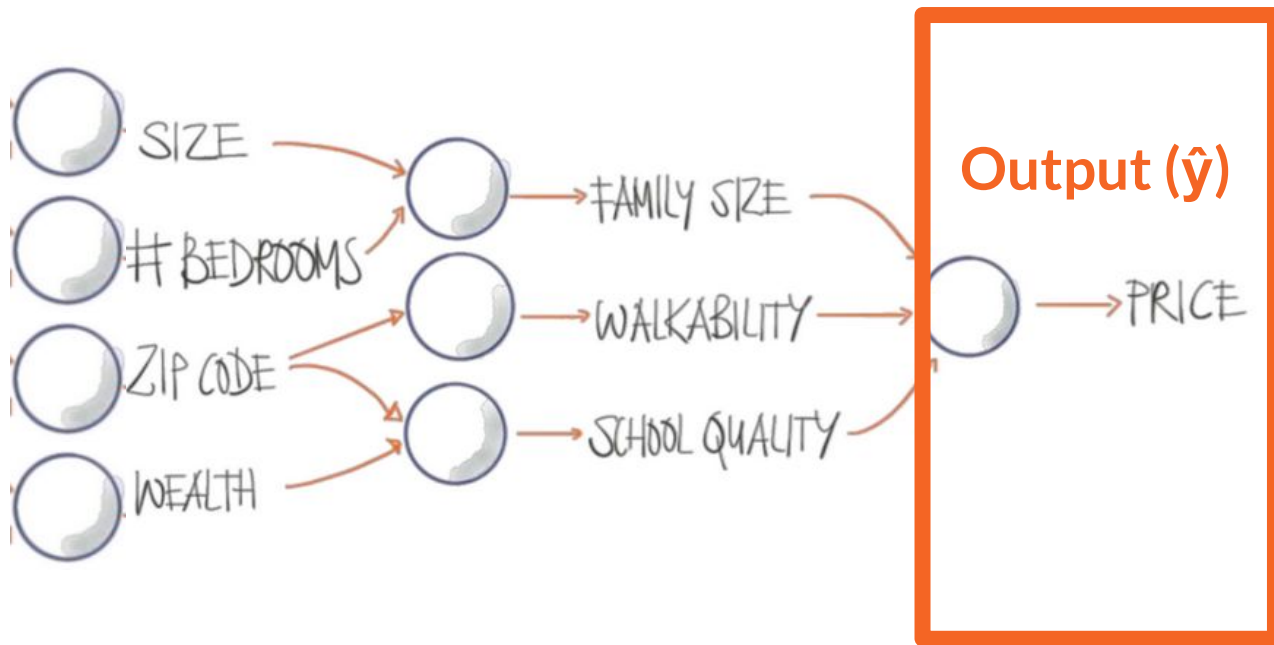


Interpreting neural nets

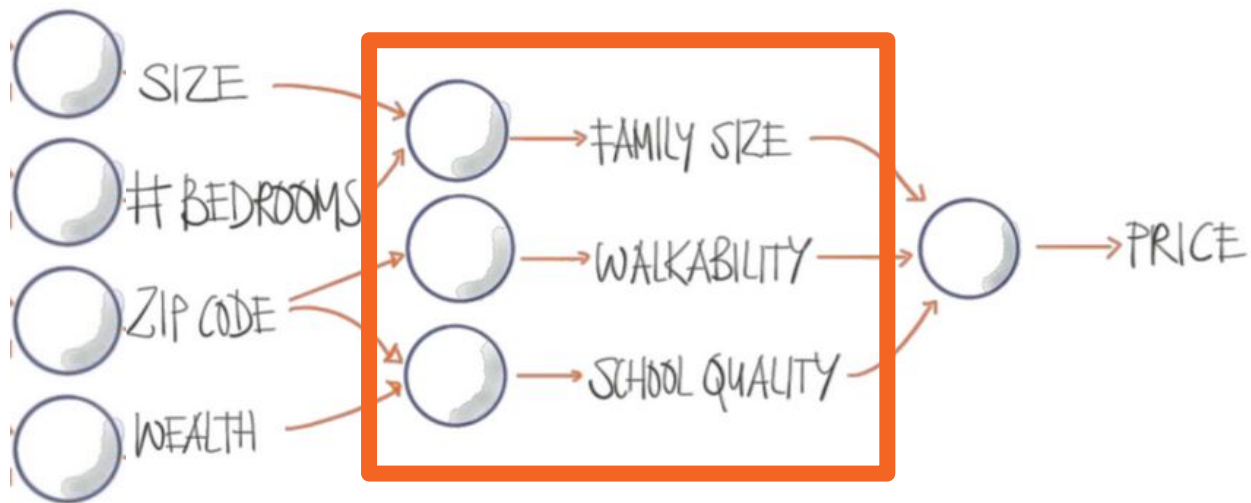
Input (x's):
known



Interpreting neural nets



Interpreting neural nets

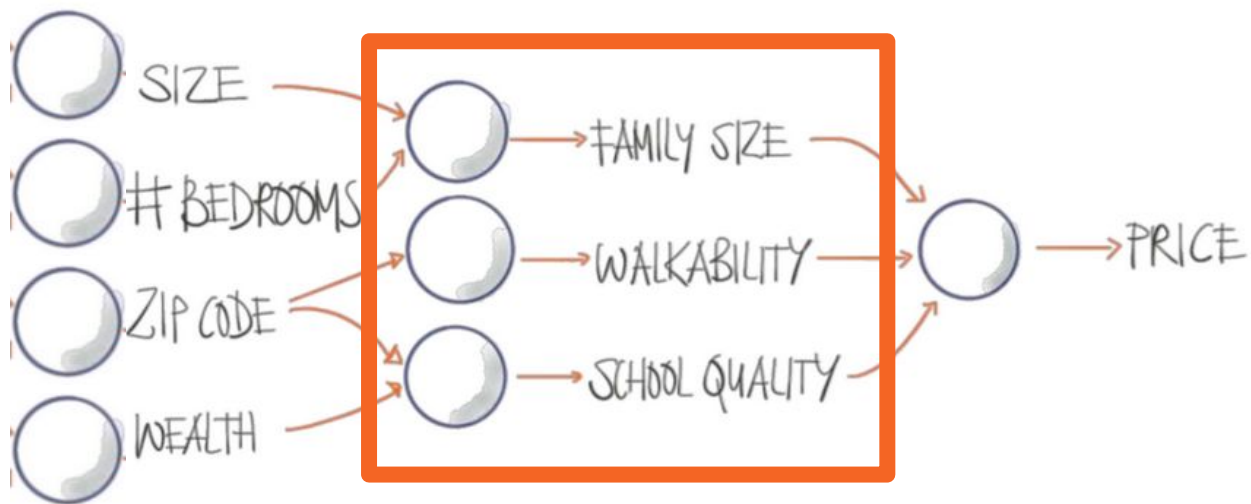


Hidden layer: unknown!

