

Week 1:

# Introduction to Deep Learning

## PCLUB CS2020

Ankur Bhatia

Jan 25, 2020



# Course Prerequisites - Linear Algebra Calculus



# The Rise of Deep Learning

**'Deep Voice' Software Can Clone Anyone's Voice With Just 3.7 Seconds of Audio**

Using snippets of voices, Baidu's 'Deep Voice' can generate new speech, accents, and tones.



**'Creative' AlphaZero leads way for chess computers and, maybe, science**

Former chess world champion Garry Kasparov likes what he sees of computer that could be used to find cures for diseases



Complex of bacteria-infesting viral proteins modeled in CASP 13. The complex contains that were modeled individually. MIT CSAIL/CSAIL

Google's DeepMind aces protein folding

By Robert F. Service | Dec. 6, 2018, 12:35 PM

DEEPMIND IS STARSCRAFT II'S TRIUMPH FOR



How an A.I. 'Cat-and-Mouse Game' Generates Believable Fake Photos

By CHIC KETTLE AND KEITH COLLINS | JAN. 3, 2019



To create the first image in this set, the system generated 30 million iterations over 18 days.

Let There Be Sight: How Deep Learning Is Helping the Blind 'See'



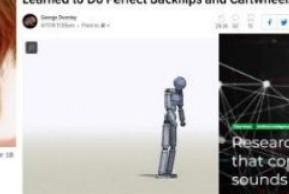
Technology outpacing security measures



Neural networks everywhere



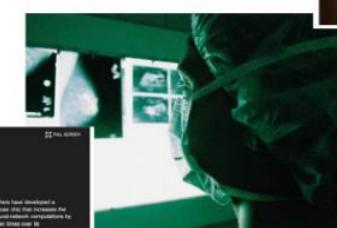
After Millions of Trials, These Simulated Humans Learned to Do Perfect Backflips and Cartwheels



6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)

AI beats docs in cancer spotting

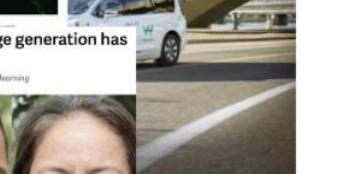
A new study provides a fresh example of machine learning as an important diagnostic tool. Paul Biagler reports.



These faces show how far AI image generation has come in just four years



AI Can Help In Predicting Cryptocurrency Value



Automation And Algorithms: De-Risking Manufacturing With Artificial Intelligence

Sarah Goshko Contributing Writer  
Manufacturing

Issue on the industrialization of additive manufacturing

TWEET THIS

The two key applications of AI in manufacturing



# What is Deep Learning?

## ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



## MACHINE LEARNING

Ability to learn without explicitly being programmed



## DEEP LEARNING

Extract patterns from data using neural networks



# Course Outcomes:

1. FINAL CLASS PROJECT
2. PAPER/ IDEA PRESENTATION
3. PROJECTS SHOWCASE ON PCLUB WEBSITE.



# Class support

DL- Wiki (Telegram Group)  
[bit.ly/dlwikiclub](https://bit.ly/dlwikiclub)



# Why Deep Learning and Why Now?

## Things we want to do with data

Images



→ Label image

Audio



→ Speech recognition

Text



→ Web search

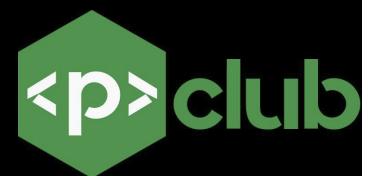
Andrew Ng



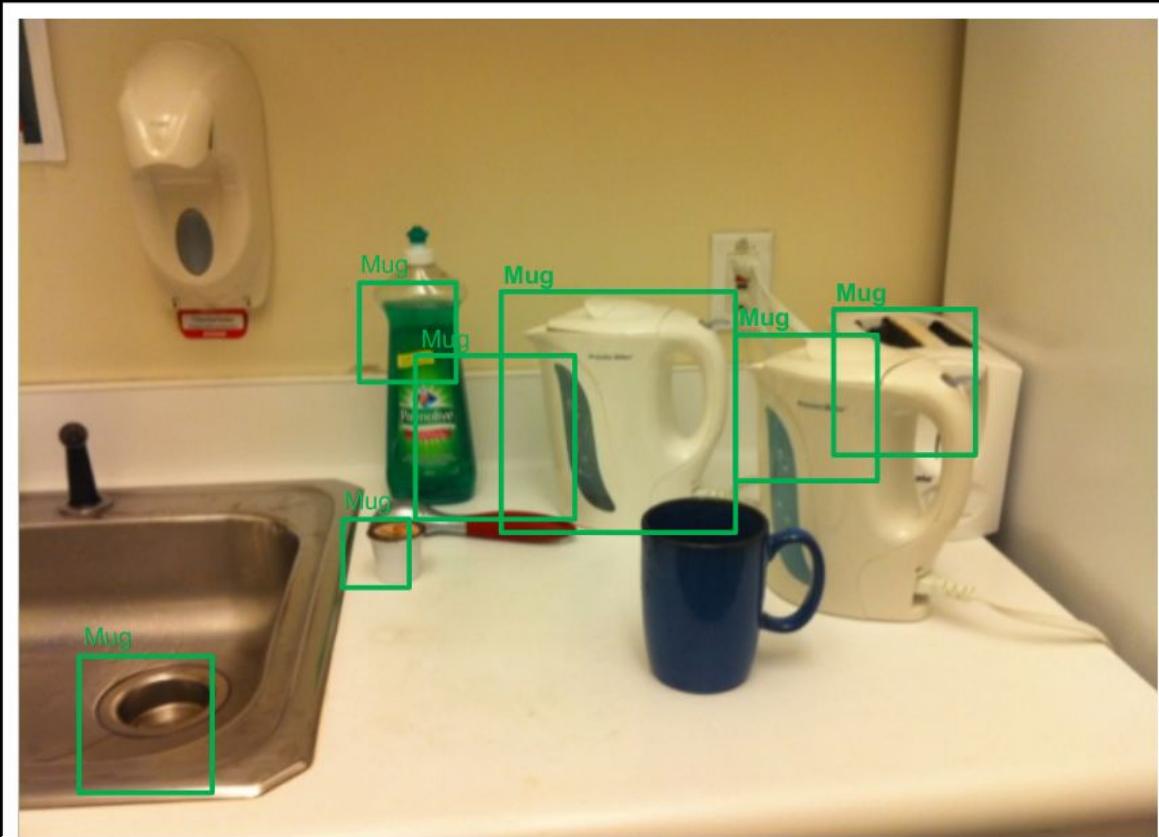
## Computer vision: Identify coffee mug



Andrew Ng



## Computer vision: Identify coffee mug



Andrew Ng



# Why is computer vision hard?



The camera sees :

194	210	201	212	199	213	215	195	178	158	182	209
180	189	190	221	209	205	191	167	147	115	129	163
114	126	140	188	176	165	152	140	170	106	78	88
87	103	115	154	143	142	149	153	173	101	57	57
102	112	106	131	122	138	152	147	128	84	58	66
94	95	79	104	105	124	129	113	107	87	69	67
68	71	69	98	89	92	98	95	89	88	76	67
41	56	68	99	63	45	60	82	58	76	75	65
20	43	69	75	56	41	51	73	55	70	63	44
50	50	57	69	75	75	73	74	53	68	59	37
72	59	53	66	84	92	84	74	57	72	63	42
67	61	58	65	75	78	76	73	59	75	69	50

