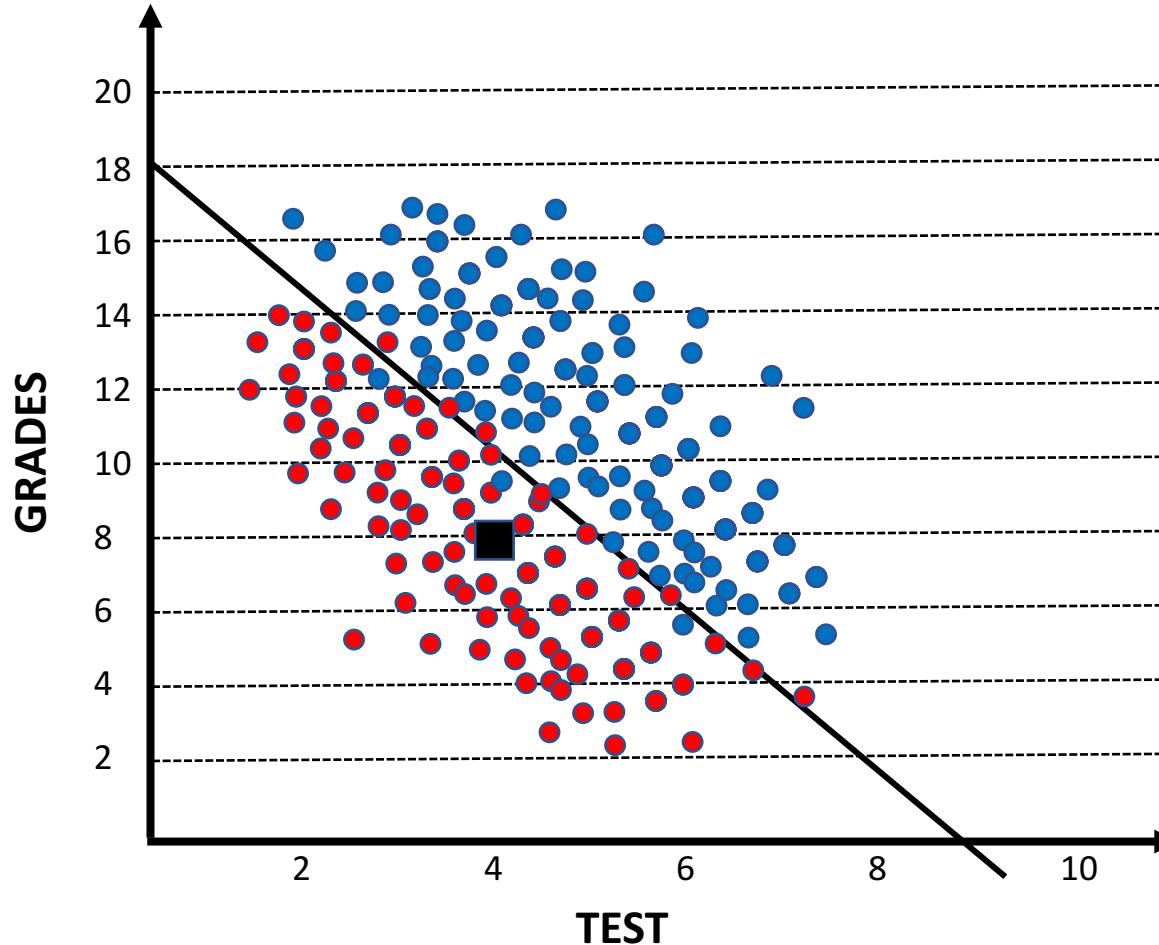


# Perceptron

# Linear Boundary / Classification



**BOUNDARY: A LINE ( $WX + B$ )**

$$2X_1 + X_2 - 18 = 0$$

$X_1$  = Test

$X_2$  = Grade

$$2 \cdot \text{Test} + \text{Grade} - 18 = \text{Score}$$

**PREDICTION:**

Score  $\geq 0$  Accept

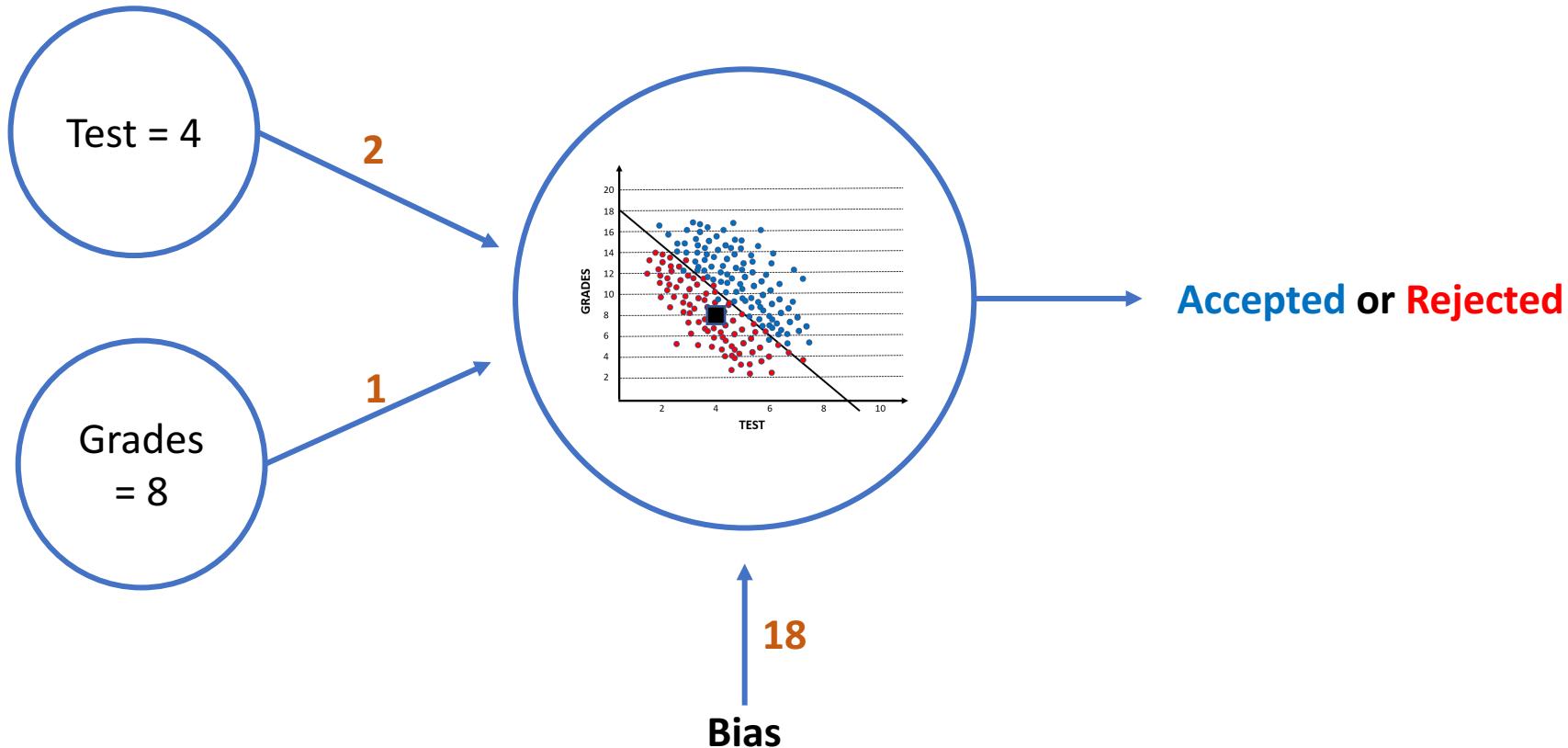
Score  $< 0$  Reject

**Example:** What is prediction of a student who got 4 on the test and 8 on the grades?

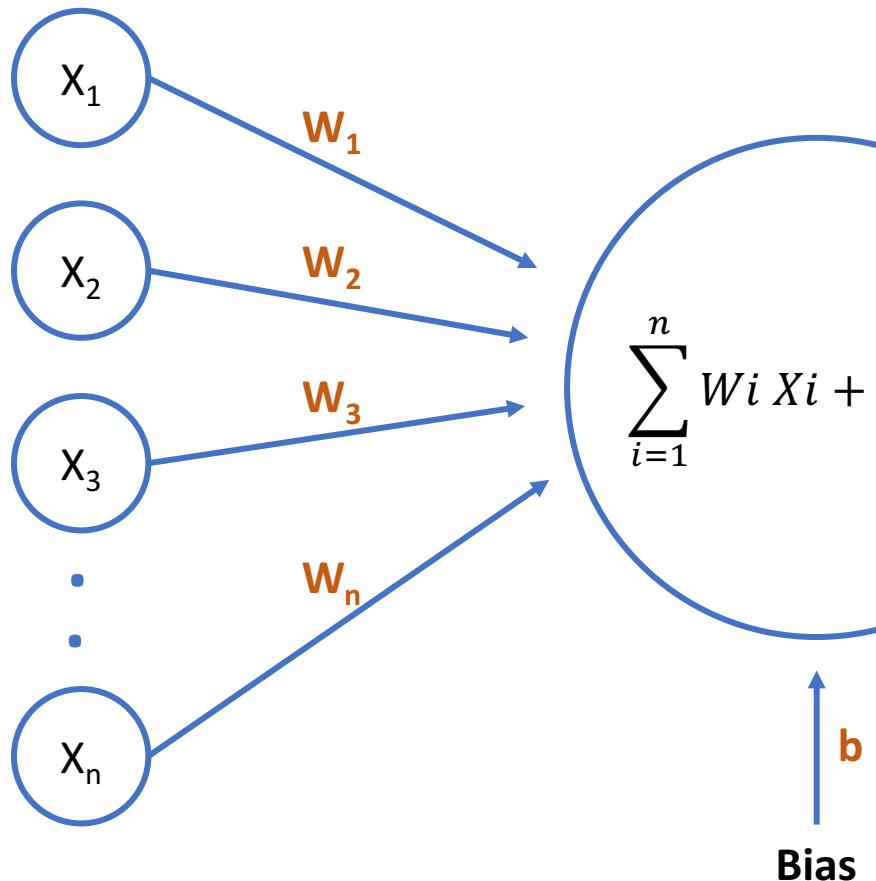
Answer:  $2 \times 4 + 8 - 18 = 16 - 18 = -2$

# Perceptron

$$2X_1 + X_2 - 18 = 0$$

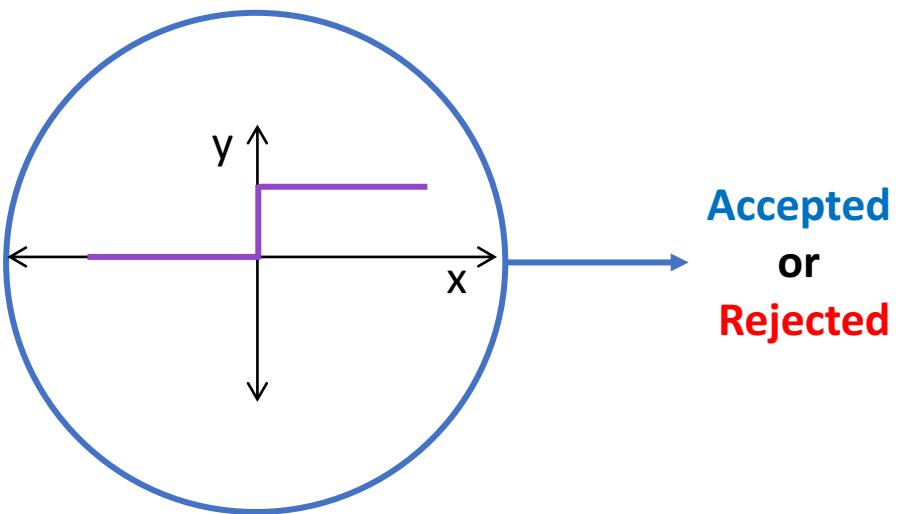
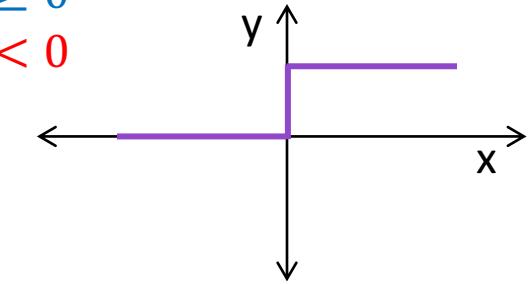


# Perceptron

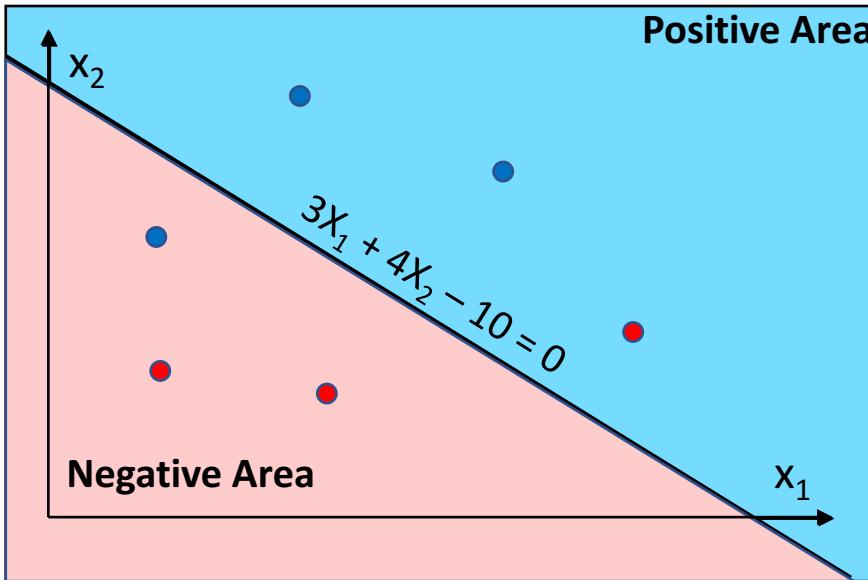


Step Function

$$Y = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$



# How to select the weights and bias – called TRAINING the Perceptron



$$3X_1 + 4X_2 - 10 = 0$$

Misclassified Point: (8,3)