We're just guessing the meanings of each node

Neural net layers

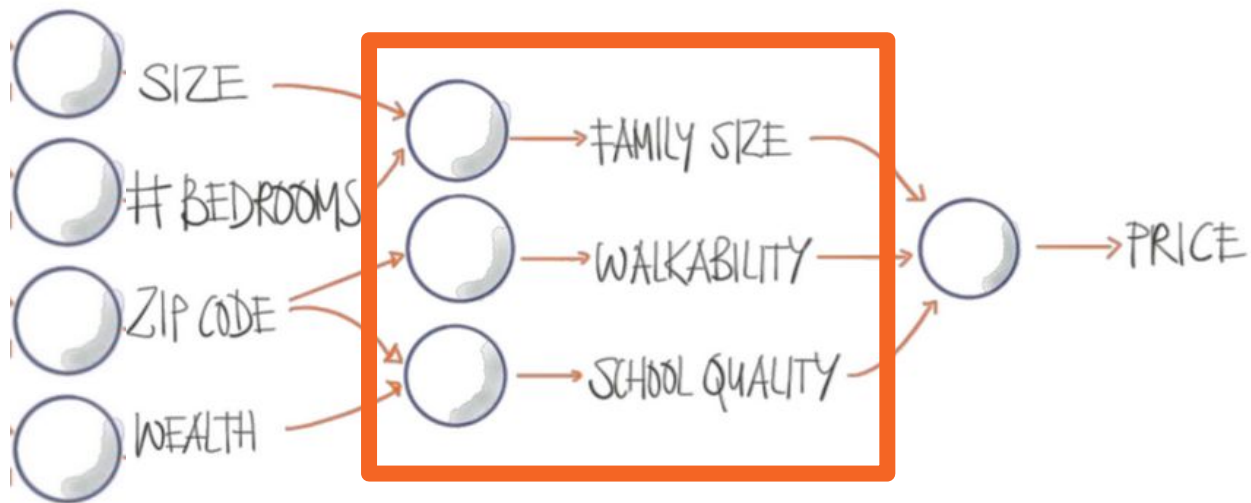
*Similar to how we had to guess concepts for SVD decompositions; in practice this is much harder for NNs



Hidden layer: unknown!