Andrew Ng



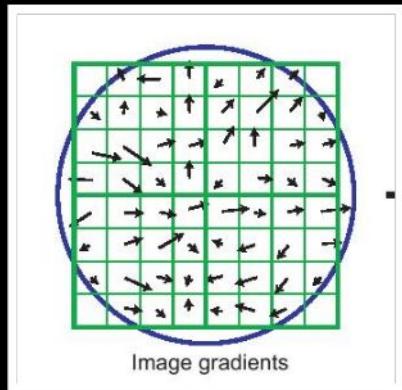
# Computer vision



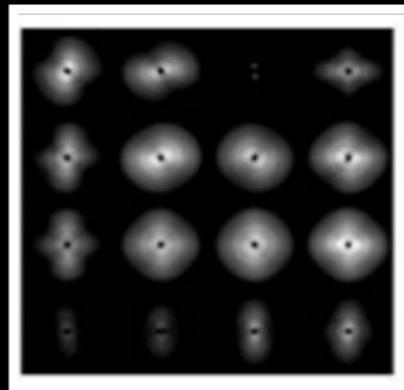
Andrew Ng



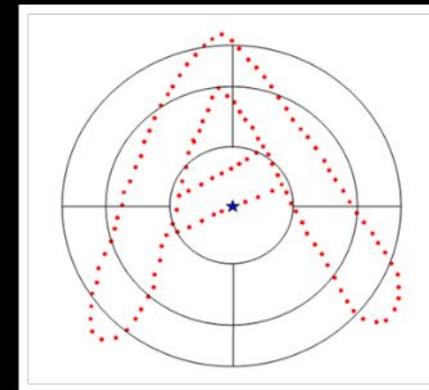
## Features for vision



SIFT

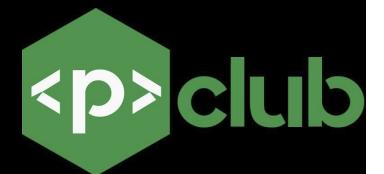


GIST



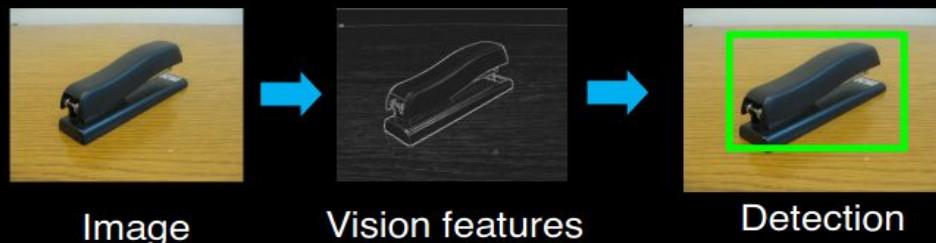
Shape context

Andrew Ng

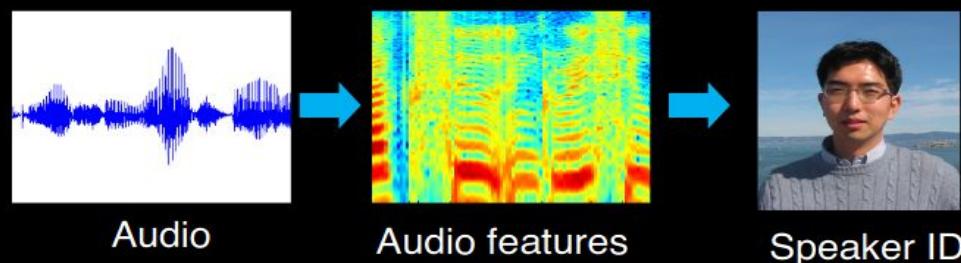


# Features for machine learning

Images

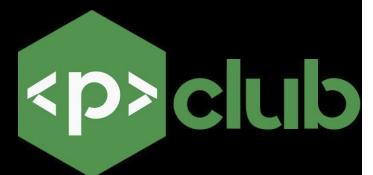
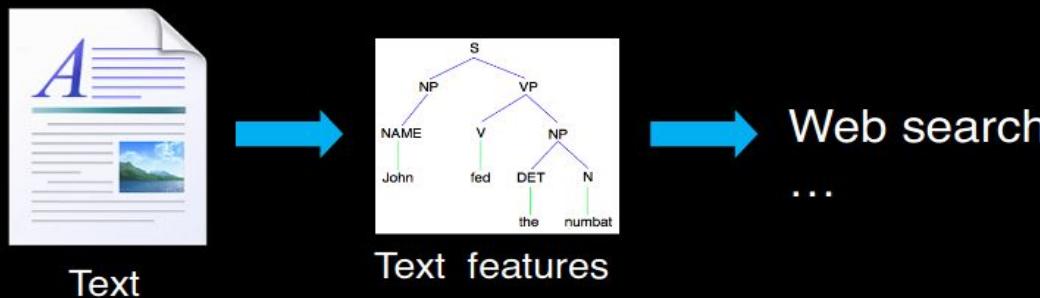


Audio



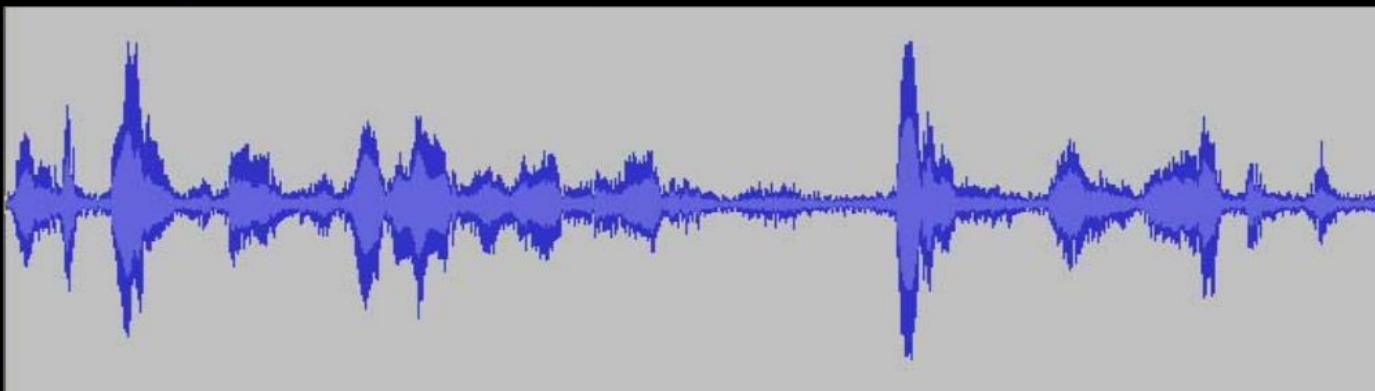
Text

Andrew Ng



# Why is speech recognition hard?

Microphone recording:



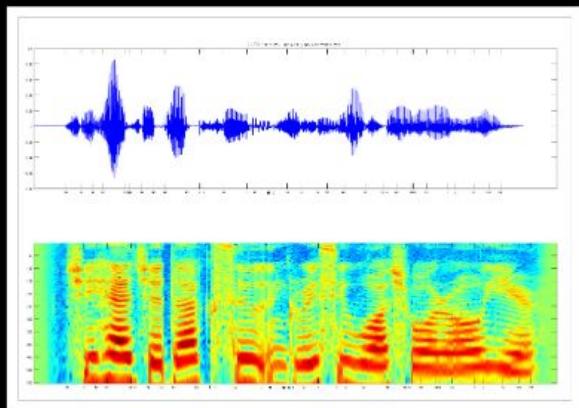
Please find the coffee mug



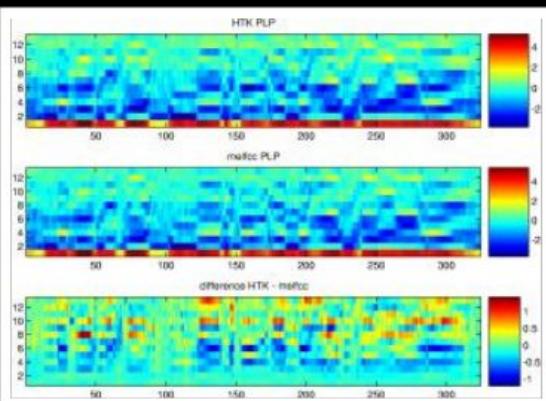
Andrew Ng



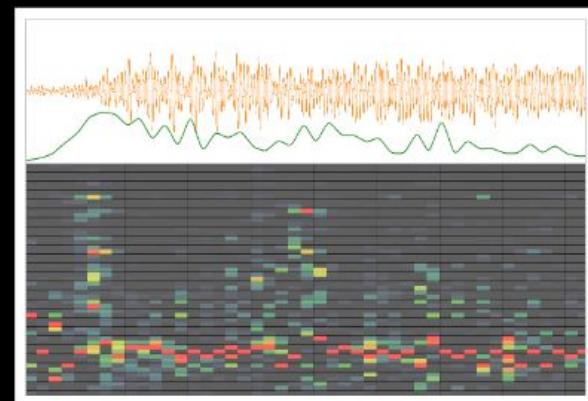
# Features for audio



Spectrogram



MFCC

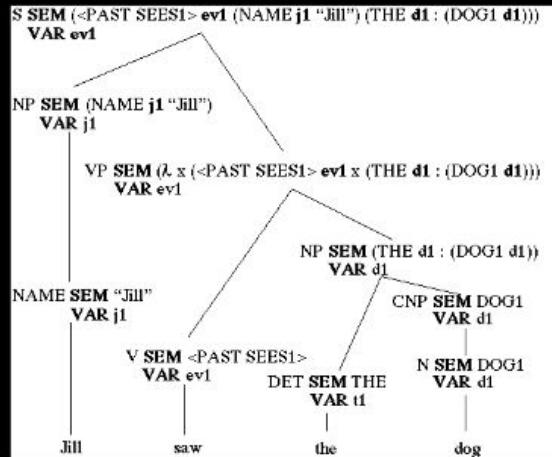


Flux

Andrew Ng



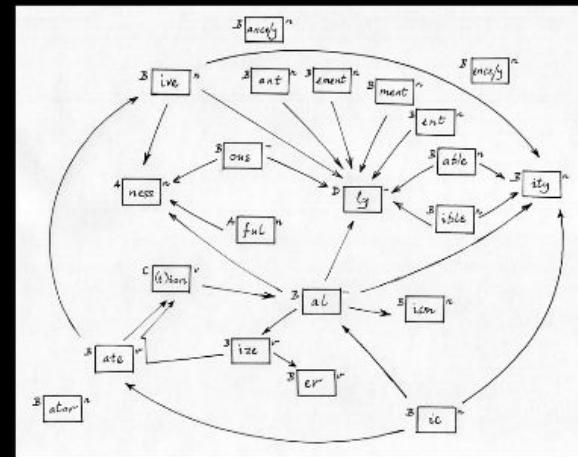
## Features for text



# Parser

<DOC>  
<DOCID> wsj94\_008.0212 </DOCID>  
<DOCNO> 940413-0062 </DOCNO>  
<HL> Who's News:  
@ Burns Fry Ltd. </HL>  
<DD> 04/13/94 </DD>  
<SO> WALL STREET JOURNAL (J), PAGE B10 </SO>  
<CO> MER </CO>  
<IN> SECURITIES (SCR) </IN>  
<TXT>  
<p> BURNS FRY Ltd. (Toronto) -- Donald Wright, 4  
named executive vice president and director of  
brokerage firm. Mr. Wright resigned as president  
Canada Inc., a unit of Merrill Lynch & Co., to  
Kassirer, 48, who left Burns Fry last month. A  
spokeswoman said it hasn't named a successor to  
expected to begin his new position by the end of  
</p>  
</TXT>  
</DOC>

## Named entity



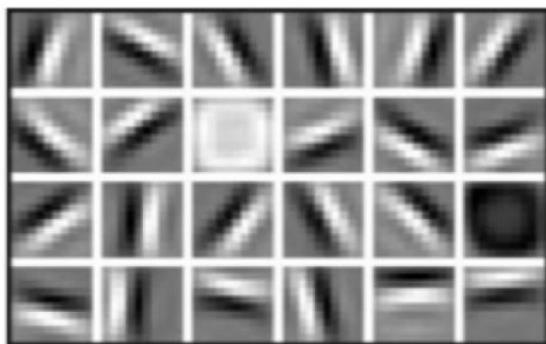
## Stemming

# Why Deep Learning?

Hand engineered features are time consuming, brittle and not scalable in practice

Can we learn the **underlying features** directly from data?

Low Level Features



Lines & Edges

Mid Level Features



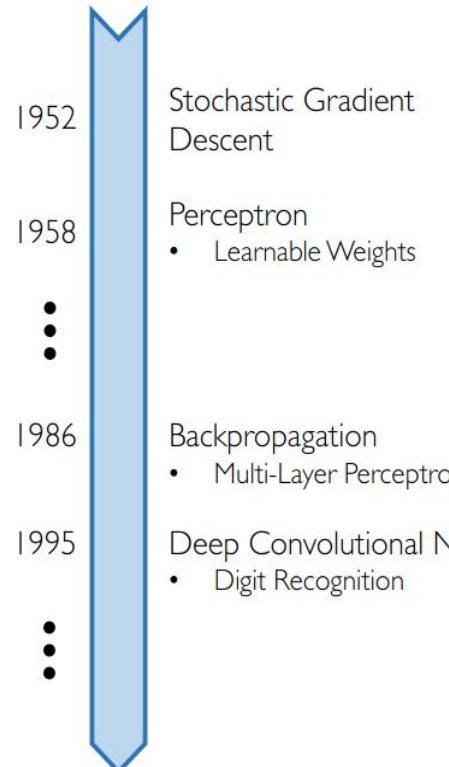
Eyes & Nose & Ears

High Level Features



Facial Structure

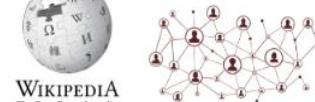
# Why Now?



Neural Networks date back decades, so why the resurgence?

## I. Big Data

- Larger Datasets
- Easier Collection & Storage



## 2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



## 3. Software

- Improved Techniques
- New Models
- Toolboxes



6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)

# Deep Learning Applications?

**Everyone knows it.**

Let's see one state of the art one:

<https://ai.googleblog.com/2018/05/duplex-ai-system-for-natural-conversation.html>



# Deep Learning in One Slide

- **What is it:**  
Extract useful patterns from data.
- **How:**  
Neural network + optimization
- **How (Practical):**  
Python + TensorFlow & friends
- **Hard Part:**  
Good Questions + Good Data
- **Why now:**  
Data, hardware, community, tools, investment
- **Where do we stand?**  
Most big questions of intelligence have not been answered nor properly formulated

## Exciting progress:

- Face recognition
- Image classification
- Speech recognition
- Text-to-speech generation
- Handwriting transcription
- Machine translation
- Medical diagnosis
- Cars: drivable area, lane keeping
- Digital assistants
- Ads, search, social recommendations
- Game playing with deep RL

# Deep Learning

## What now in 2020?



# Deep Learning & AI in Context of Human History



Perspective:

- Universe created  
13.8 billion years ago
- Earth created  
4.54 billion years ago
- Modern humans  
300,000 years ago
- Civilization  
12,000 years ago
- Written record  
5,000 years ago



**1700s and beyond:** Industrial revolution, steam engine, mechanized factory systems, machine tools

6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)



# Artificial Intelligence in Context of Human History



## Perspective:

- Universe created  
13.8 billion years ago
- Earth created  
4.54 billion years ago
- Modern humans  
300,000 years ago
- Civilization  
12,000 years ago
- Written record  
5,000 years ago



## Dreams, mathematical foundations, and engineering in reality.

**Alan Turing, 1951:** "It seems probable that once the machine thinking method had started, it would not take long to outstrip our feeble powers. They would be able to converse with each other to sharpen their wits. At some stage therefore, we should have to expect the machines to take control."

- 1943: Neural networks
- 1957-62: Perceptron
- 1970-86: Backpropagation, RBM, RNN
- 1979-98: CNN, MNIST, LSTM, Bidirectional RNN
- 2006: “Deep Learning”, DBN
- 2009: ImageNet + AlexNet
- 2014: GANs
- 2016-17: AlphaGo, AlphaZero
- 2017: 2017-19: Transformers

\* Dates are for perspective and not as definitive historical record of invention or credit



• Yann LeCun  
 • Geoffrey Hinton  
 • Yoshua Bengio

Turing Award given for:

- “The conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing.”