-	3	4	-10
8	3	1	
-	5	1	-11

$$3X_1 + 4X_2 - 10 = 0$$

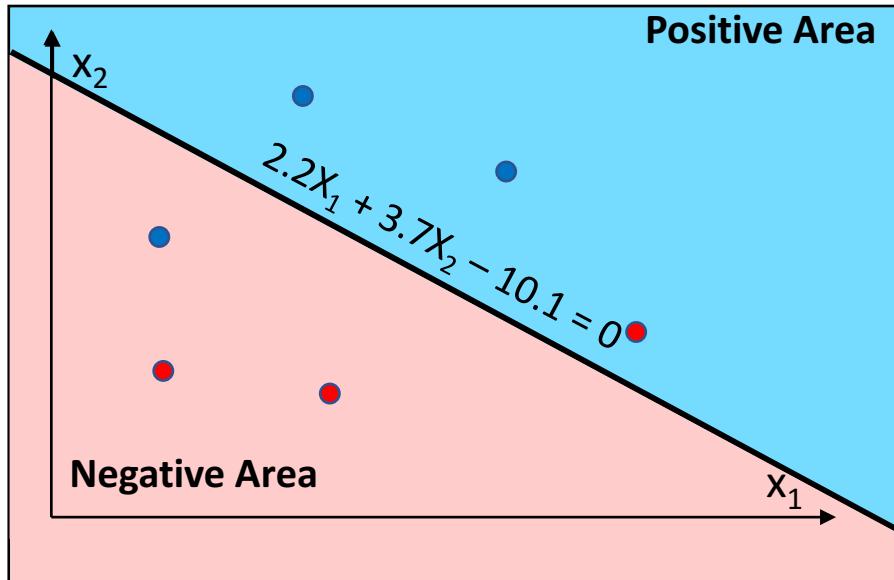
Misclassified Point: (8,3)

**Learning Rate = 0.1**

-	3	4	-10
8x0.1	3x0.1	1x0.1	
-	2.2	3.7	-10.1

New Line:  $2.2X_1 + 3.7X_2 - 10.1 = 0$

# How to select the weights and bias – called TRAINING the Perceptron



$$3X_1 + 4X_2 - 10 = 0$$

Misclassified Point: (8,3)

-	3	4	-10
8	3	1	
			-11
-	5	1	

$$3X_1 + 4X_2 - 10 = 0$$

Misclassified Point: (8,3)

**Learning Rate = 0.1**

-	3	4	-10
	$8 \times 0.1$	$3 \times 0.1$	$1 \times 0.1$
			-10.1
-	2.2	3.7	

New Line:  $2.2X_1 + 3.7X_2 - 10.1 = 0$

# How to select the weights and bias – called TRAINING the Perceptron

$$2.2X_1 + 3.7X_2 - 10.1 = 0$$

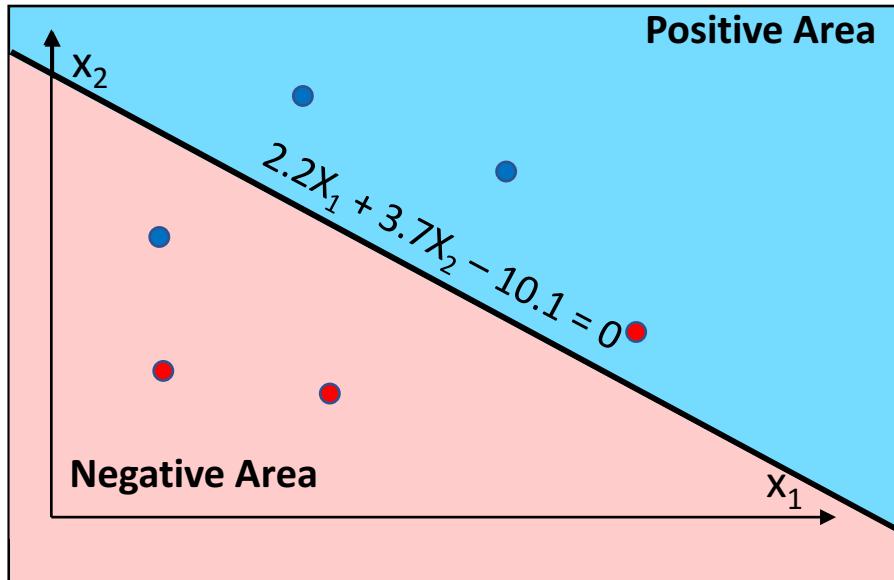
Misclassified Point: (2,6)

Learning Rate = 0.1

$$\begin{array}{r} + \\ \begin{array}{rrr} 2.2 & 3.7 & -10.1 \\ 2 \times 0.1 & 6 \times 0.1 & 1 \times 0.1 \\ \hline 2.4 & 4.3 & -10 \end{array} \end{array}$$

New Line:

$$2.4X_1 + 4.3X_2 - 10 = 0$$



# How to select the weights and bias – called TRAINING the Perceptron

$$2.2X_1 + 3.7X_2 - 10.1 = 0$$

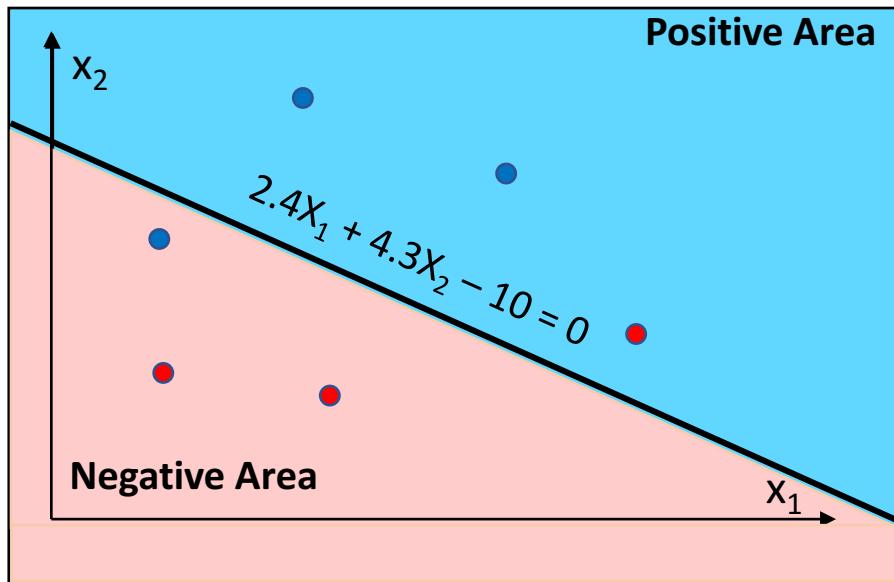
Misclassified Point: (2,6)

Learning Rate = 0.1

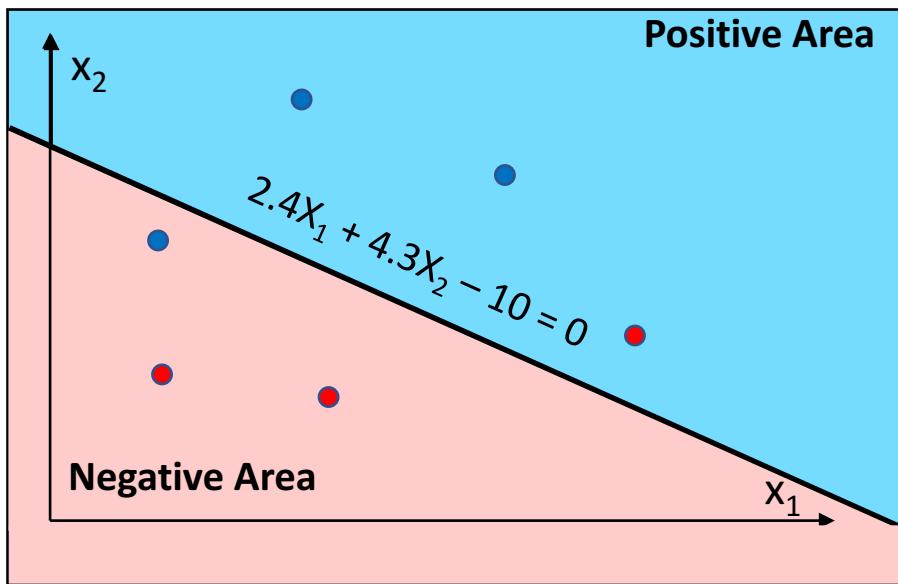
$$\begin{array}{r} + \\ \begin{array}{rrr} 2.2 & 3.7 & -10.1 \\ 2 \times 0.1 & 6 \times 0.1 & 1 \times 0.1 \\ \hline 2.4 & 4.3 & -10 \end{array} \end{array}$$

New Line:

$$2.4X_1 + 4.3X_2 - 10 = 0$$



# Pseudocode for Training – Adjusting the Weights & Bias



1. Start with random weights:  $(W_1, \dots, W_n, b)$
2. For every misclassified point  $(X_1, \dots, X_n)$ :
  - If **Prediction = 0**
    - For  $i = 1 \dots n$ 
      - $W_i = W_i + \alpha X_i$
      - $b = b + \alpha$
    - If **Prediction = 1**
      - For  $i = 1 \dots n$ 
        - $W_i = W_i - \alpha X_i$
        - $b = b - \alpha$

# Perceptron – Basic Functions

- Step Function

$$Y' = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

- Predict Function

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

- Error Function

$$Error = y - \hat{y}$$

- Update Weights Function

$$b \leftarrow b + \alpha \cdot x \cdot Error$$

$$w_i \leftarrow w_i + \alpha \cdot x_i \cdot Error$$

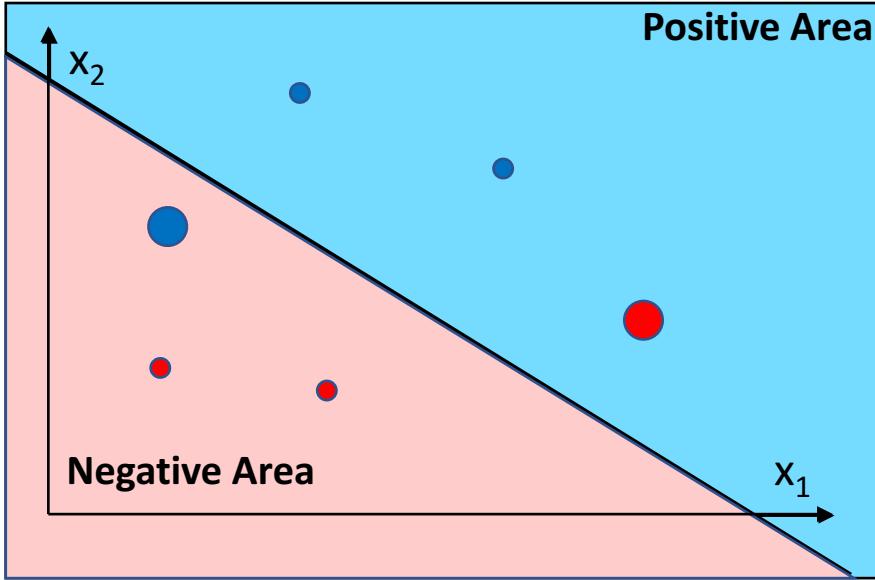
# Run Code

- `perceptron_1.py`
- How can we initialize **weights** and **bias**?
- How **learning rate** effects on model training?

# Run Code

- **perceptron\_2.py**
- What is **epoch**?
- What does epochs=100 do?
- Is our model better when we have high epoch?