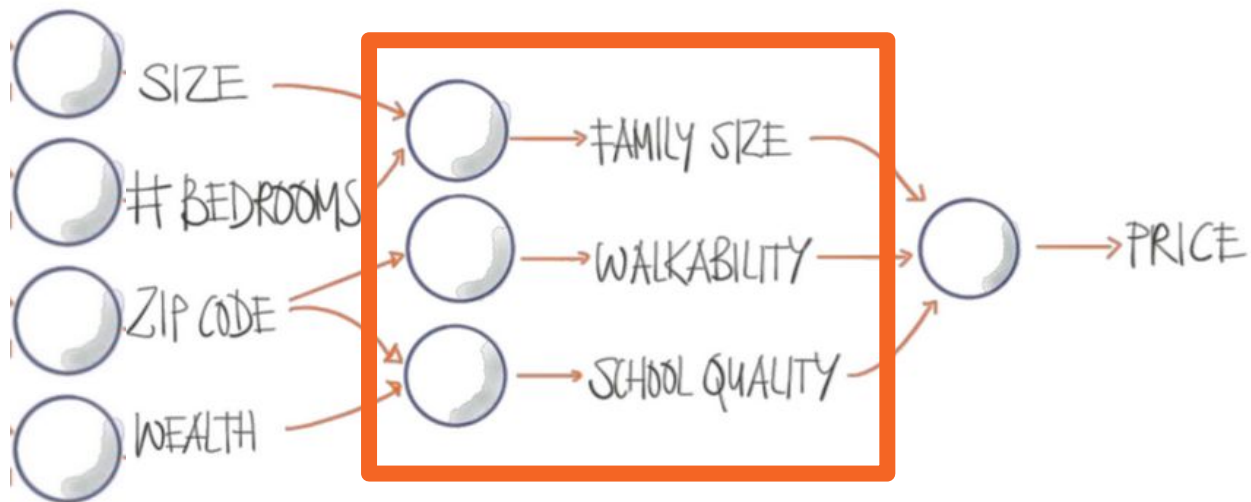
We're just guessing the meanings of each node

Neural net hidden layer



What we have control over: # hidden layers, # neurons in each layer, and what “activation” each layer uses

Neural net hidden layer

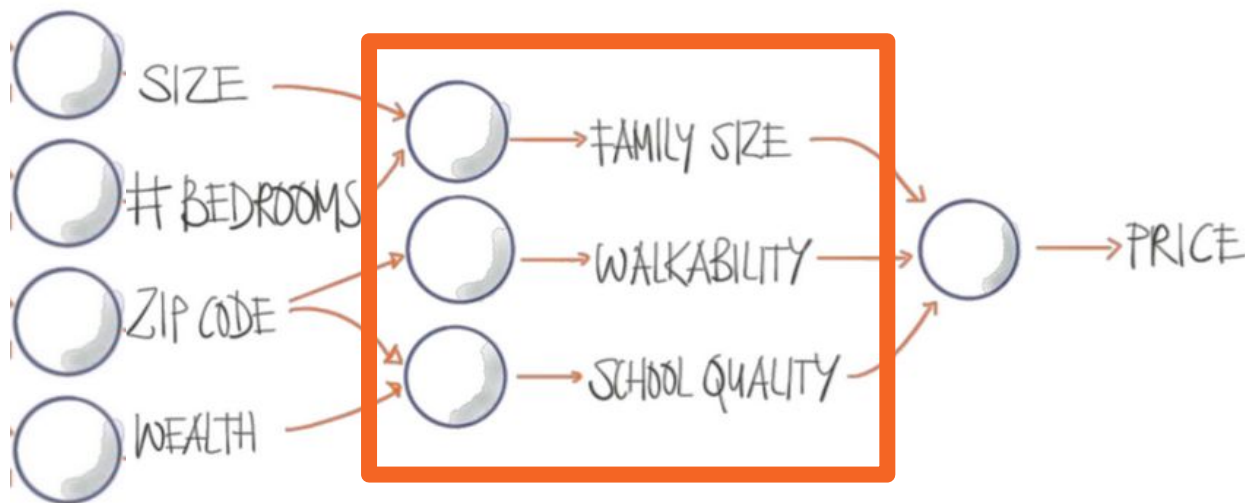


“Neural net architecture”



What we have control over: # hidden layers, # neurons in each layer, and what “activation” each layer uses

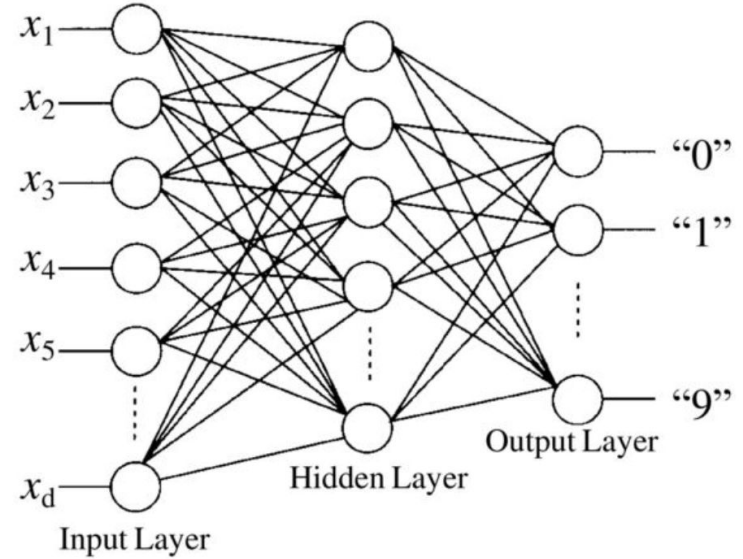
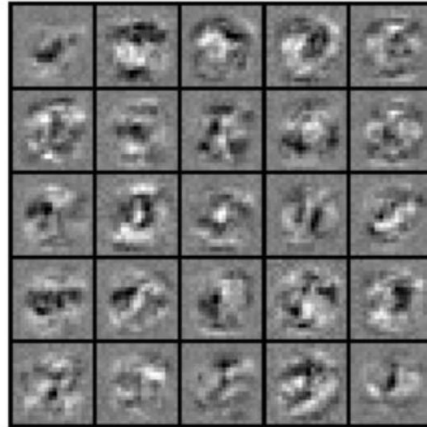
Neural net hidden layer