## Deep Learning Growth, Celebrations, and Limitations Hopes for 2020

- **Less Hype & Less Anti-Hype:** Less tweets on how there is too much hype in AI and more solid research in AI.
- **Hybrid Research:** Less contentious, counter-productive debates, more open-minded interdisciplinary collaboration.
- **Research topics:**
  - Reasoning
  - Active learning and life-long learning
  - Multi-modal and multi-task learning
  - Open-domain conversation
  - Applications: medical, autonomous vehicles
  - Algorithmic ethics
  - Robotics

# Now Let's Dive in?



## Tagged vs. untagged data



Coffee mug



Coffee mug



Coffee mug



Coffee mug



Coffee mug



Coffee mug

Andrew Ng



## Untagged data (unsupervised learning)



Unknown



Unknown



Unknown



Unknown



Unknown



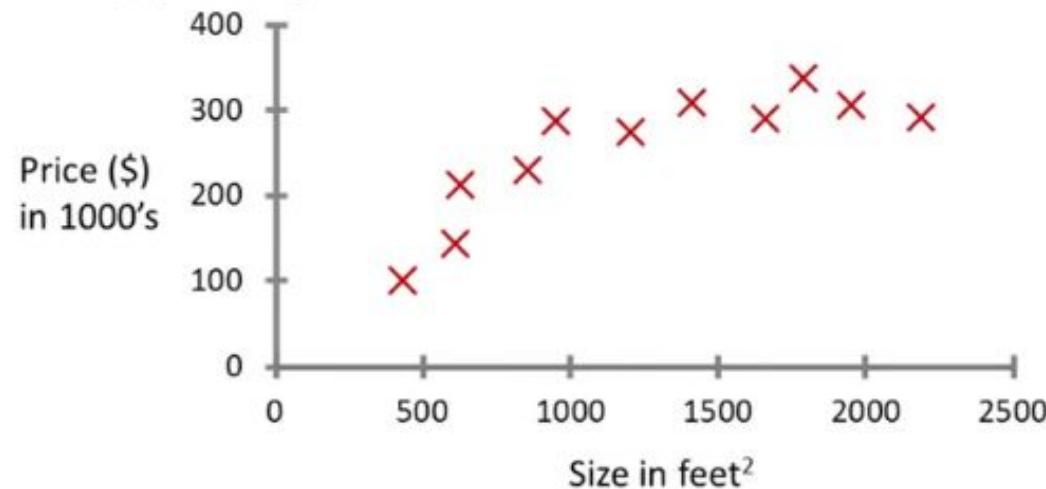
Unknown

Andrew Ng

# Supervised Learning

## EXAMPLE 1

Housing price prediction.



# Supervised Learning

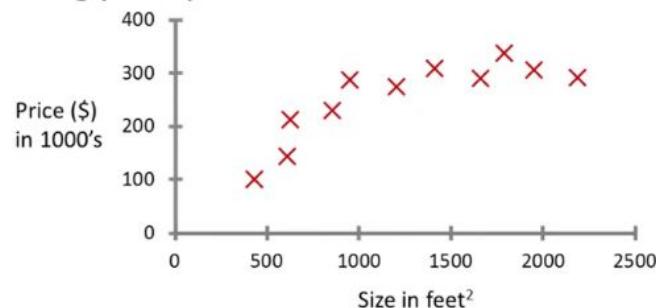
## EXAMPLE 2

Breast cancer (malignant, benign)



# Linear Regression:

Housing price prediction.



Hypothesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x$

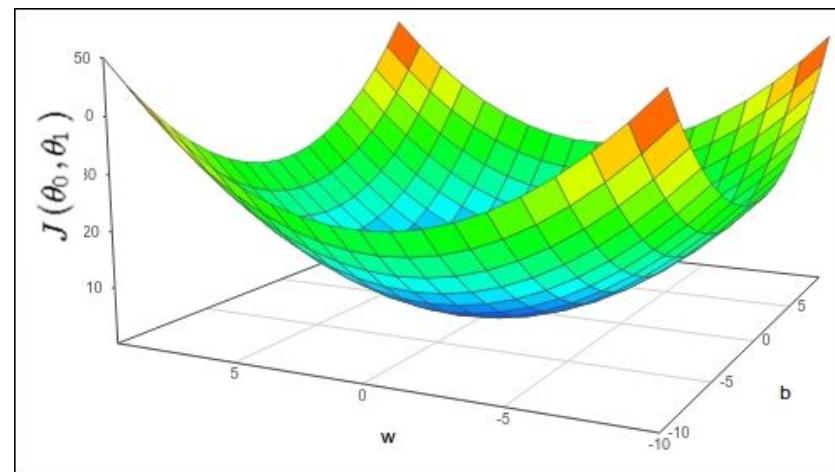
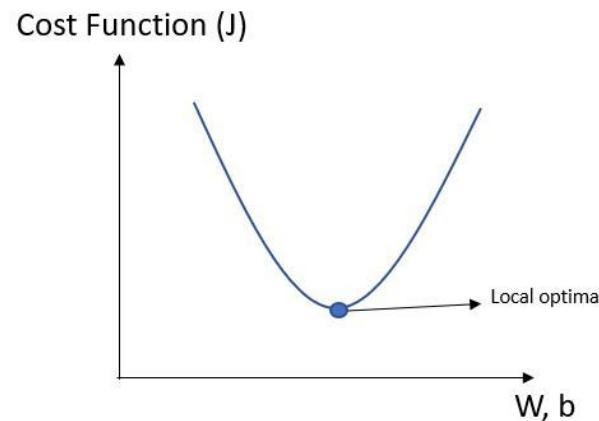
Parameters:  $\theta_0, \theta_1$

Cost Function:  $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal:  $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

1. Notations
2. Hypothesis function and its parameters
3. Cost Function
  - a. Single parameter
  - b. Multiple parameters
4. Contour plots points of various  $h(\theta_0, \theta_1)$
5. Minimize cost function

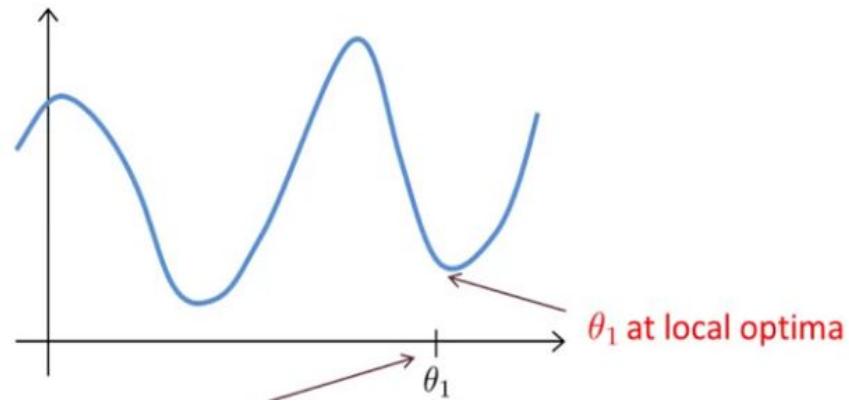
# Cost Function:



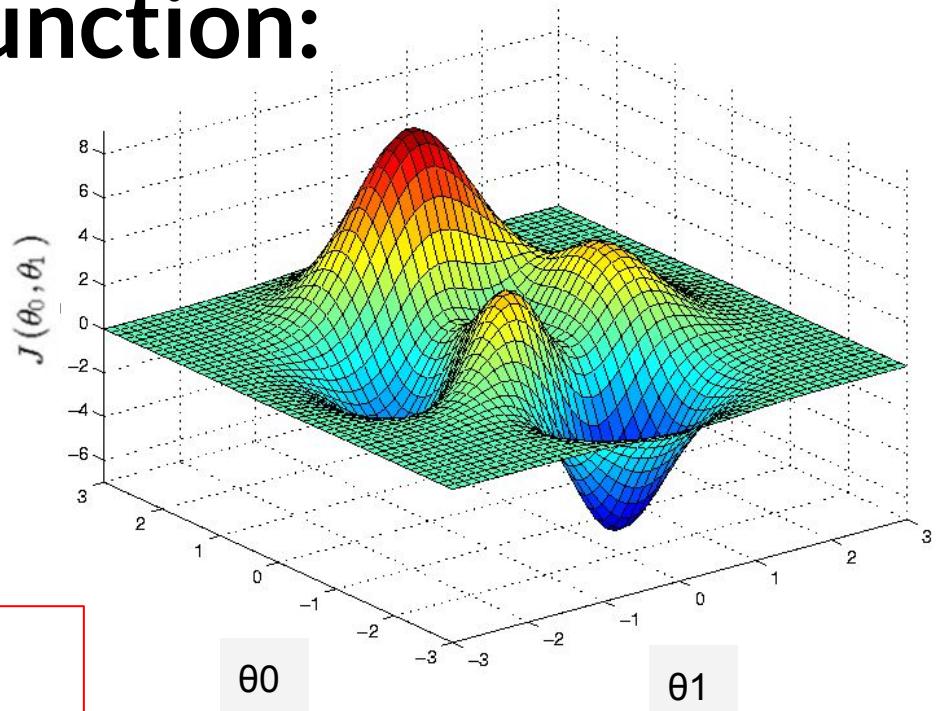
1. Notations
2. Hypothesis function and its parameters
3. Cost Function
  - a. Single parameter
  - b. Multiple parameters
4. Contour plots points of various  $h(\theta_0, \theta_1)$
5. Minimize cost function