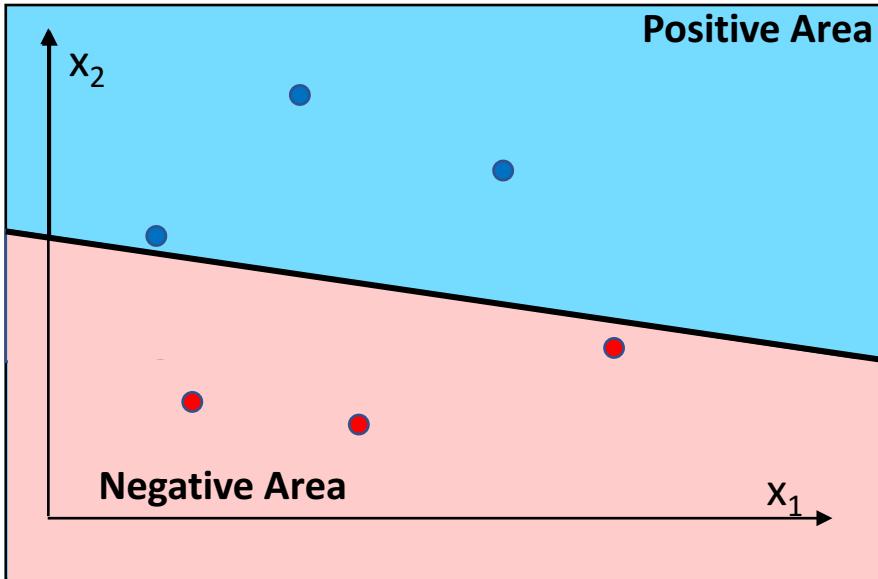
# Error Function & Log-Loss Function



$$\text{Error} = \bullet + \bullet + \textcolor{blue}{\bullet} + \textcolor{red}{\bullet} + \bullet + \bullet$$

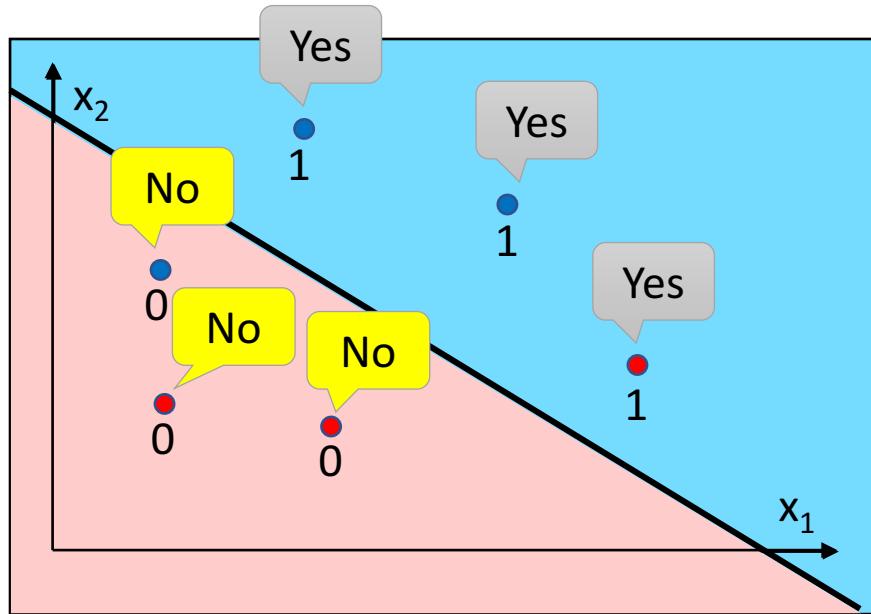
- What happened if error function is **discrete**?
- Error function should be **continuous** and **differentiable**

# Error Function & Log-Loss Function

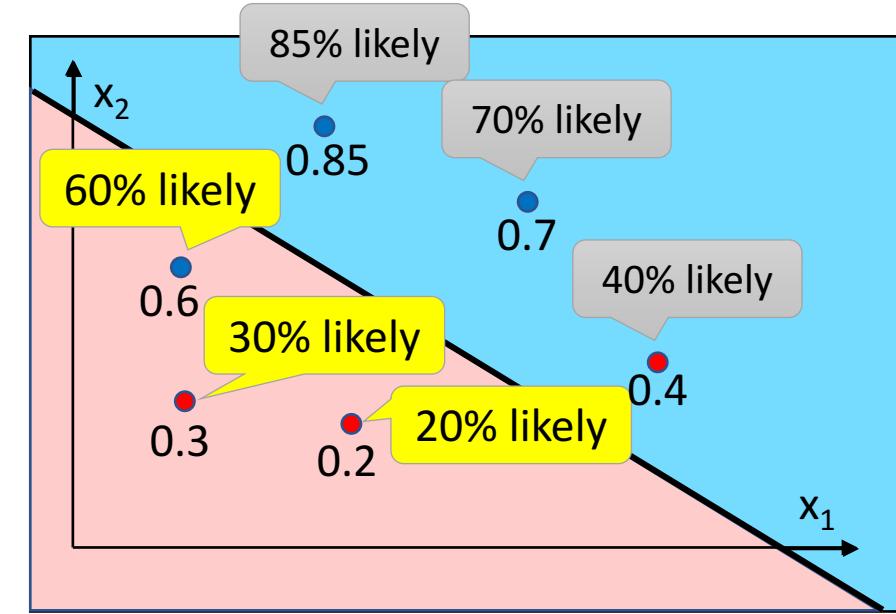


$$\text{Error} = \bullet + \bullet + \bullet + \bullet + \bullet + \bullet + \bullet$$

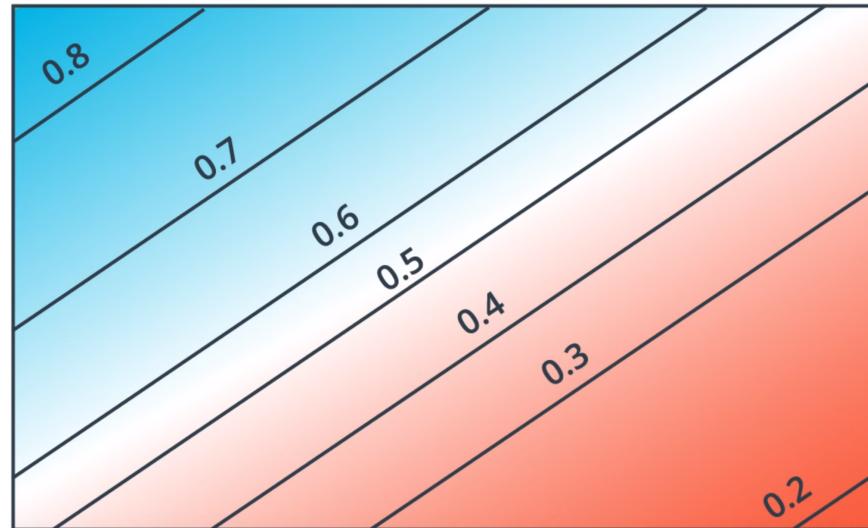
# Predictions



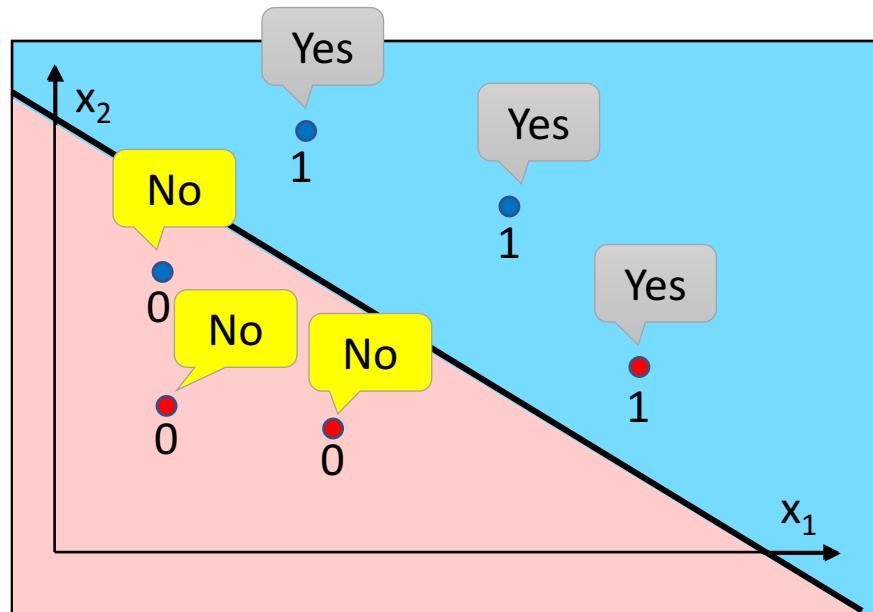
Discrete



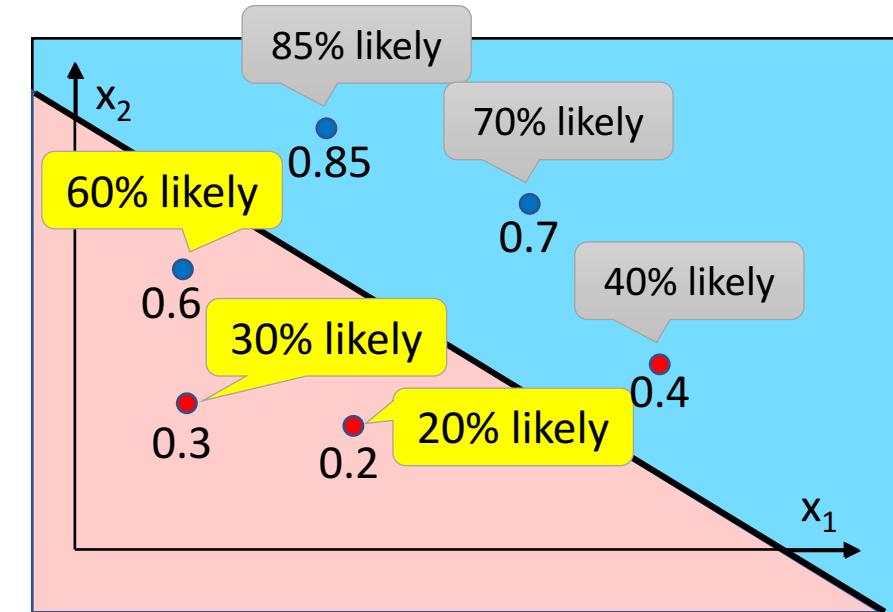
Continuous



# Predictions & Activation Functions

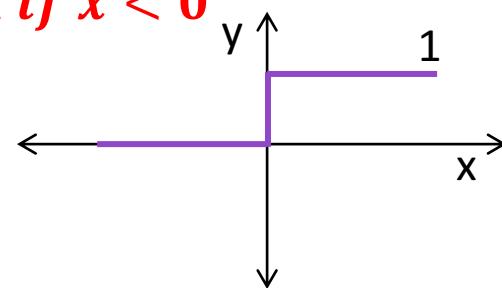


Discrete

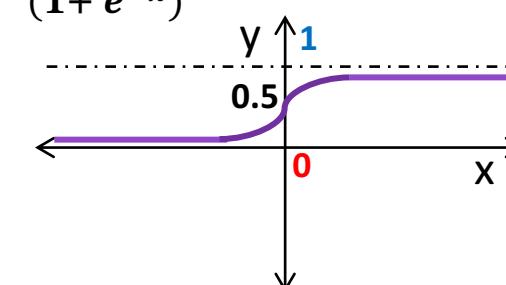


Continuous

Step Function  $Y = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$



Sigmoid Function  $Y = \frac{1}{(1 + e^{-x})}$



# Predictions & Activation Functions

## Binary Classification

Sigmoid Function  $Y = \frac{1}{(1 + e^{-x})}$

Animal	Value
	1
	0

## Multi-Class Classification

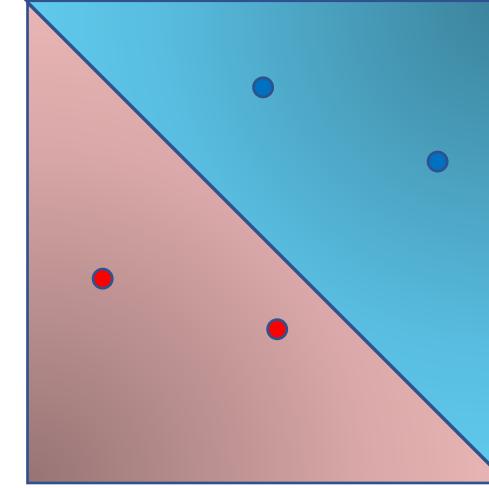
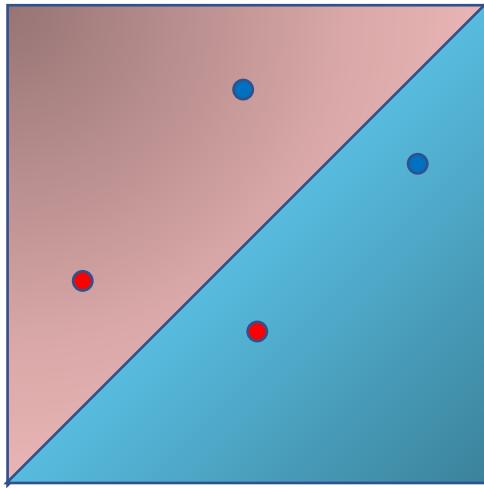
Softmax Function  $= \frac{e^{Z_i}}{e^{Z_1} + e^{Z_2} + \dots + e^{Z_n}}$

## One-Hot Encoding

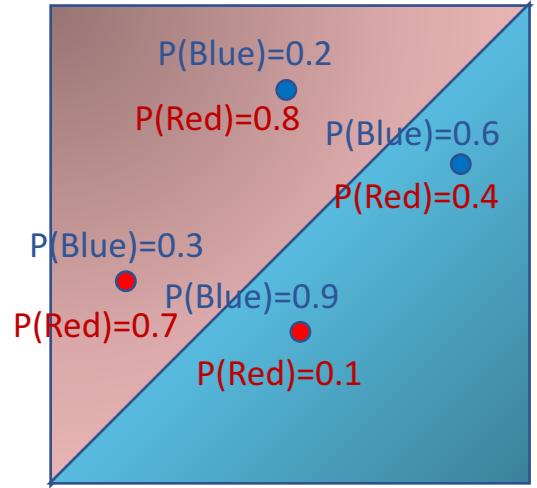
Animal	Dog	Cat	Duck
	1	0	0
	0	1	0
	0	0	1

# Maximum Likelihood

Which model is better?