Should always be reported in your research for reproducibility!

What we have control over: # hidden layers, # neurons in each layer, and what “activation” each layer uses

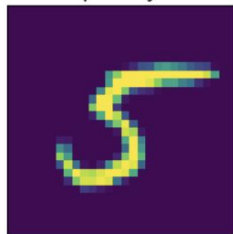
7	9	6	5	8	7	4	4	1	0
0	7	3	3	2	4	8	4	5	7
6	6	3	2	9	2	3	3	2	6
1	3	7	1	5	6	5	2	4	4
7	0	9	0	7	5	8	9	5	4
4	6	6	5	0	2	1	3	6	9
8	5	1	8	9	7	8	7	3	6
1	0	2	8	2	3	0	5	1	5
6	7	8	2	5	3	9	7	0	0
7	9	3	9	8	5	7	2	9	8



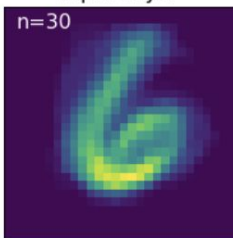
Visualization of Hidden Layer

MNIST data

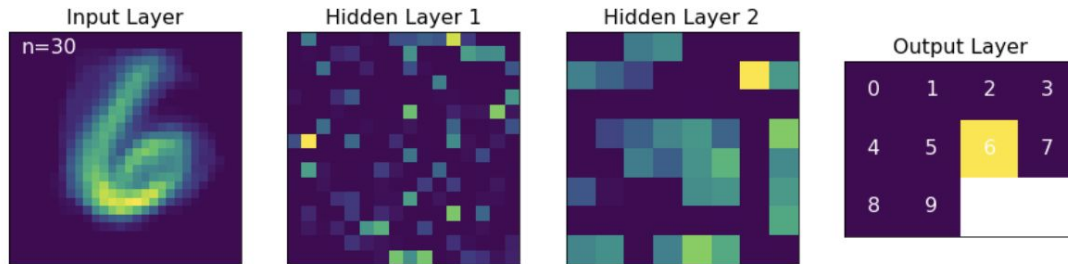
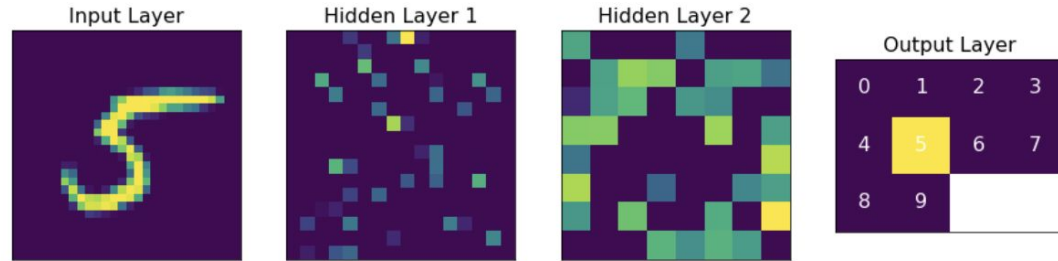
Input Layer



Input Layer



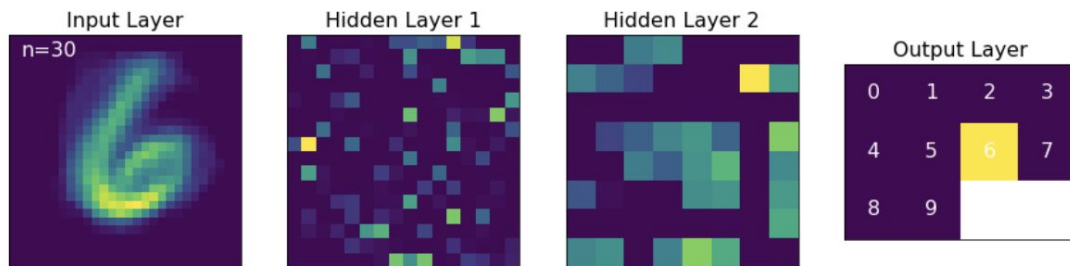
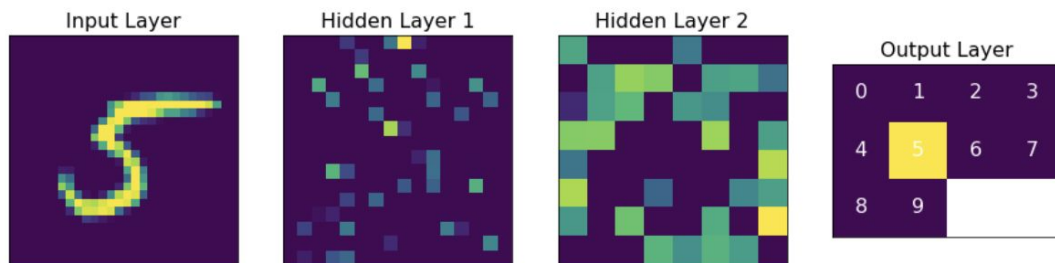
Which neurons are triggered?