Convex Function

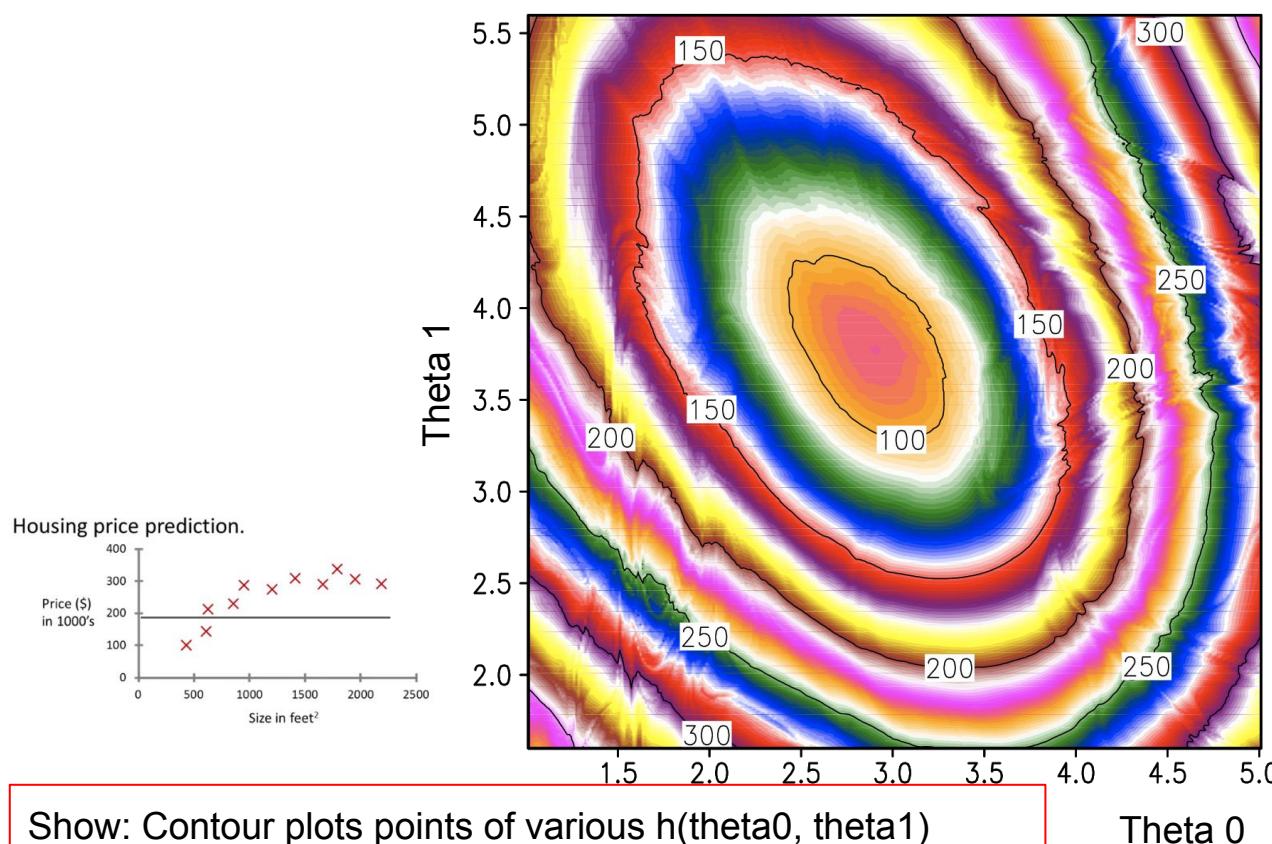
# Cost Function:



- 1. Notations
- 2. Hypothesis function and its parameters
- 3. Cost Function
  - a. Single parameter
  - b. Multiple parameters
- 4. Contour plots points of various  $h(\theta_0, \theta_1)$
- 5. Minimize cost function



# Cost Function contour:



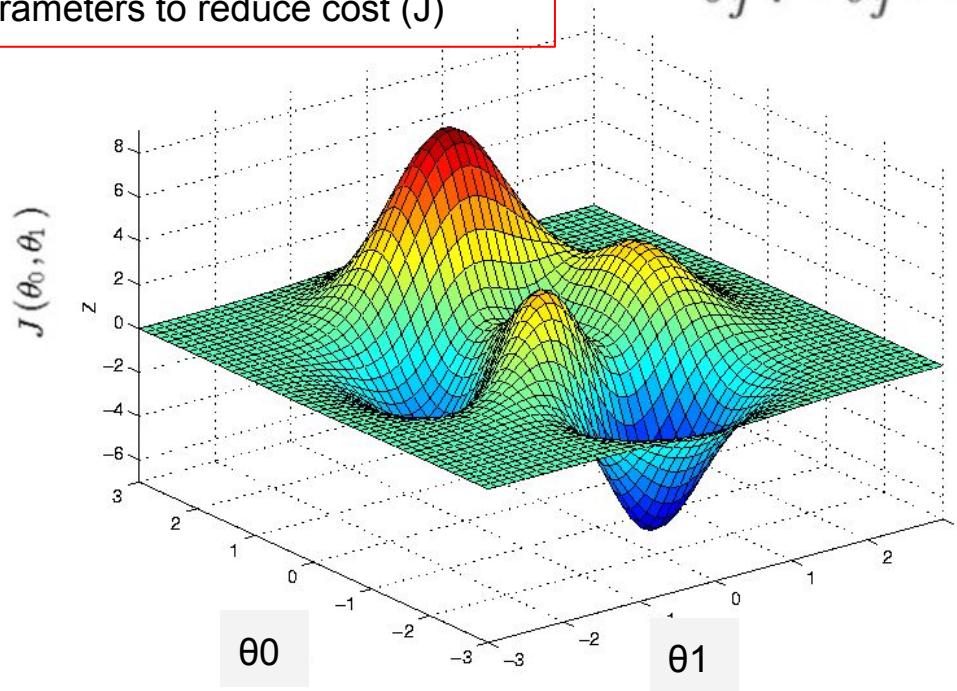
What does a point denote in this?

# Gradient Descent “Batch”:

Task: Minimize  $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

1. Start with some initial parameter value  $\theta_1, \theta_2$
2. Change the parameters to reduce cost ( $J$ )

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$



# Gradient Descent for Linear Reg:

Task: Minimize  $J(\theta_0, \theta_1) = \frac{1}{2m} (\sum(h(x) - y)^2)$

1. Compute the derivative of cost function wrt. each of the parameters.
2. Put into the gradient descent algorithm.
3. Stop when get reached the convergence.

$$j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

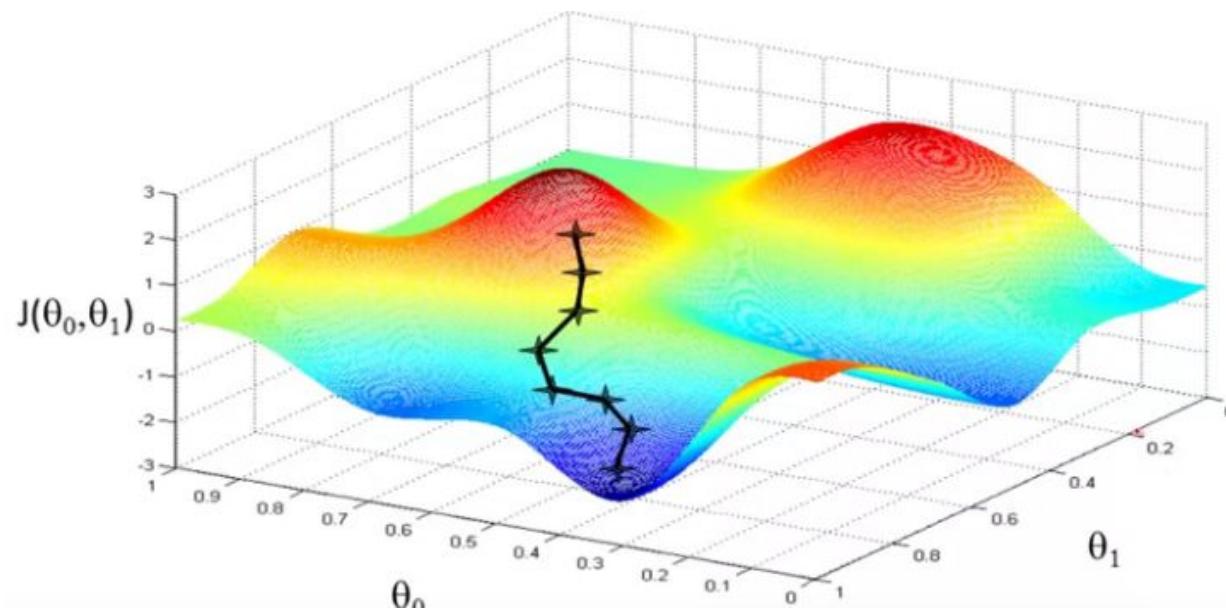
```
repeat until convergence {  
     $\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$   
     $\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$   
}
```