# Maximum Likelihood

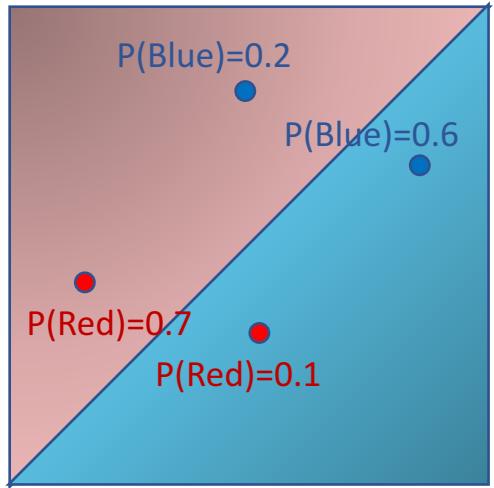


$$\hat{Y} = \sigma(Wx + b)$$

$$P(\text{Blue}) = \sigma(Wx + b)$$

$$P(\text{Red}) = 1 - P(\text{Blue})$$

# Maximum Likelihood



X

P(Blue)=0.2  
P(Blue)=0.6  
P(Red)=0.7  
P(Red)=0.1

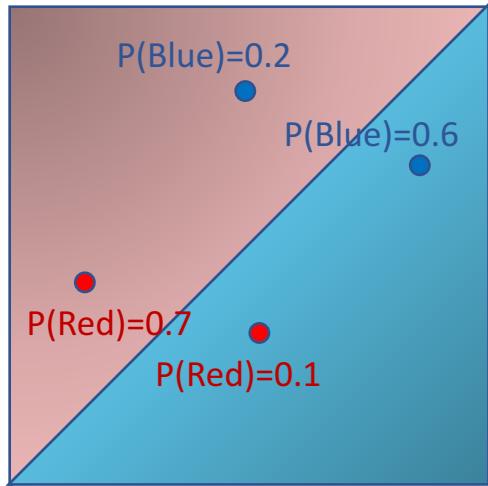
---

P(all)=0.0084

$$\hat{Y} = \sigma(Wx + b)$$
$$P(\text{Blue}) = \sigma(Wx + b)$$
$$P(\text{Red}) = 1 - P(\text{Blue})$$

# Maximum Likelihood

Goal: Maximize the Probability



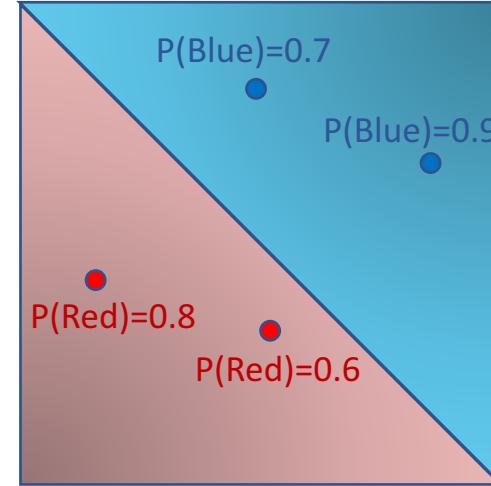
$$0.6 \times 0.2 \times 0.1 \times 0.7 = 0.0084$$

$$\ln(0.6) + \ln(0.2) + \ln(0.1) + \ln(0.7)$$

$$(-0.51) + (-1.61) + (-2.3) + (-0.36) = -4.78$$

$$-\ln(0.6) - \ln(0.2) - \ln(0.1) - \ln(0.7)$$

$$0.51 + 1.61 + 2.3 + 0.36 = 4.78$$



$$0.7 \times 0.9 \times 0.8 \times 0.6 = 0.3024$$

$$\ln(0.7) + \ln(0.9) + \ln(0.8) + \ln(0.6)$$

$$(-0.36) + (-0.1) + (-0.22) + (-0.51) = -1.19$$

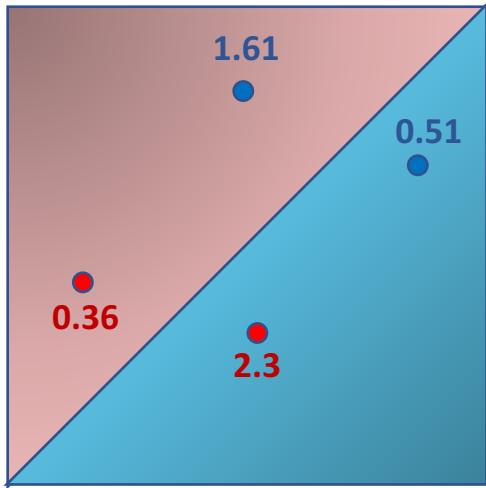
$$-\ln(0.2) - \ln(0.6) - \ln(0.7) - \ln(0.1)$$

$$0.36 + 0.1 + 0.22 + 0.51 = 1.19$$

Cross Entropy

# Maximum Likelihood

**Goal: Maximize the Probability by minimize the cross entropy**



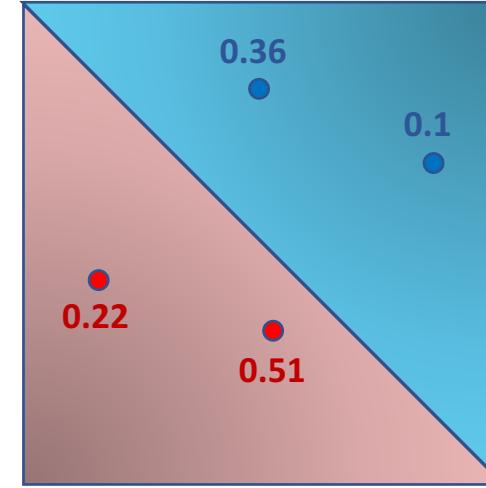
$$0.6 \times 0.2 \times 0.1 \times 0.7 = 0.0084$$

$$\ln(0.6) + \ln(0.2) + \ln(0.1) + \ln(0.7)$$

$$(-0.51) + (-1.61) + (-2.3) + (-0.36) = -4.78$$

$$-\ln(0.6) - \ln(0.2) - \ln(0.1) - \ln(0.7)$$

$$0.51 + 1.61 + 2.3 + 0.36 = 4.78$$



$$0.7 \times 0.9 \times 0.8 \times 0.6 = 0.3024$$

$$\ln(0.7) + \ln(0.9) + \ln(0.8) + \ln(0.6)$$

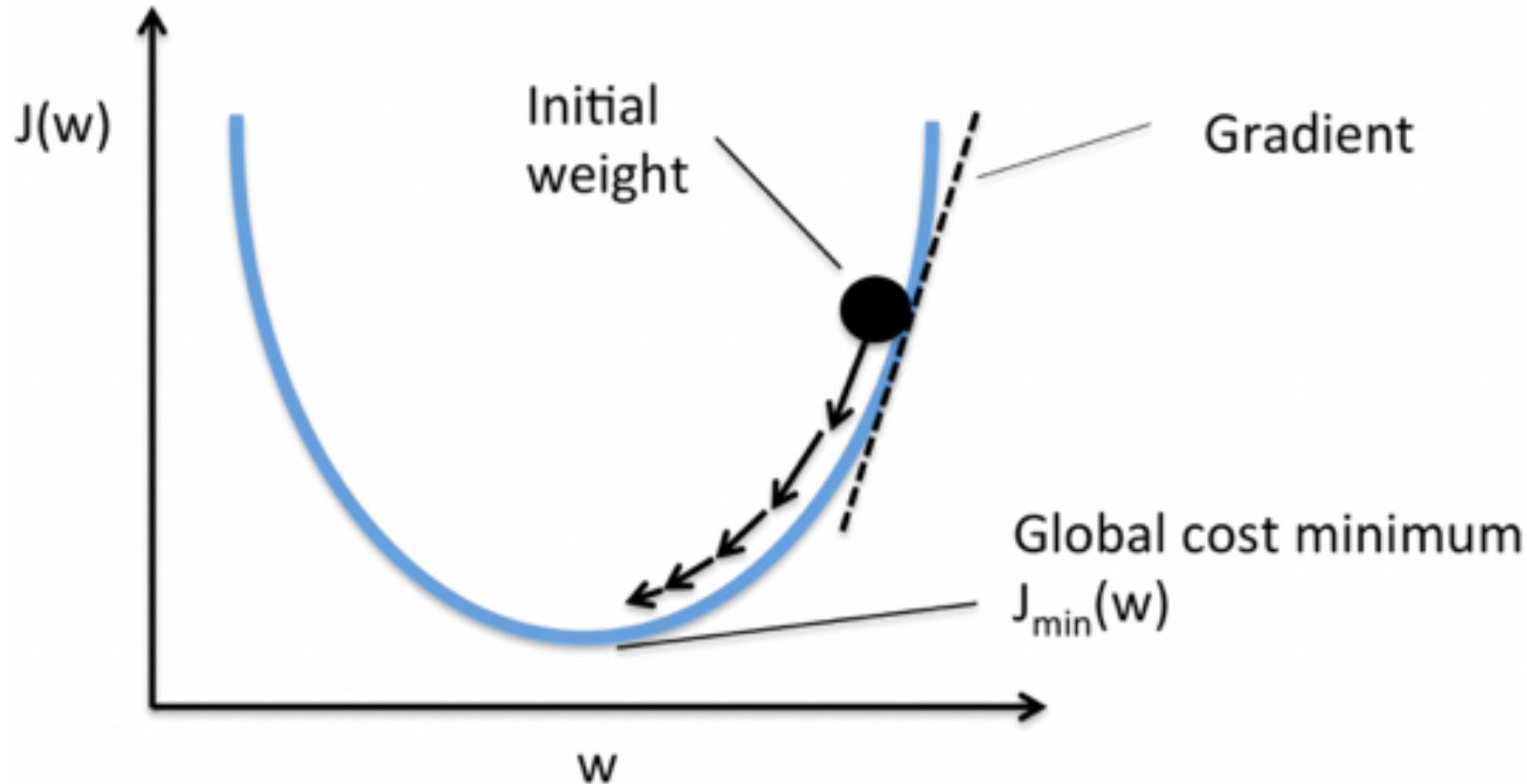
$$(-0.36) + (-0.1) + (-0.22) + (-0.51) = -1.19$$

$$-\ln(0.2) - \ln(0.6) - \ln(0.7) - \ln(0.1)$$

$$0.36 + 0.1 + 0.22 + 0.51 = 1.19$$

# Gradient Decent

What are we optimizing?



# Gradient Descent – Basic Functions

- Sigmoid Function

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

- Predict Function

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

- Error Function (Cross Entropy)

$$Error(y, \hat{y}) = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

- Update Weights Function

$$b \leftarrow b + \alpha \cdot Error(y, \hat{y})$$

$$w_i \leftarrow w_i + \alpha \cdot x_i \cdot Error(y, \hat{y})$$

# Run Code

- `gradient_descent.py`
- How gradient descent different from perceptron?

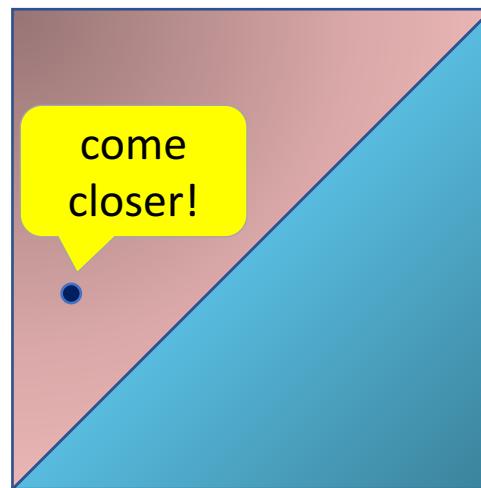
# Run Code

- **perceptron\_keras\_scikit\_learn.py**
- How can we evaluate the model with cross validation?