Which neurons are triggered?

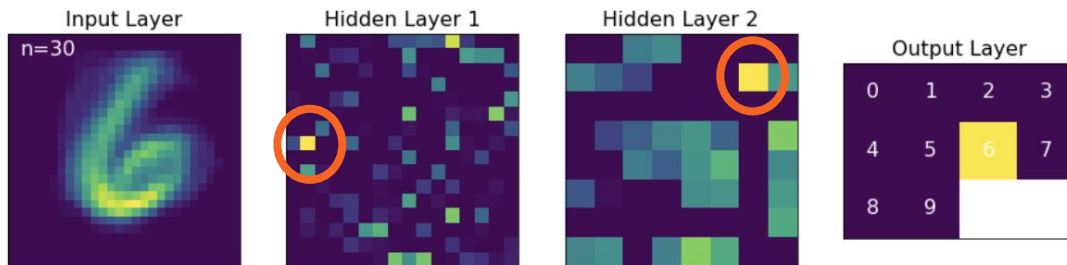
Yellow = more important

Blue = less important

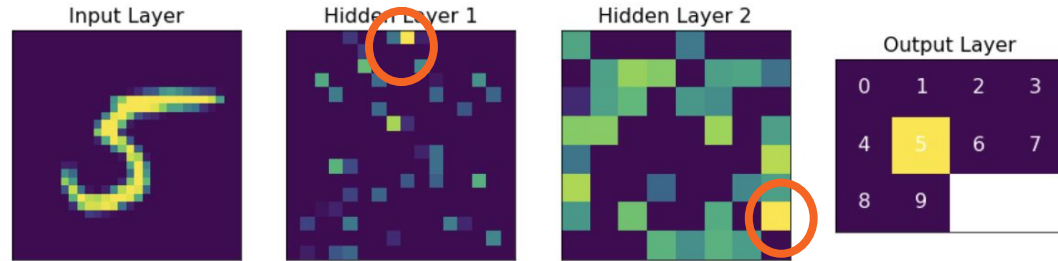


Which neurons are triggered?

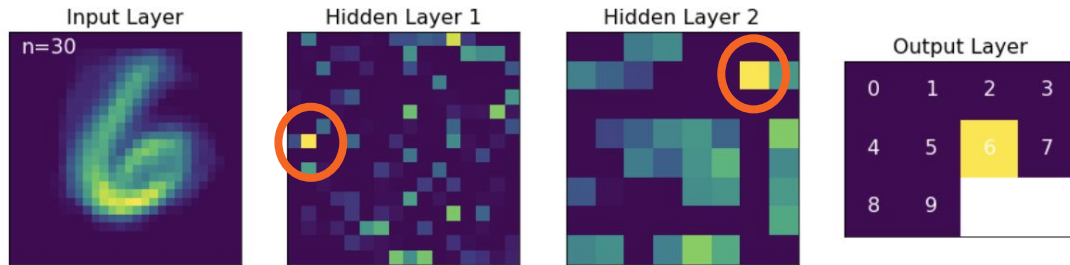
The “important”
neurons are different
for different inputs!



Which neurons are triggered?



We can't really tell what each hidden layer's "concept" is, but we tell the NN is using them to correctly noticed that $5! = 6$






How do we do this in Python?

- *More libraries!*
 - Tensorflow (v1 and v2)
 - Keras
 - Jax
 - PyTorch
 - Theano

Running neural nets

- A lot of these steps are the same!

numpy, scipy,
scikit-learn, etc.