# Gradient Descent for Linear Reg:

Task: Minimize  $J(\theta_0, \theta_1) = 1/2m (\sum(h(x) - y)^2)$

Q. What happens if algo gets stuck in a local minima?

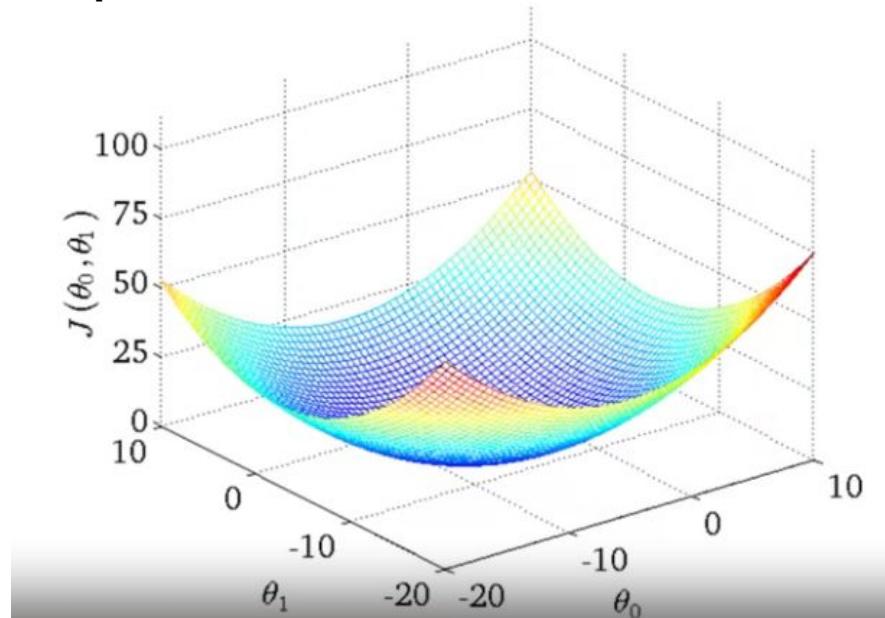


# Gradient Descent for Linear Reg:

Task: Minimize  $J(\theta_0, \theta_1) = 1/2m (\sum(h(x) - y)^2)$

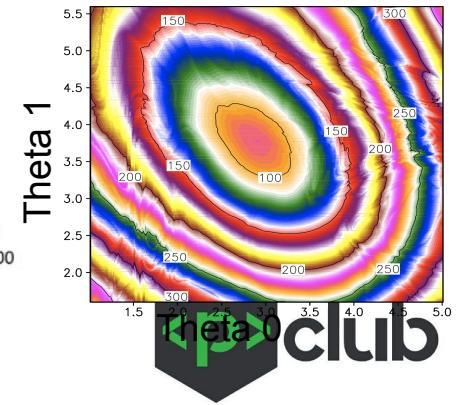
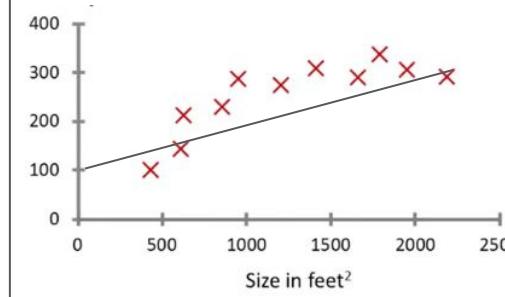
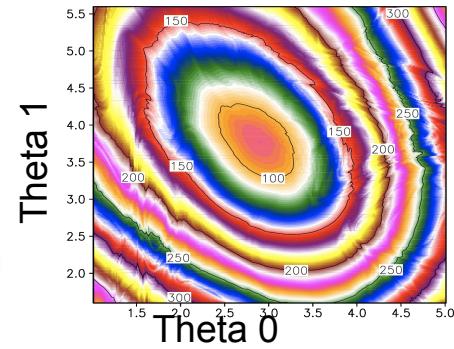
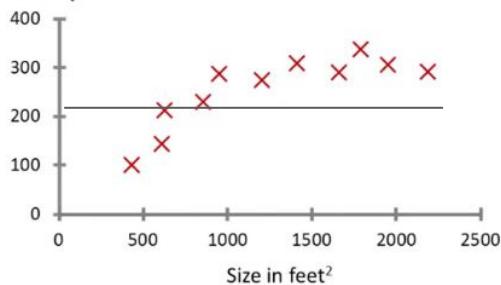
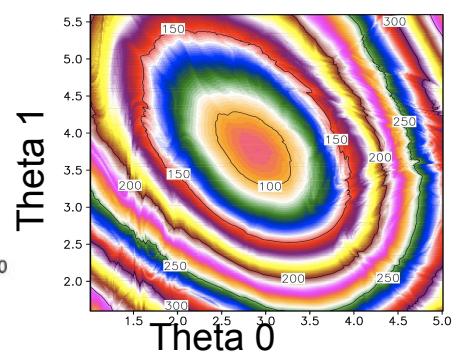
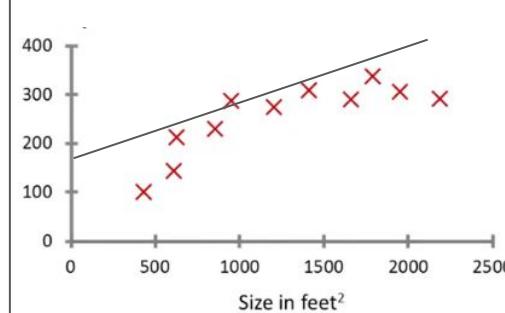
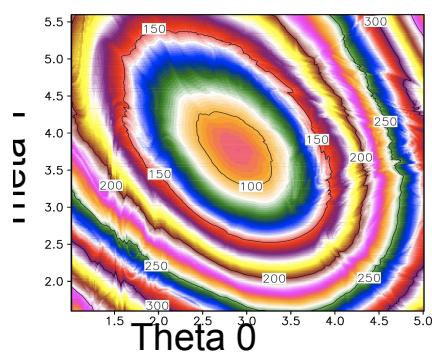
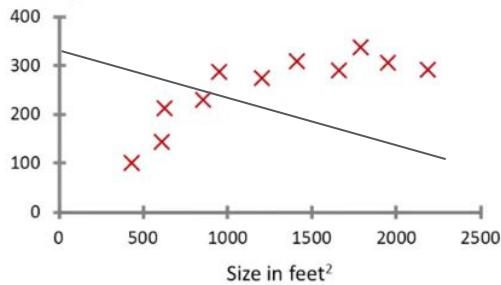
Q. What happens if algo gets stuck in a local minima?

Ans. For Linear Reg. , we have a bowl type cost function (Convex function).  
Optimization problem is convex here.