# Perceptron vs Gradient Descent

## Perceptron

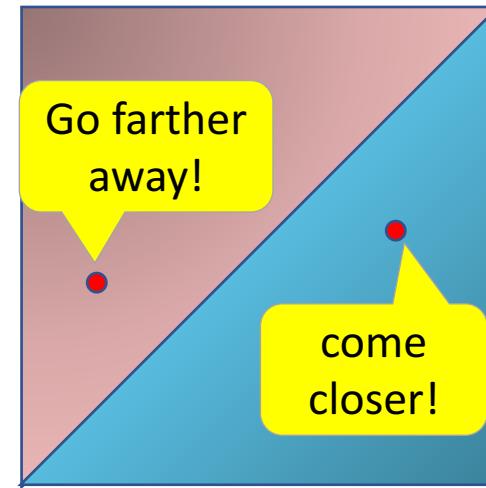
- Use step activation function
- $\hat{y}$  value is only 0 or 1
- If  $x_i$  is misclassified then only change its weight  $w_i$



Incorrectly classified

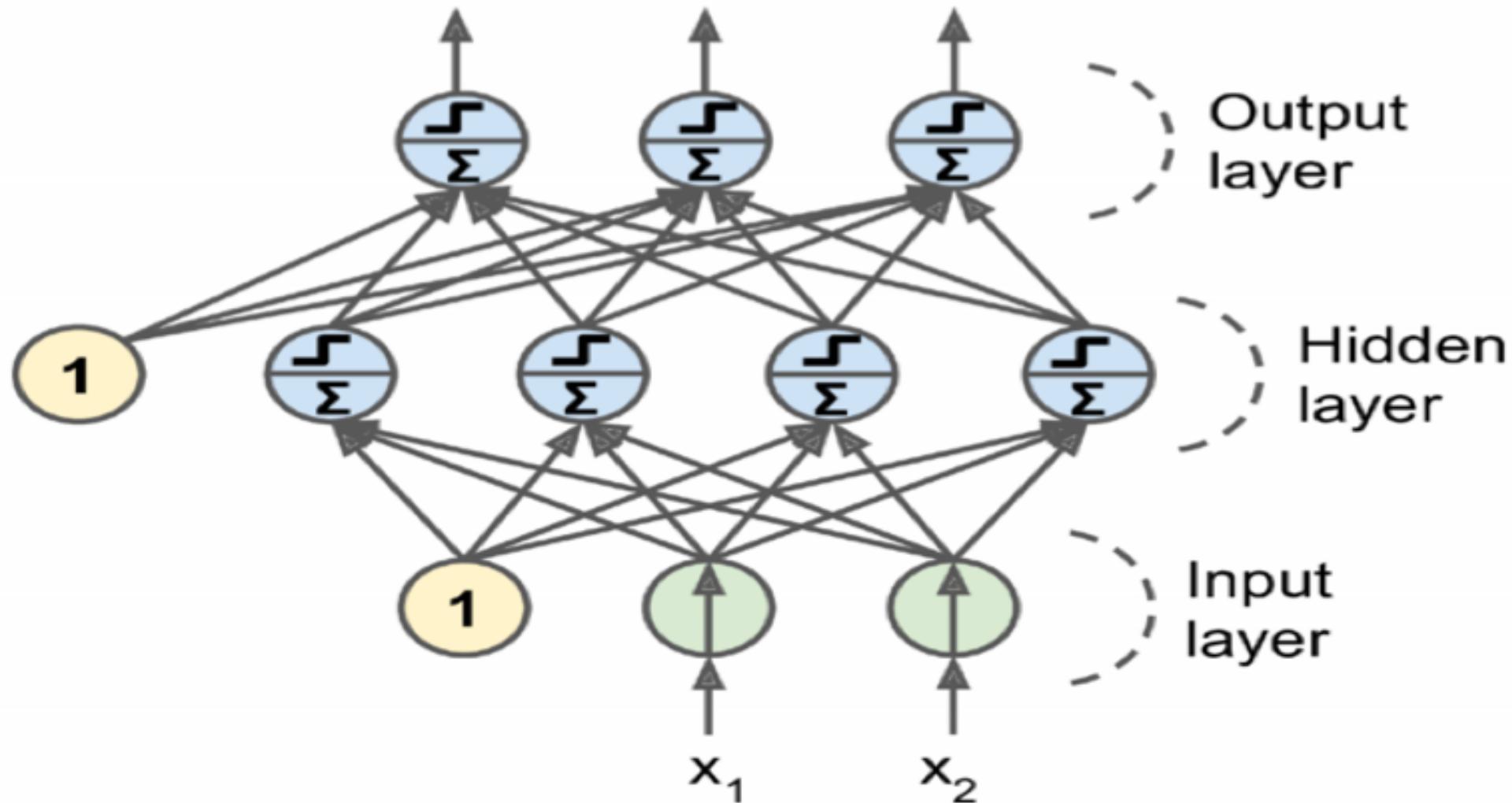
## Gradient Descent

- Use sigmoid activation function
- $\hat{y}$  can be any value between 0 and 1
- Change all weights as  $b \leftarrow b + \alpha \cdot Error(y, \hat{y})$



Correctly classified

# Two Layers of Perceptron



# Cross Entropy

## Binary Classification

- Cross-Entropy =  $-\sum_{i=1}^m y_i \ln P_i + (1 - y_i) \ln(1 - P_i)$
- Keras has “**binary\_crossentropy**” function

## Multiclass Classification

- Cross-Entropy =  $-\sum_{i=1}^n \sum_{j=1}^m y_{ij} \ln P_{ij}$
- Keras has “**categorical\_crossentropy**” function.

# Run Code

- **gradient\_descent\_2.py**
  - What is the data type of X\_train?
  - What does **Flatten** do?
  - What does activation=“**softmax**” do?
- Why use a loss function=“**category\_crossentropy**” at all?
  - What does optimizer=“**adam**” do?
  - How can we get model performance (metrics)?
  - How can we save best model (**ModelCheckPoint**)?

# Run Code

- `gradient_descent_load.py`
- How can we load saved model?

# Run Code

- **model\_saving.py**
- How to save a model as json file?
- How to save a model as yaml file?

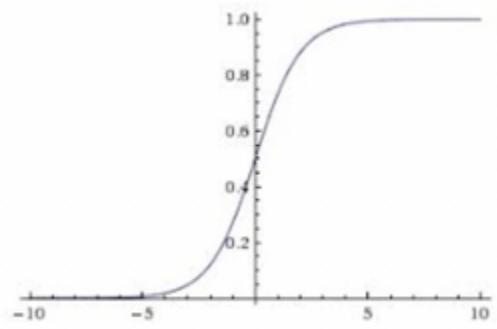
# Run Code

- `gradient_descent_predict.py`
- How can we make inference on unseen data?

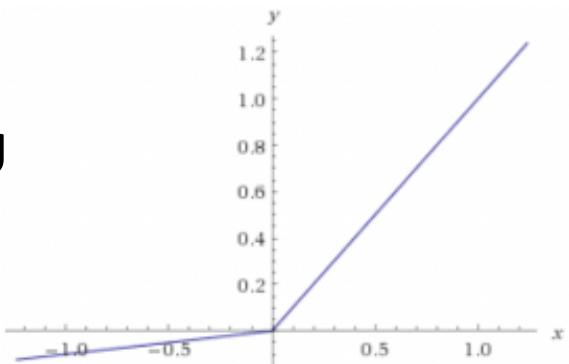
# Multi Layers Networks

# Activation Functions

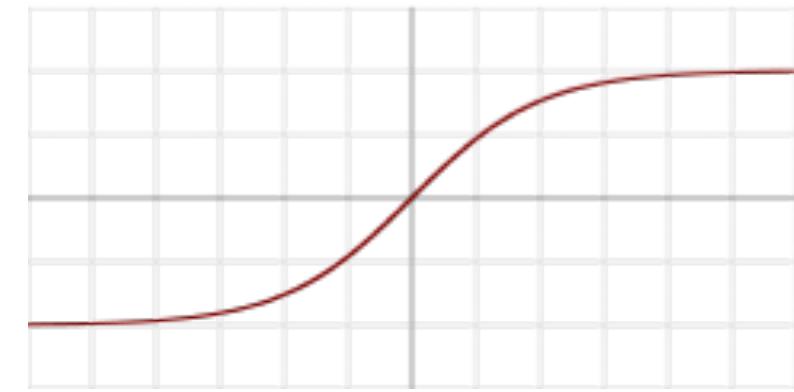
Sigmoid



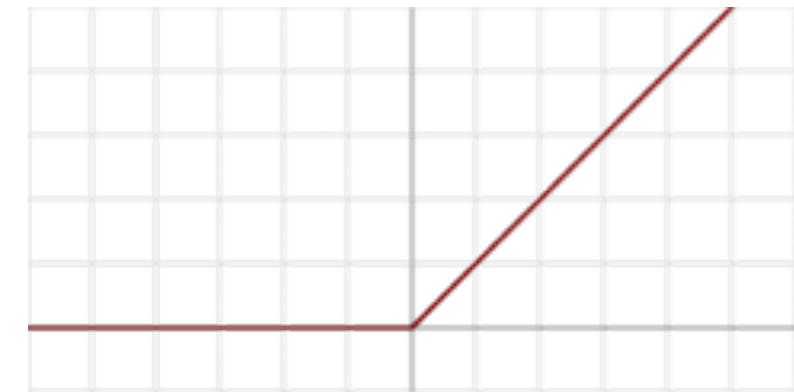
Leaky ReLU



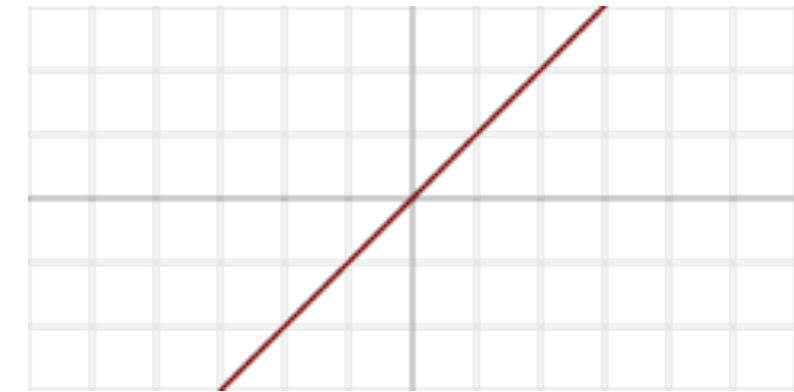
Tanh



ReLU



Linear



# Loss Functions

- Binary Cross Entropy
- Categorical Cross Entropy
- Sparse Categorical Cross Entropy
- Mean Squared Error – MSE
- Mean Absolute Error – MAE
- Mean Absolute Percentage Error – MAPE
- Mean Squared Logarithmic Error – MSLE

# Optimization Functions

- SGD – Stochastic Gradient Descent optimizer
  - Learning rate
  - Momentum
  - Decay
  - Nesterov momentum
- RMSProp
- Adagrad
- Adadelta
- Adam
- Adamax
- Nadam
- TFOptimize

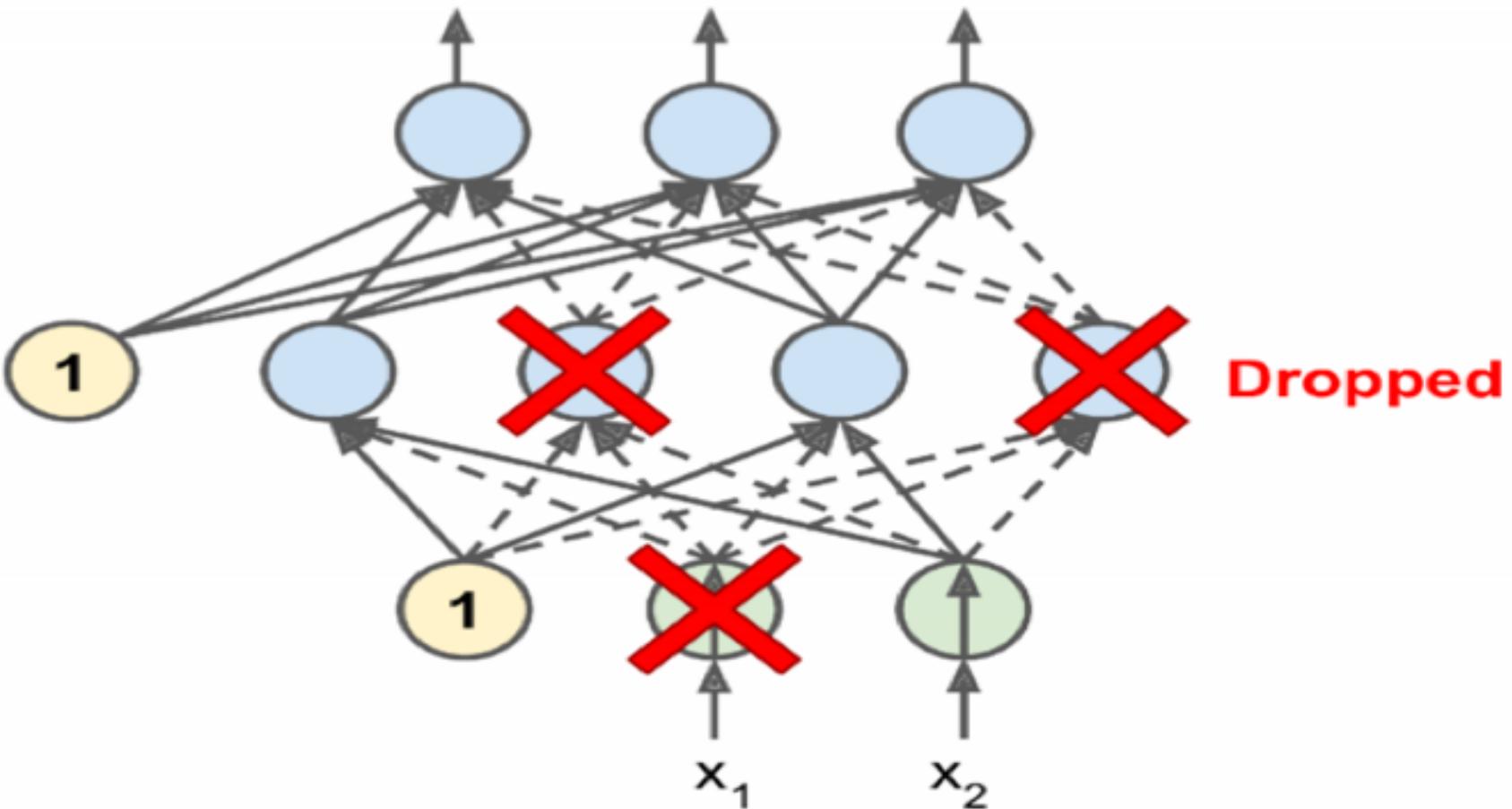
# Metrics

- Accuracy
- MAE – Mean Absolute Error

# Run Code

- `gradient_descent_multilayers.py`
- How did this change run-time and accuracy?
- Why do we use activation=“`relu`” in one layer?
  - How could we make a deeper model?