- 
- Step 1: decide on inputs / outputs
 - Step 2: data preprocessing
 - Step 3: decide on an evaluation metric
 - Step 4: split your data

Can use, e.g., Keras

- 
- Step 5: run the model
 - Step 6: interpret results

**Import relevant
packages**

E.g., in Keras:

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Activation
from tensorflow.keras.optimizers import Adam
```

Import relevant
packages

E.g., in Keras:

Deal with your NN architecture


```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Activation
from tensorflow.keras.optimizers import Adam
```

Deals with gradient descent

E.g., in Keras:

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Activation
from tensorflow.keras.optimizers import Adam
```

Define NN



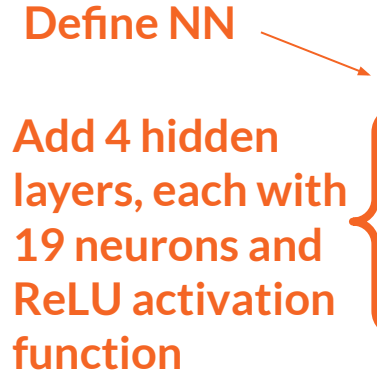
```
model = Sequential()
model.add(Dense(19, activation='relu'))
model.add(Dense(19, activation='relu'))
model.add(Dense(19, activation='relu'))
model.add(Dense(19, activation='relu'))
model.add(Dense(1)) model.compile(optimizer='Adam', loss='mes')
```

E.g., in Keras:

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Activation
from tensorflow.keras.optimizers import Adam
```

Define NN

Add 4 hidden
layers, each with
19 neurons and
ReLU activation
function



```
model = Sequential()
model.add(Dense(19, activation='relu'))
model.add(Dense(19, activation='relu'))
model.add(Dense(19, activation='relu'))
model.add(Dense(19, activation='relu'))
model.add(Dense(1))
model.compile(optimizer='Adam', loss='mse')
```

E.g., in Keras:

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Activation
from tensorflow.keras.optimizers import Adam
```

Define NN

```
model = Sequential()
{
    model.add(Dense(19, activation='relu'))
    model.add(Dense(19, activation='relu'))
    model.add(Dense(19, activation='relu'))
    model.add(Dense(19, activation='relu'))
}
model.add(Dense(1)) model.compile(optimizer='Adam', loss='mes')
```

Add output layer
(final single
neuron)

E.g., in Keras:

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Activation
from tensorflow.keras.optimizers import Adam
```

Define NN

```
model = Sequential()
{
    model.add(Dense(19, activation='relu'))
    model.add(Dense(19, activation='relu'))
    model.add(Dense(19, activation='relu'))
    model.add(Dense(19, activation='relu'))
}
model.add(Dense(1)) model.compile(optimizer='Adam', loss='mse')
```

Add output layer
(final single
neuron)

Define loss function

E.g., in Keras:

Fit the model by taking
"batches" of 128 data
rows at a time,
sweeping over the full
data set 400 times
("epochs")

```
model.fit(x=X_train,y=y_train,  
          validation_data=(X_test,y_test),  
          batch_size=128,epochs=400)  
model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
dense (Dense)	multiple	418

dense_1 (Dense)	multiple	380

dense_2 (Dense)	multiple	380

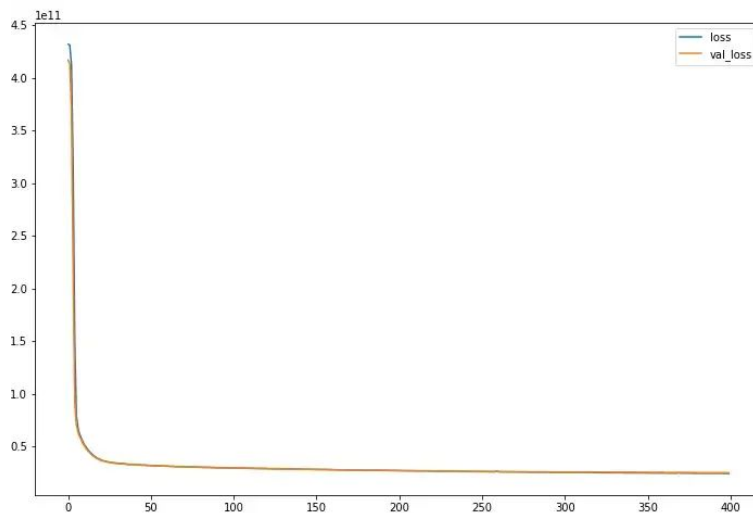
dense_3 (Dense)	multiple	380

dense_4 (Dense)	multiple	20
=====		
Total params: 1,578		
Trainable params: 1,578		
Non-trainable params: 0		

E.g., in Keras:

```
loss_df = pd.DataFrame(model.history.history)
loss_df.plot(figsize=(12,8))
```

Plot loss over epochs



E.g., in Keras:

```
y_pred = model.predict(X_test) Predict y-hats
```

```
from sklearn import metrics Use sklearn to check evaluation metrics  
print('MAE:', metrics.mean_absolute_error(y_test, y_pred))  
print('MSE:', metrics.mean_squared_error(y_test, y_pred))  
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))  
print('VarScore:', metrics.explained_variance_score(y_test, y_pred))
```

Interpret: Regression vs. NN

Which is better?