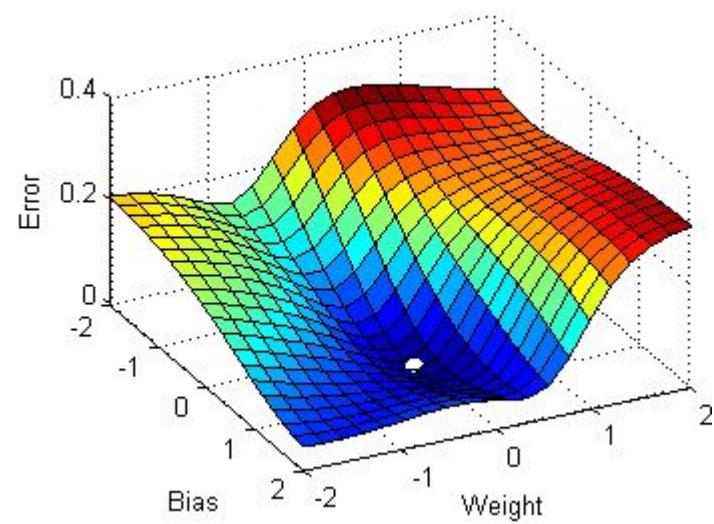
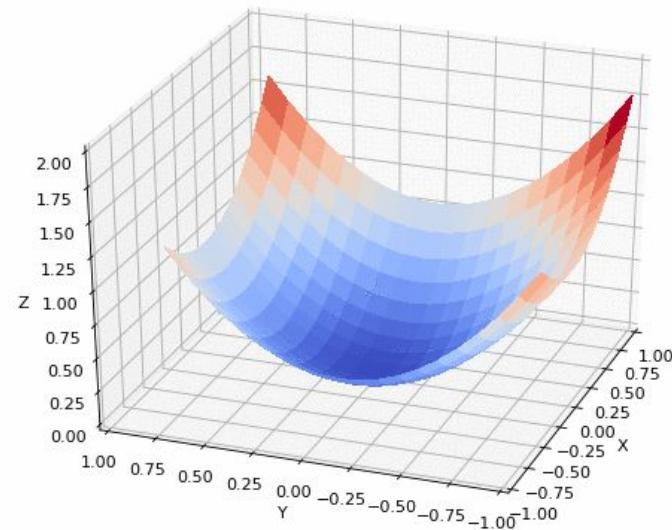


# Gradient Descent for Linear Reg:

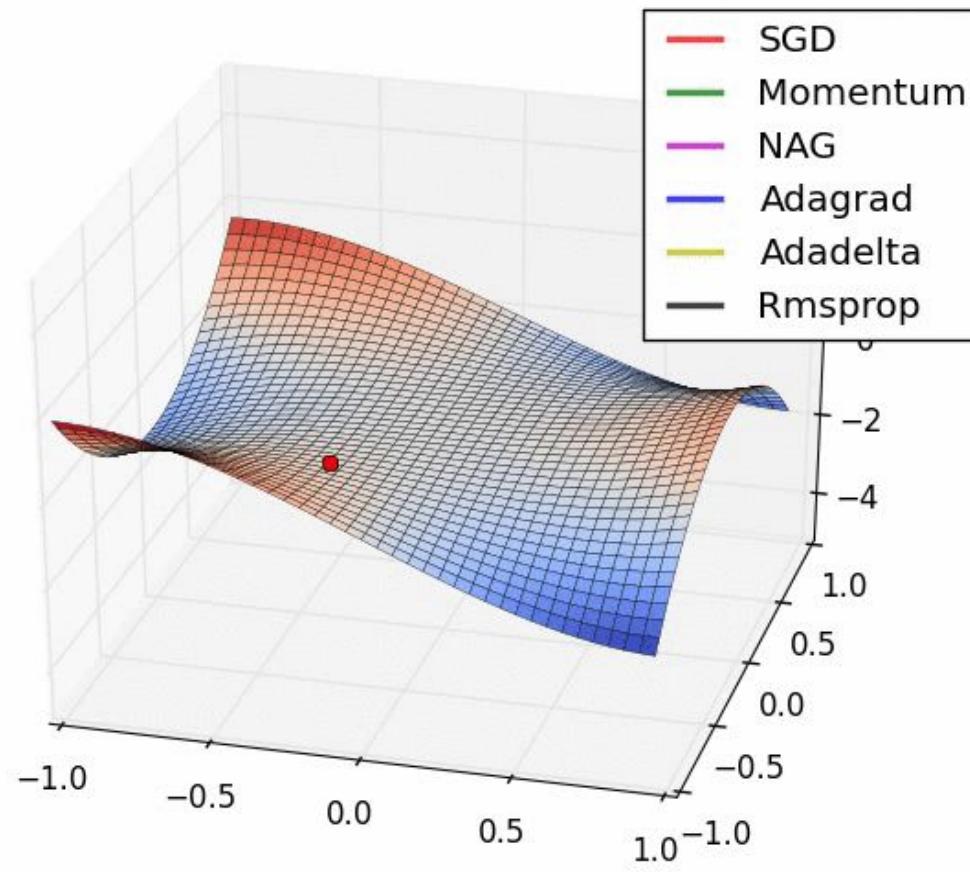
Task: Minimize  $J(\theta_0, \theta_1) = \frac{1}{2m} (\sum(h(x) - y)^2)$