# Dropout



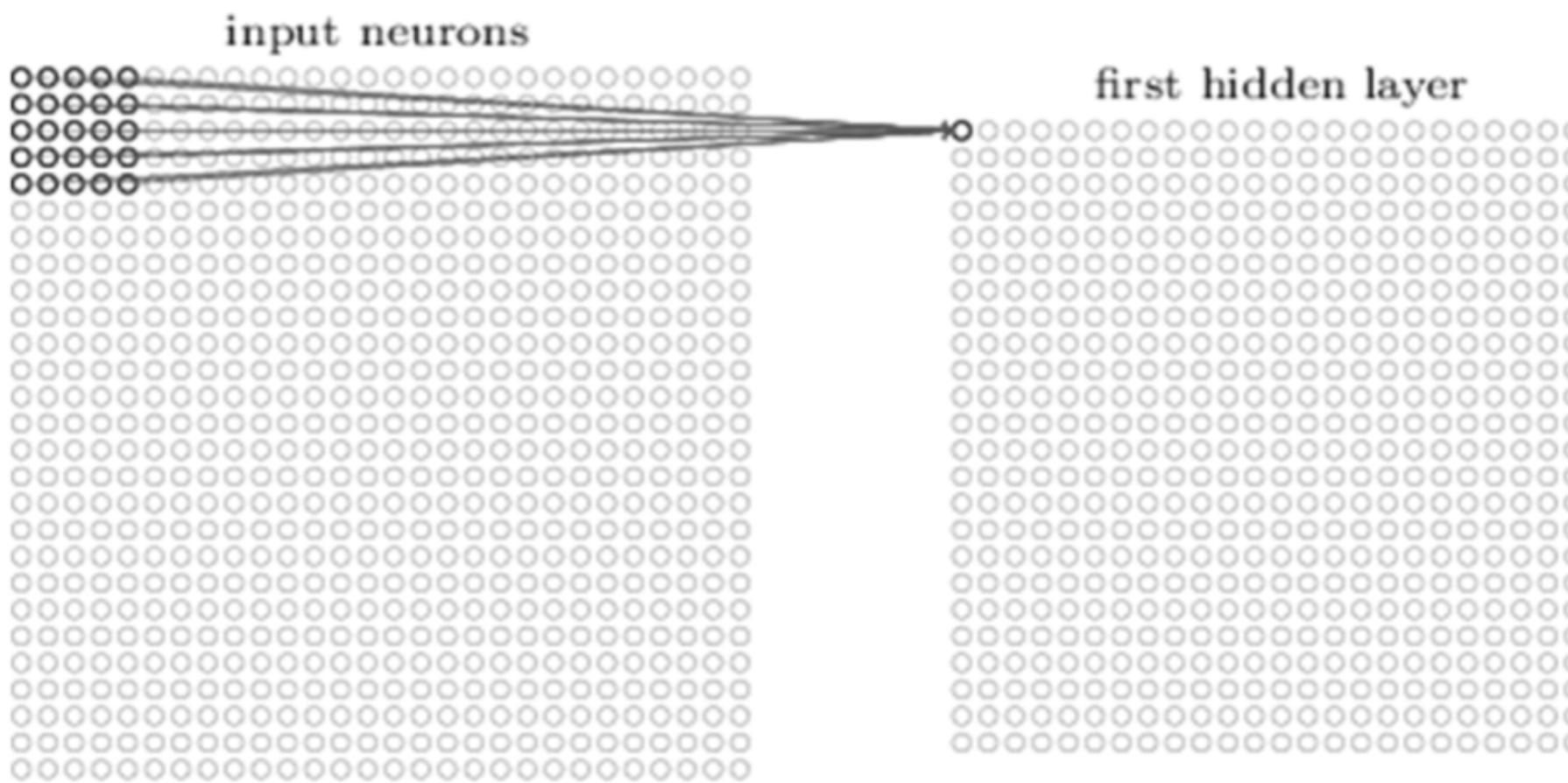
# Run Code

- **gradient\_descent\_dropout.py**
- How does this change training accuracy and validation accuracy?
  - When would we consider using dropout

# Convolutional Neural Net

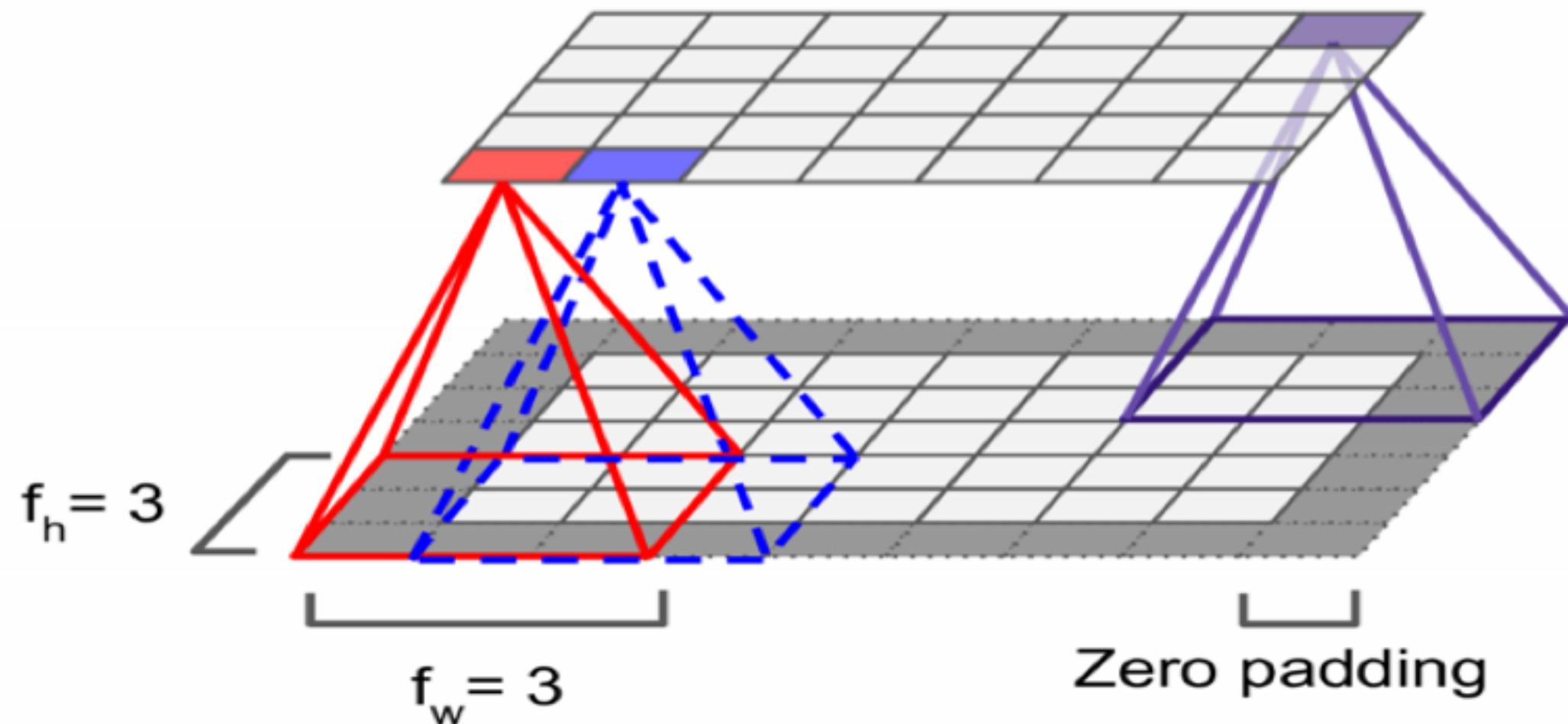
## CNN

# Convolution



Visualization of  $5 \times 5$  filter convolving around an input volume and producing an activation map

# Convolutions



# Run Code

- **convolution-demo.py**
- What happens if all the numbers are zeros?
- What happens if the numbers are large?
- What convolution would leave the image unaffected?