Note: “Variance score”
refers to “explained variance
score” (similar to R^2)

Model: Keras Neural Net

```
Mean Absolute Error(MAE): 96667.89
Mean Squared Error(MSE): 24912134897.75
Root Mean Squared Error(RMSE): 157835.78
Variance score: 80.84
```

```
*****
```

Model: Sklearn Multiple Linear Regression

```
Mean Absolute Error(MAE): 124516.17
Mean Squared Error(MSE):39763621927.16
Root Mean Squared Error(RMSE):199408.18
Variance score: 69.42
```

Interpret: Regression vs. NN

Lower error and
higher explainability
is better!



Model: Keras Neural Net

```
Mean Absolute Error(MAE): 96667.89
Mean Squared Error(MSE): 24912134897.75
Root Mean Squared Error(RMSE): 157835.78
Variance score: 80.84
```

Model: Sklearn Multiple Linear Regression

```
Mean Absolute Error(MAE): 124516.17
Mean Squared Error(MSE): 39763621927.16
Root Mean Squared Error(RMSE): 199408.18
Variance score: 69.42
```

Make sure to think about the broader implications of a poorly-performing model!

People with no idea
about AI, telling me my AI
will destroy the world

Me wondering why my
neural network is
classifying a cat as a dog...



“Parameters” in machine learning

- Model **parameters** are the α/β s/weights that are assigned to each variable
 - These are internally set through the model’s learning process

“Parameters” in machine learning

- Model **parameters** are the α/β s/weights that are assigned to each variable
 - These are internally set through the model's learning process
- **Hyperparameters** are specific model settings
 - These are established prior to model learning/training
 - Hyperparameters do not change during model training; they guide training

Hyperparameters

- Hyperparameters are the *arguments* you can pass into a model when instantiating it
- There can be **many** hyperparameters for a single model

Hyperparameters

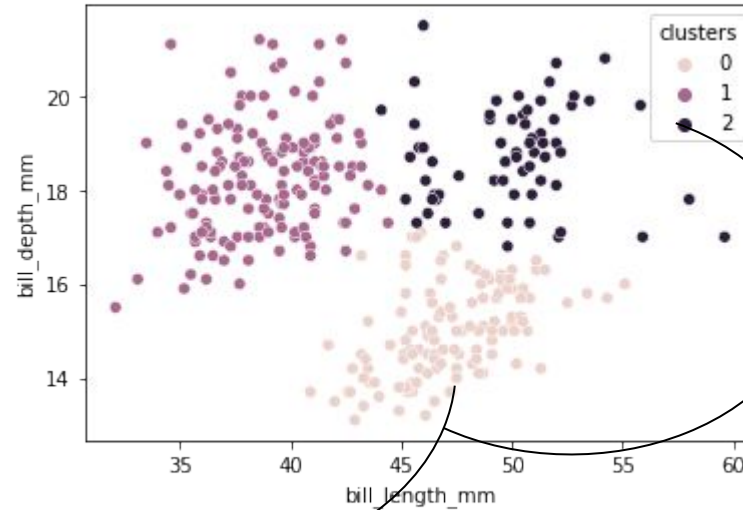
- Hyperparameters are the *arguments* you can pass into a model when instantiating it
- There can be **many** hyperparameters for a single model
- Choosing the best hyperparameters is an important part of the machine learning process
- Different hyperparameters can affect each other
 - Sort of like interaction effects, but with model performance!

Hyperparameters

- Hyperparameters are the *arguments* you can pass into a model when instantiating it
- There can be **many** hyperparameters for a single model
- Choosing the best hyperparameters is an important part of the machine learning process
- Different hyperparameters can affect each other
 - Sort of like interaction effects, but with model performance!

Can you think of any hyperparameters we've already encountered?

Output: k-means clustering



Examples of hyperparameters

- K-Means Clustering:
 - Number of clusters
 - `clustering = KMeans(n_clusters=5)`
- TF-IDF
 - Max/min number of documents a word can appear in
 - `tfidf_vectorizer = TfidfVectorizer(min_df=5)`
- SGD
 - Learning rate
- Neural networks
 - Model architecture (hidden layers, nodes, etc.)

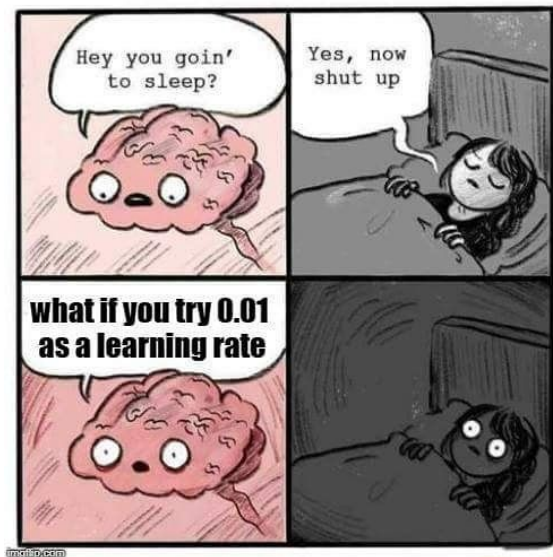
How do we choose these numbers??