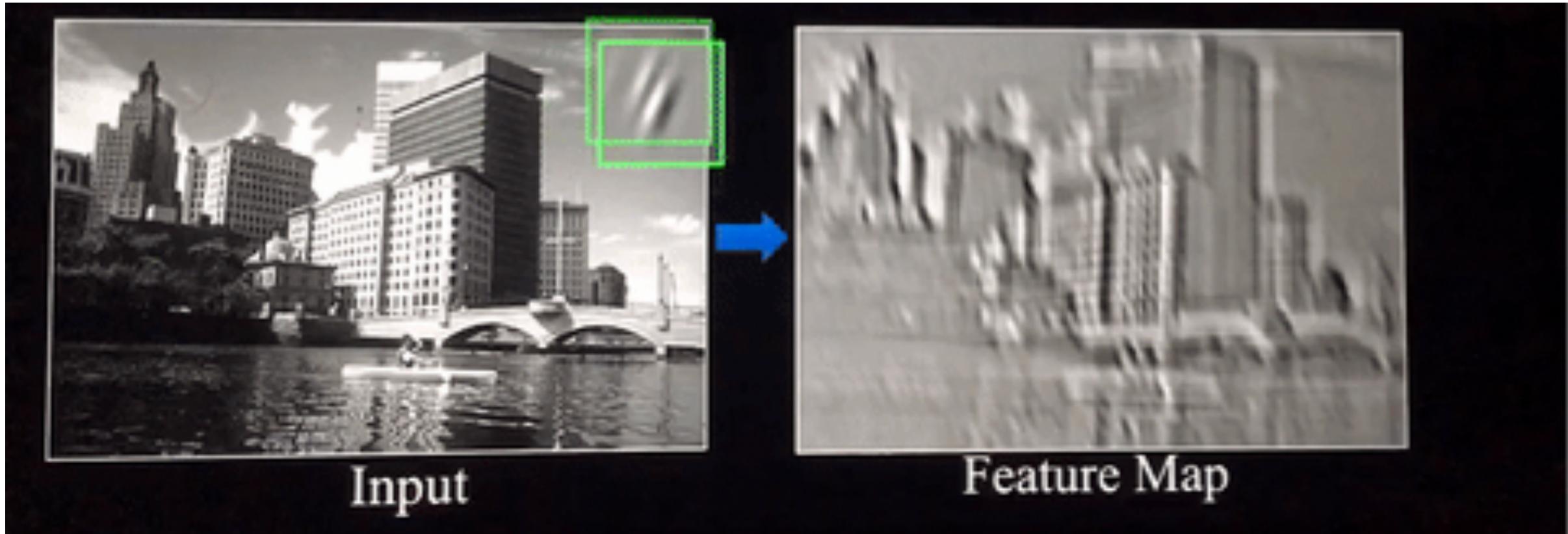
# Convolutions



Horizontal edge

Vertical edge

# Convolutions

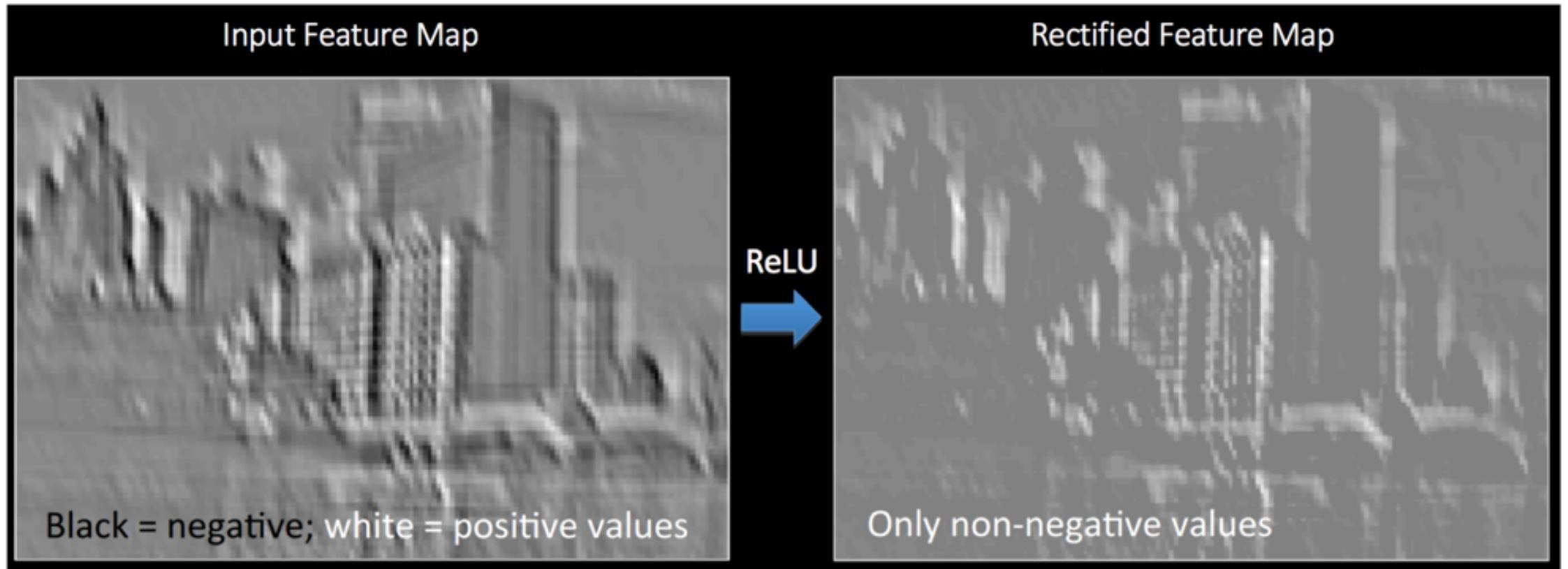


Input

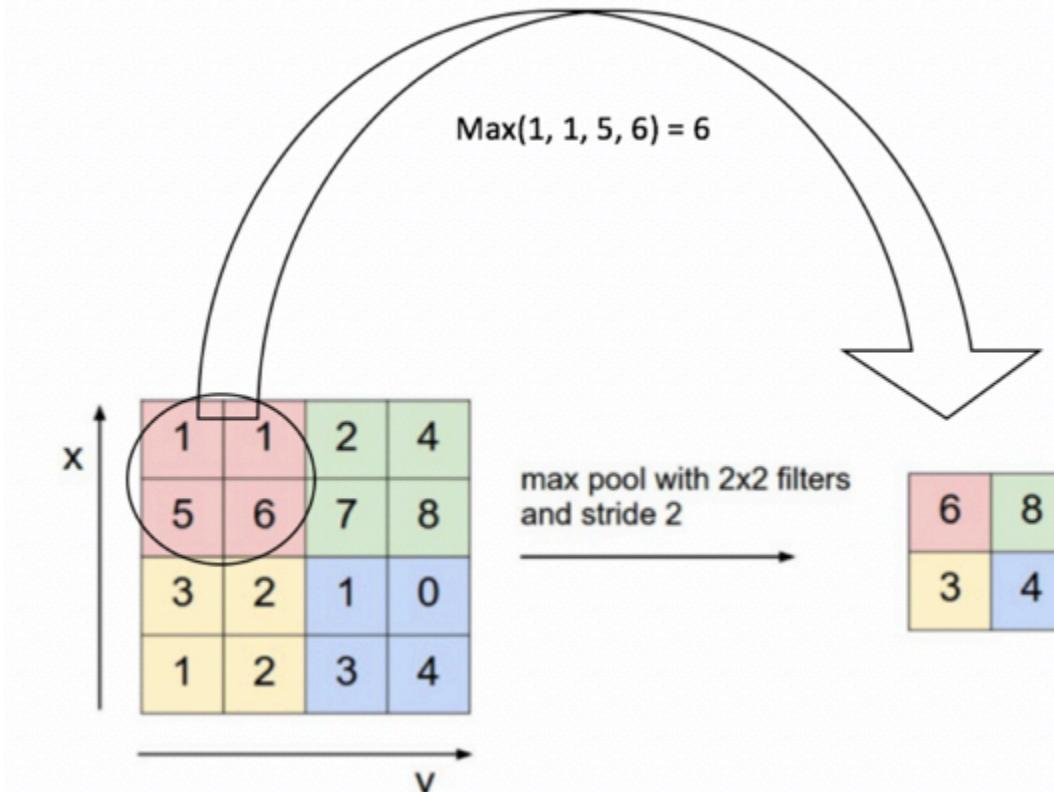
Feature Map

# Activation Function

## ReLU Activation



# Max Pooling



# Run Code

- `maxpool_demo.py`
- What is happening here?

# Max Pooling



3 x 3



5 x 5

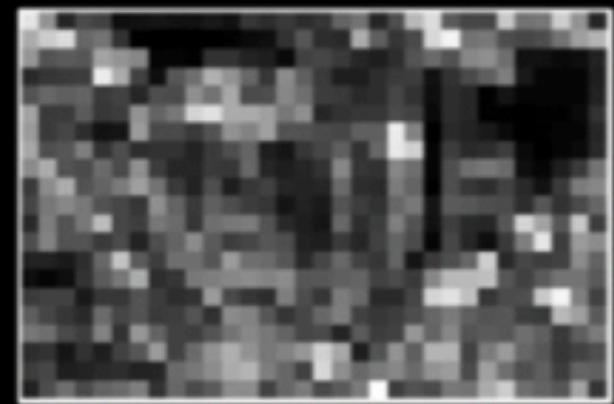
# Max Pooling



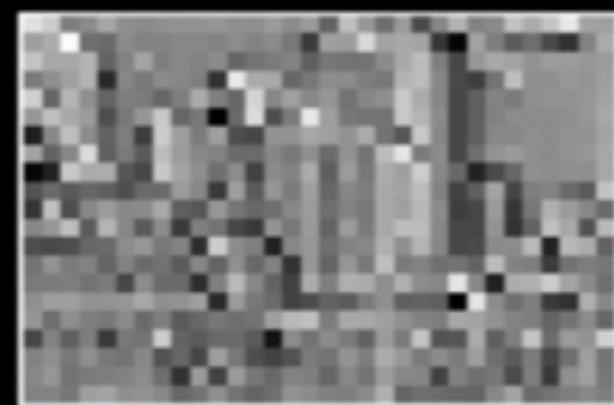
Rectified Feature Map

Pooling  
→

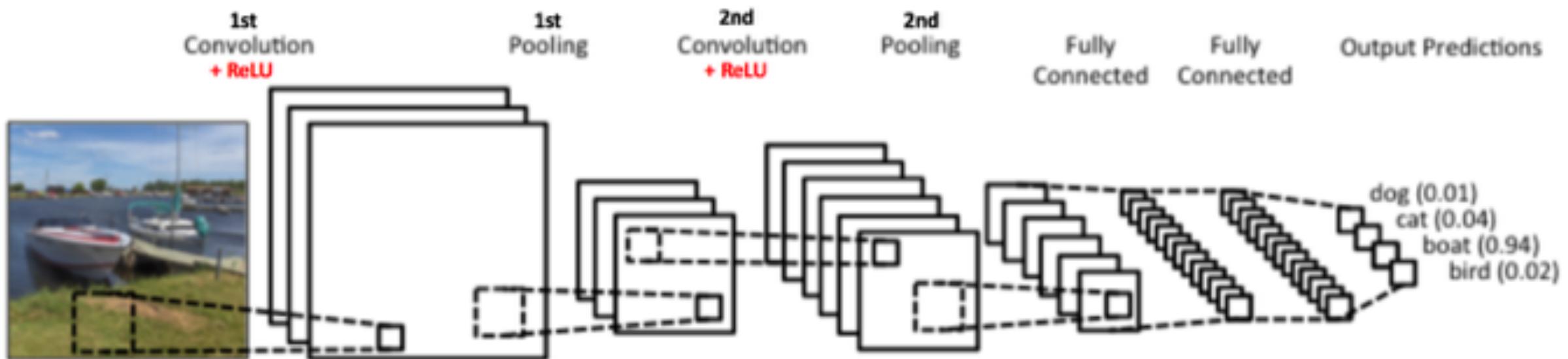
Max



Sum



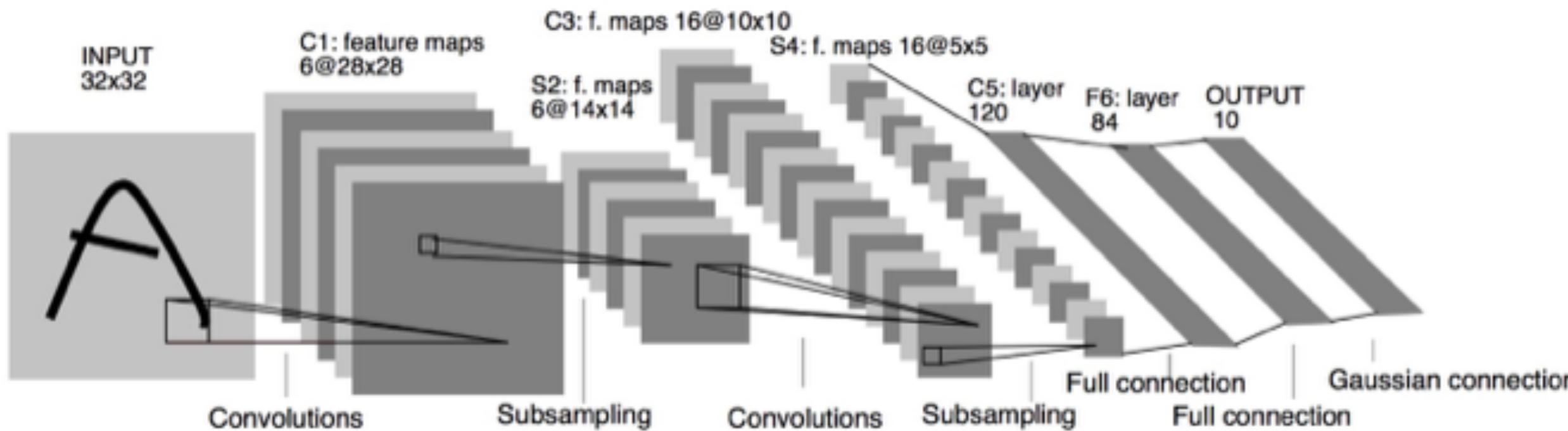
# CNN Architecture



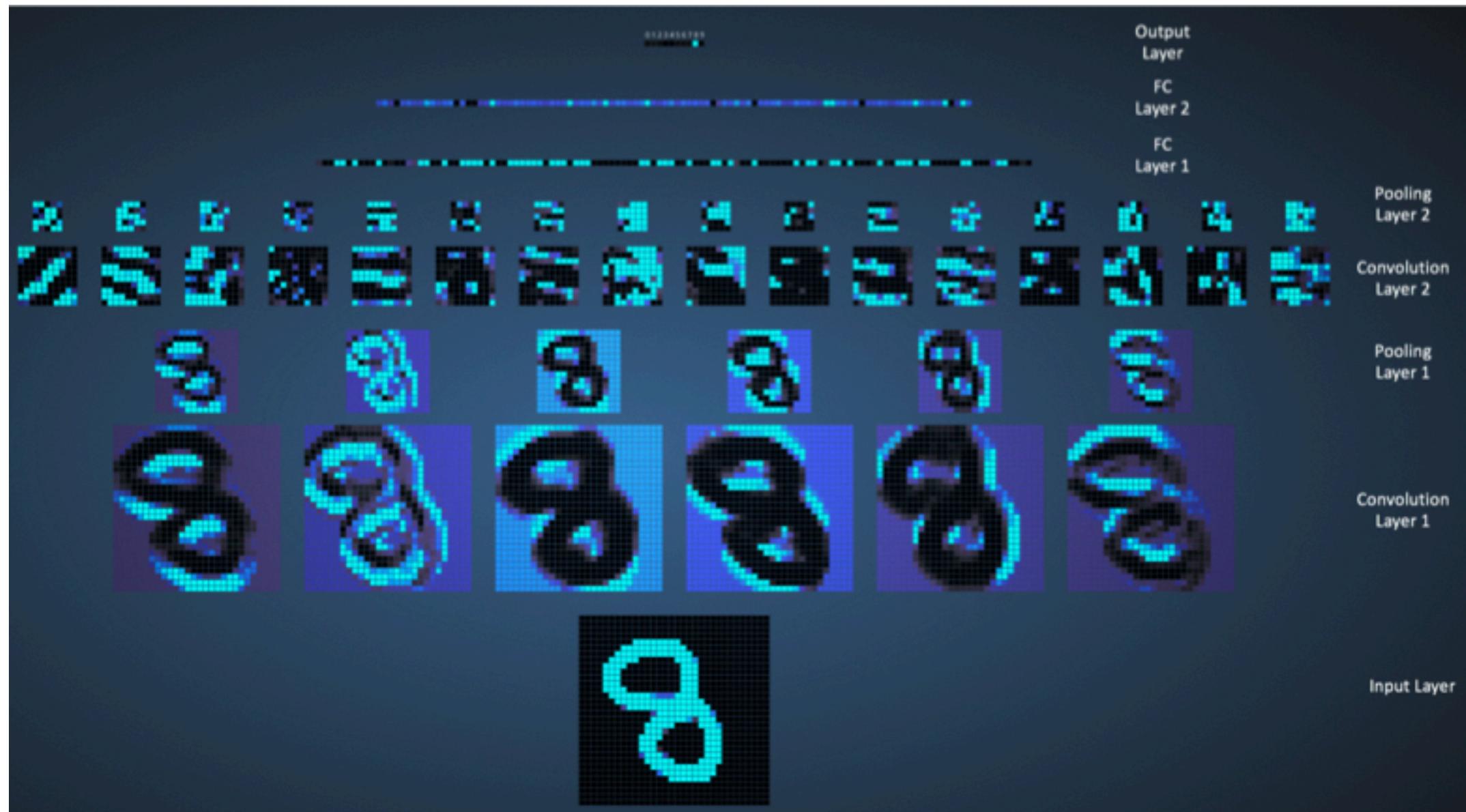
# Run Code

- **cnn.py**
  - What is happening here?
  - Why am I reshaping the input data?
  - How could I add more convolutional and pooling layers?

# LeNet Architecture 1994



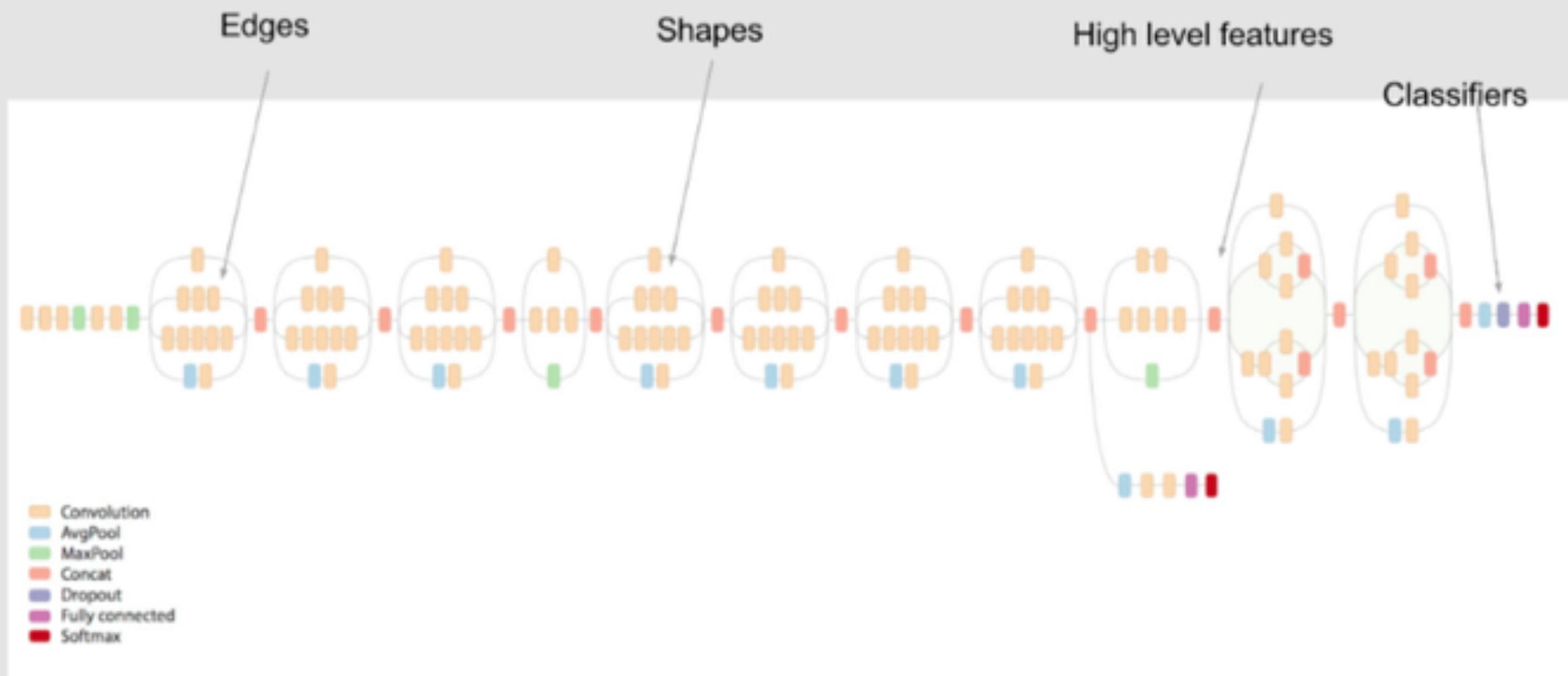
# CNN



# Run Code

- LeNet.py

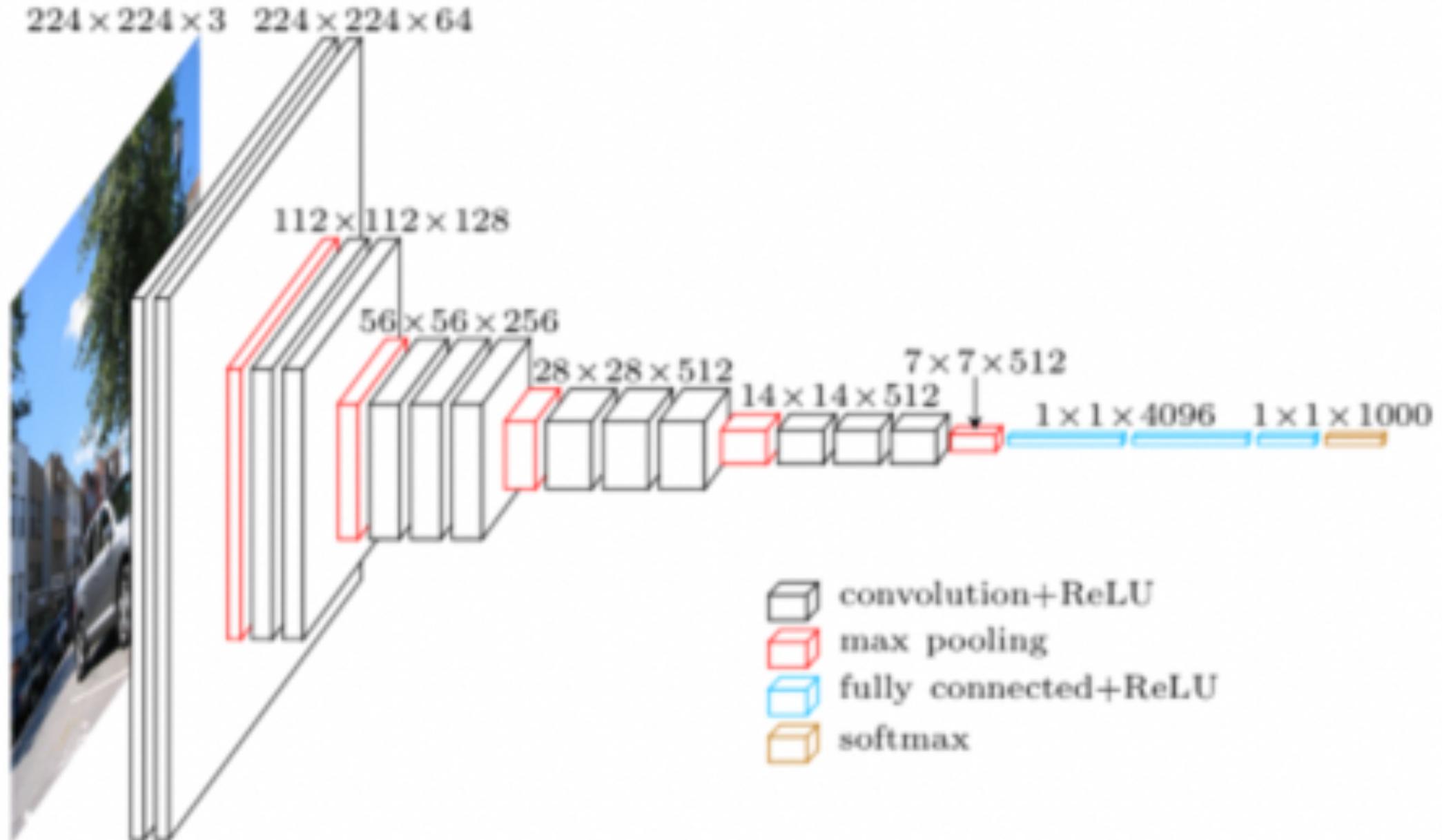
# What does the layers learn?



# Takeaways

- What is the difference between MLPs and CNNs?
- When might one or the other be more useful?

# VGG Architecture 2014



# Run Code

- **vgg\_inspect.py**
- **How many parameters does this model have?**
- **Why doesn't it overfit the training data?**

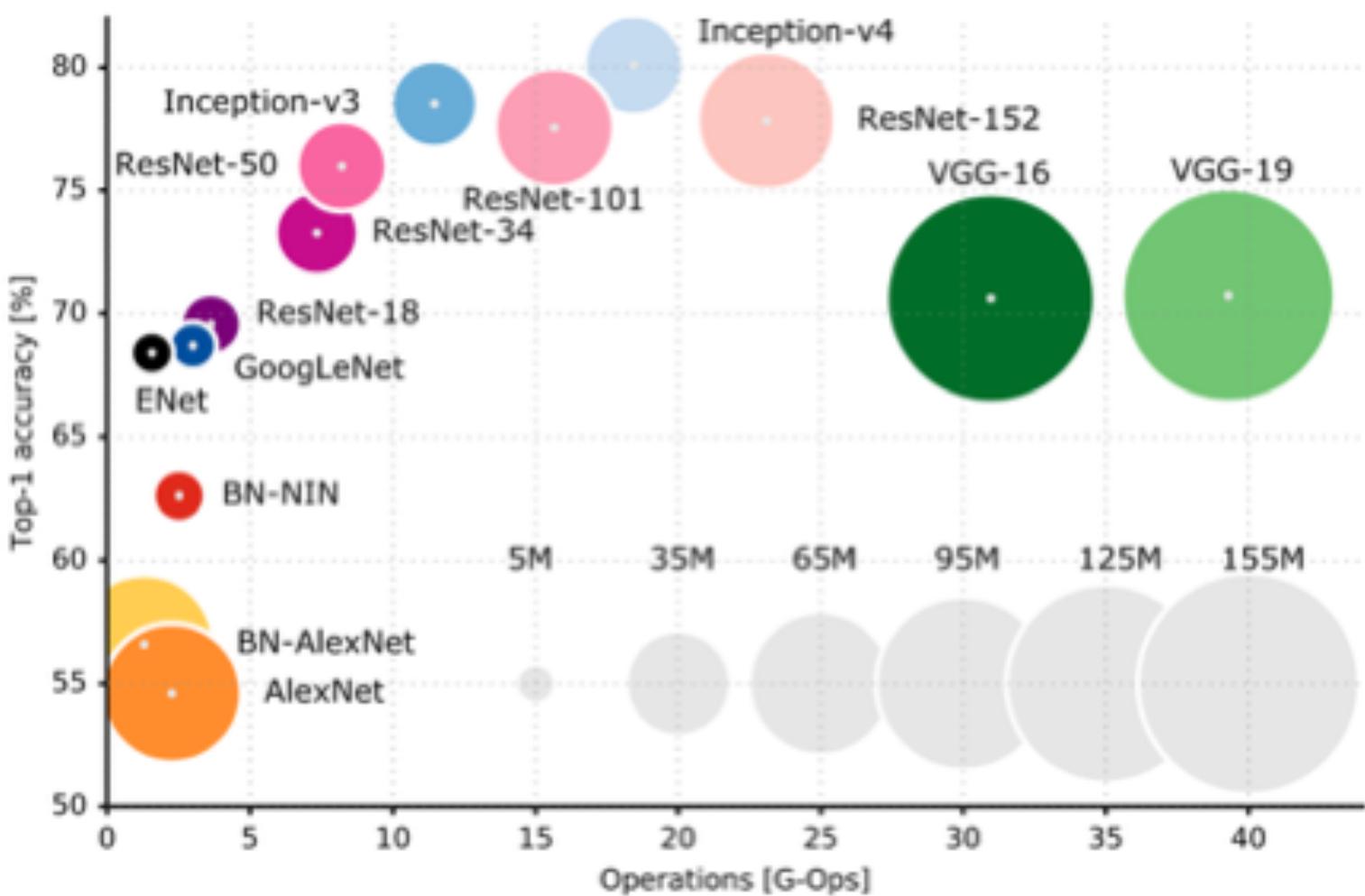
# Run Code

- `resnet50_inspect.py`

# Run Code

- `inception_inspect.py`

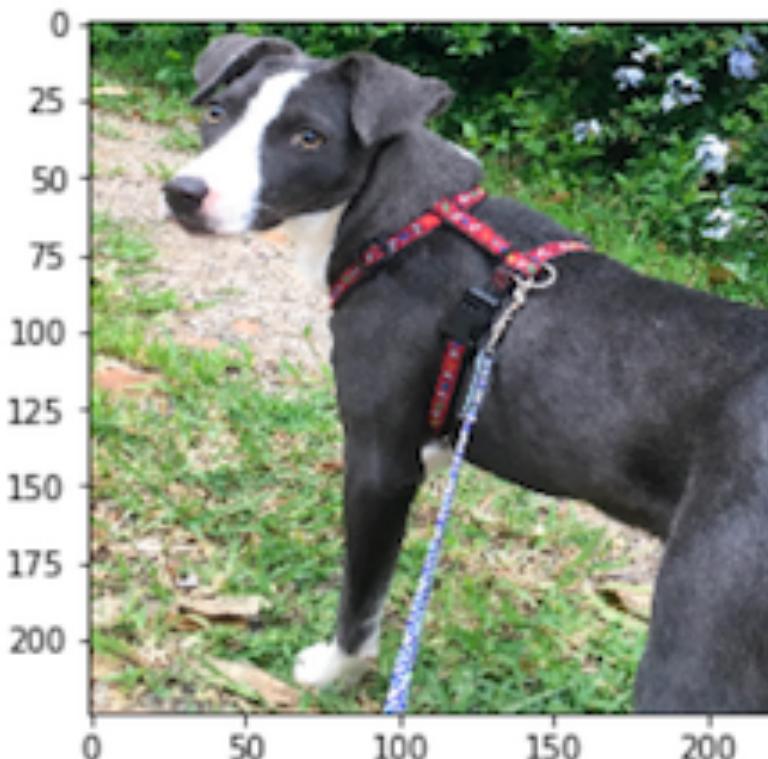
# Overview of Neural Networks



# Project- Dog Breed Classification

Given an image of a dog, the algorithm will identify an estimate of the canine's breed. If supplied an image of a human, the code will identify the resembling dog breed.

```
hello, dog!
your predicted breed is ...
American Staffordshire terrier
```



# Project- Dog Breed Classification

[Step 0](#): Import Datasets

[Step 1](#): Detect Humans

[Step 2](#): Detect Dogs

[Step 3](#): Create a CNN to Classify Dog Breeds (from Scratch)

[Step 4](#): Use a CNN to Classify Dog Breeds (using Transfer Learning)

[Step 5](#): Create a CNN to Classify Dog Breeds (using Transfer Learning)

[Step 6](#): Write your Algorithm

[Step 7](#): Test Your Algorithm

# Project- Dog Breed Classification

- Git clone [https://github.com/humayun/udacity\\_nanodegree/tree/master/p2\\_dog\\_project](https://github.com/humayun/udacity_nanodegree/tree/master/p2_dog_project)
- [https://github.com/humayun/udacity\\_nanodegree/blob/master/p2\\_dog\\_project/dog\\_app.ipynb](https://github.com/humayun/udacity_nanodegree/blob/master/p2_dog_project/dog_app.ipynb)