The best hyperparameter values vary based on model, dataset, and evaluation metric

Examples of hyperparameters

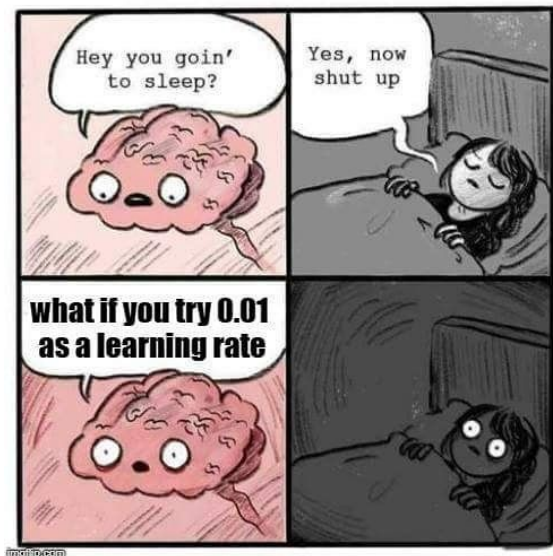
- K-Means Clustering:
 - Number of clusters
 - `clustering = KMeans(n_clusters=5)`
- TF-IDF
 - Max/min number of documents a word can appear in
 - `tfidf_vectorizer = TfidfVectorizer(min_df=5)`
- SGD
 - Learning rate
- Neural networks
 - Model architecture (hidden layers, nodes, etc.)

Hyperparameter tuning



- We always want to have the best-performing model
- How do we choose the hyperparameters that will “optimize” performance?
 - One option: manually choose a combination of a few hyperparameters and then see which model performs best

Hyperparameter tuning



- We always want to have the best-performing model
- How do we choose the hyperparameters that will “optimize” performance?
 - One option: manually choose a combination of a few hyperparameters and then see which model performs best
 - **Better option:** algorithmically “search” for best combinations of hyperparameters

Hyperparameter tuning

- **Better option:** algorithmically “search” for best combinations of hyperparameters
 - This is known as **hyperparameter optimization** or **hyperparameter tuning**
- Hyperparameter tuning tries various combinations of hyperparameters and how they affect model performance

Main steps of hyperparameter tuning

1. Choose a model and a set of relevant hyperparameters

Main steps of hyperparameter tuning

1. Choose a model and a set of relevant hyperparameters
2. Define an evaluation function (like F1 or MSE)

Main steps of hyperparameter tuning

1. Choose a model and a set of relevant hyperparameters
2. Define an evaluation function (like F1 or MSE)
3. Iterate through hyperparameters

Main steps of hyperparameter tuning

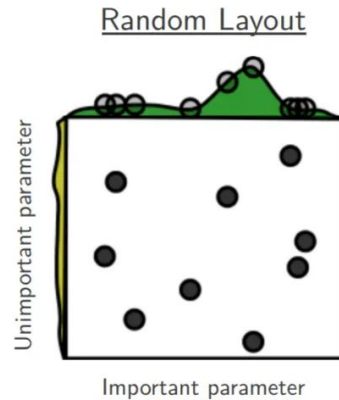
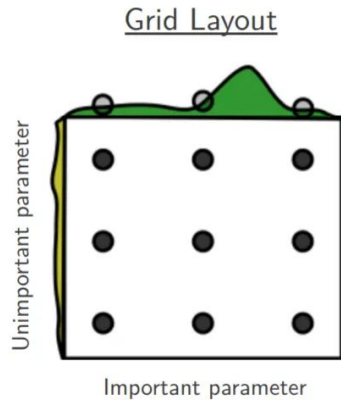
1. Choose a model and a set of relevant hyperparameters
2. Define an evaluation function (like F1 or MSE)
3. Iterate through hyperparameters
4. For each combination of hyperparameter values, evaluate model performance (using step 2!)

Main steps of hyperparameter tuning

1. Choose a model and a set of relevant hyperparameters
2. Define an evaluation function (like F1 or MSE)
3. Iterate through hyperparameters
4. For each combination of hyperparameter values, evaluate model performance (using step 2!)
5. Select combination of hyperparameters that resulted in best performance

Main steps of hyperparameter tuning

1. Choose a model and a set of relevant hyperparameters
2. Define an evaluation function (like F1 or MSE)
3. **Iterate through hyperparameters**
4. For each combination of hyperparameter values, evaluate model performance (using step 2!)
5. Select combination of hyperparameters that resulted in best performance



Most common hyperparameter tuning procedures:

- **Grid search:** evaluate every combination of hyperparameter values
- **Random search:** evaluate hyperparameter values drawn from a distribution

Demo: Neural Network Playground

<https://playground.tensorflow.org/>