# Gradient Descent for Linear Reg:



# Various Descent's:

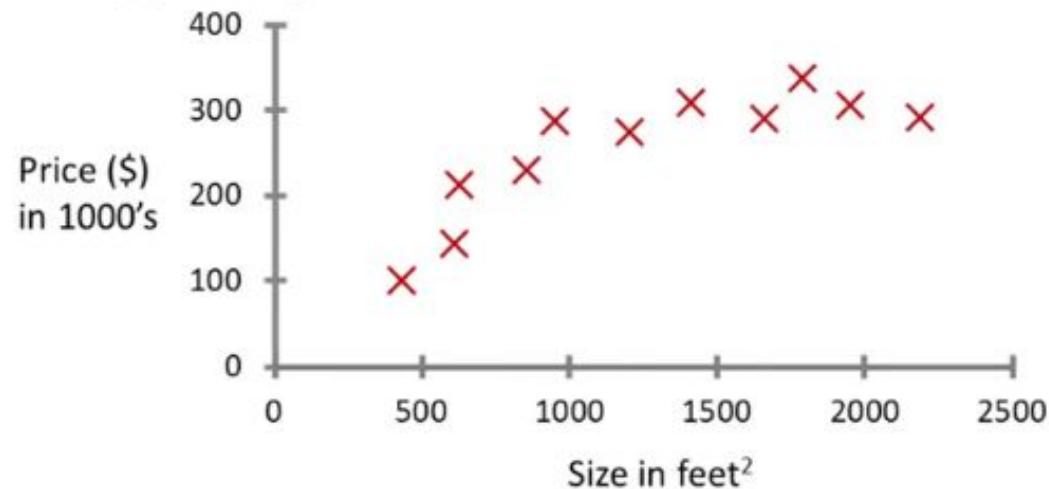


**Now Let's Dive a bit more Deep.**

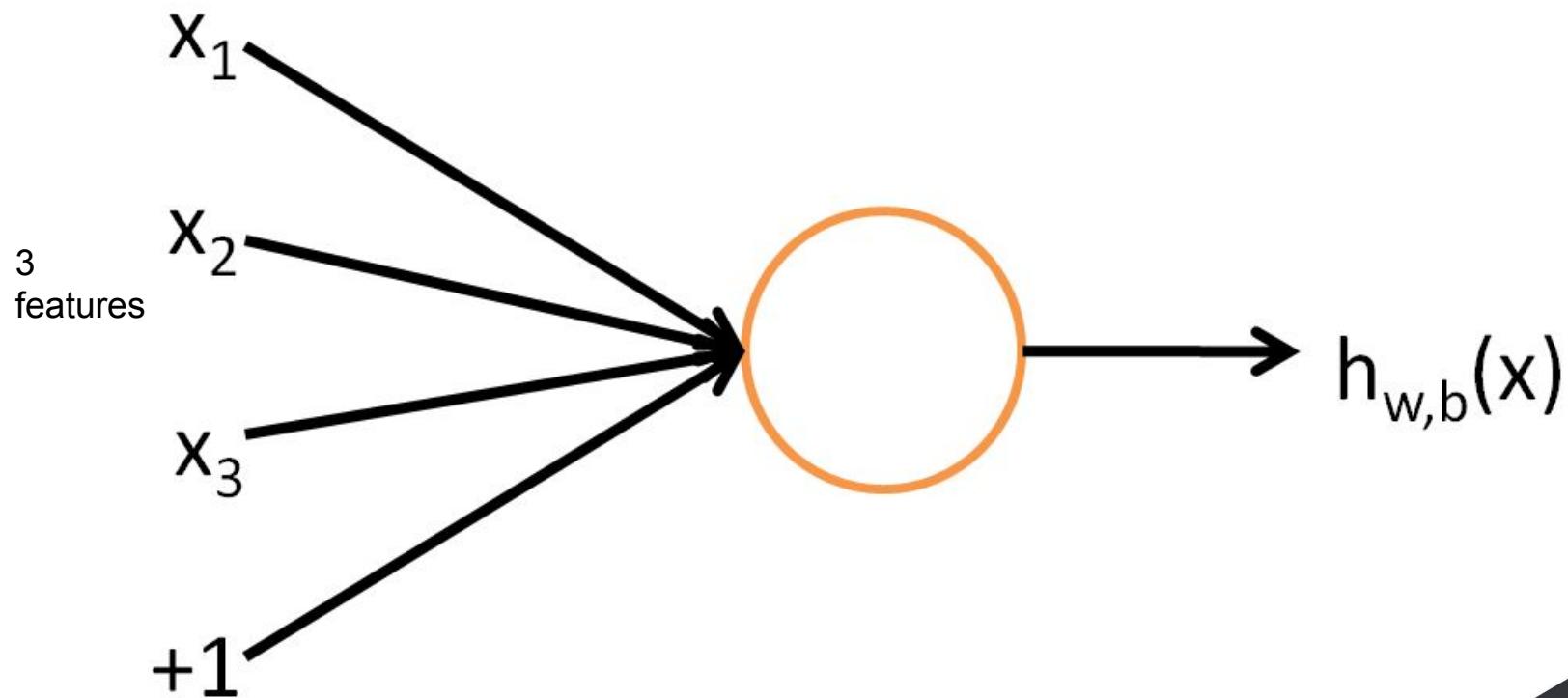


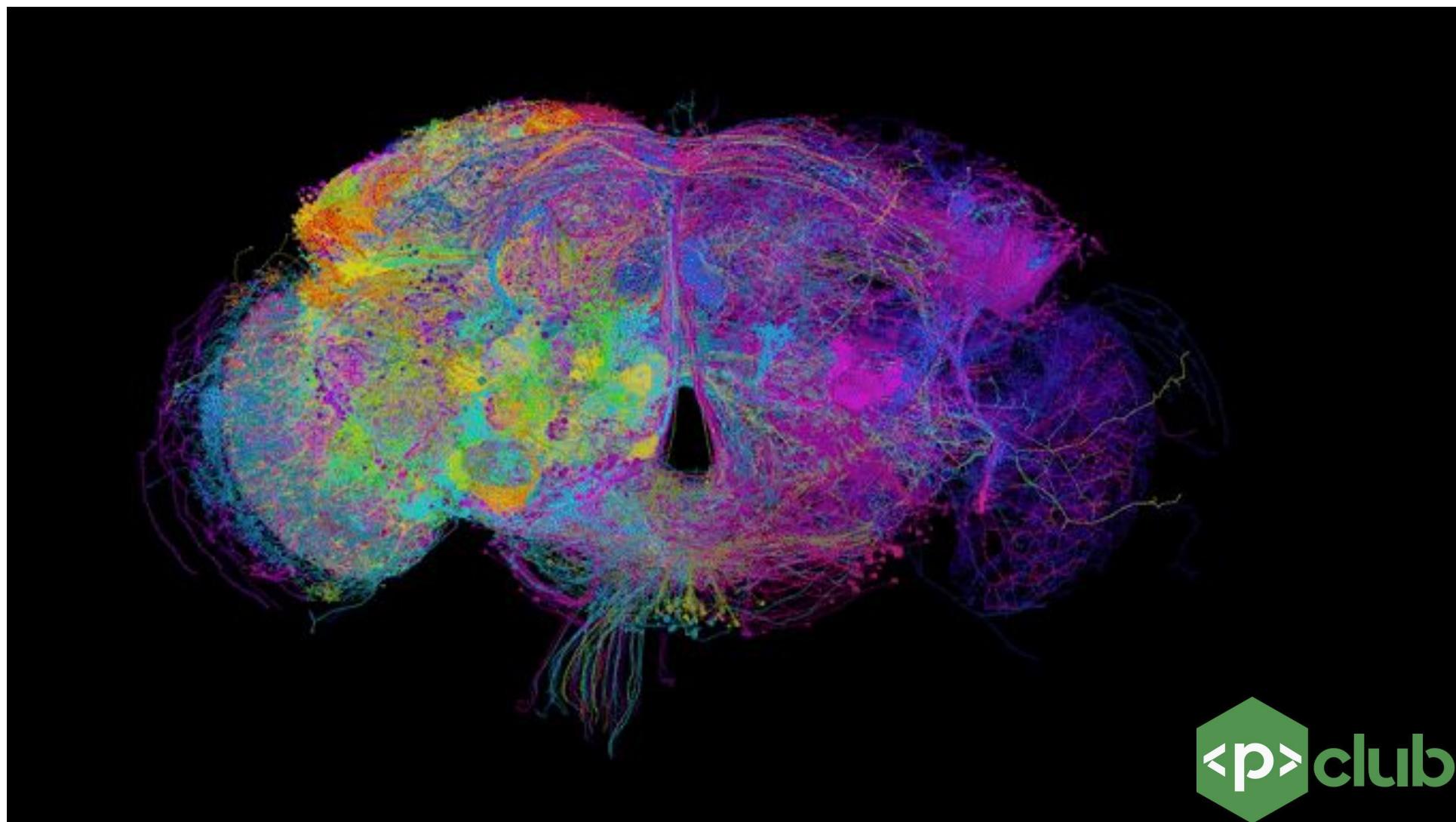
# What is a Neural Network?

Housing price prediction.



# What is a Neural Network?





<P>club

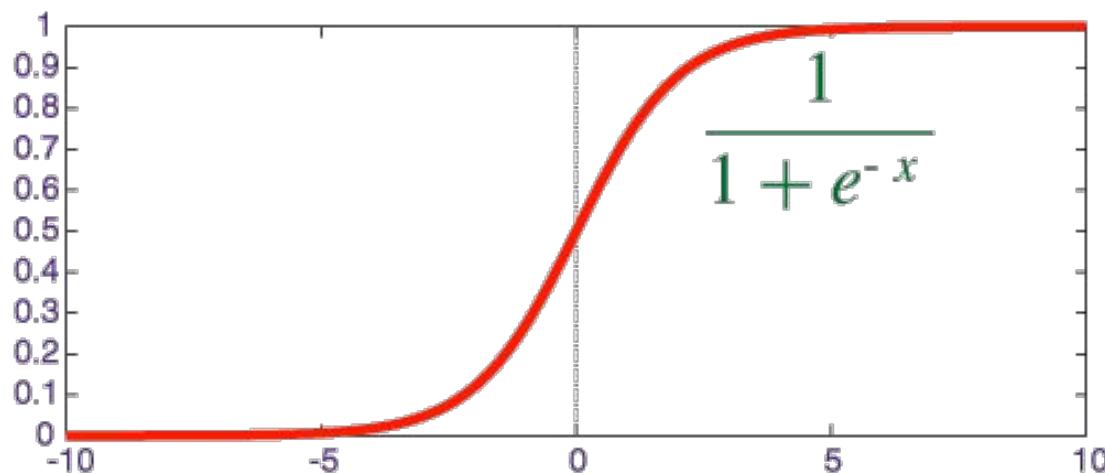
# Binary Classification problem:



Email - Spam/Not spam  
Ad - Clicked/Not clicked



# Understanding Logistic Regression problem:

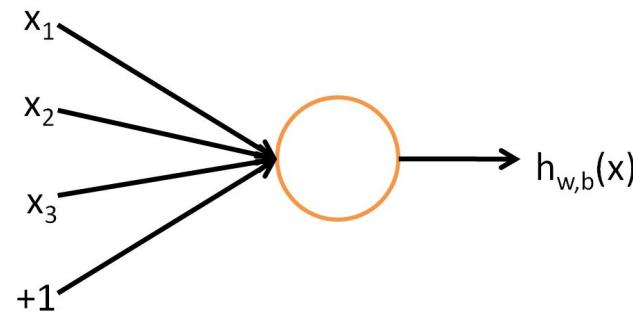


Used in Supervised Learning  
When outputs are 0/1 (BC)



# Logistic Regression loss function:

Log Loss / Binary Cross Entropy Loss



$$\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Gradient Descent may not find a global optimum -  $(\hat{y} - y)^2$



# Logistic Regression Cost function:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

This cost function allows a convex optimization.



# Logistic Regression Gradient Descent:

Backprop:

On 1 Training example



# Logistic Regression Gradient Descent:

Backprop:

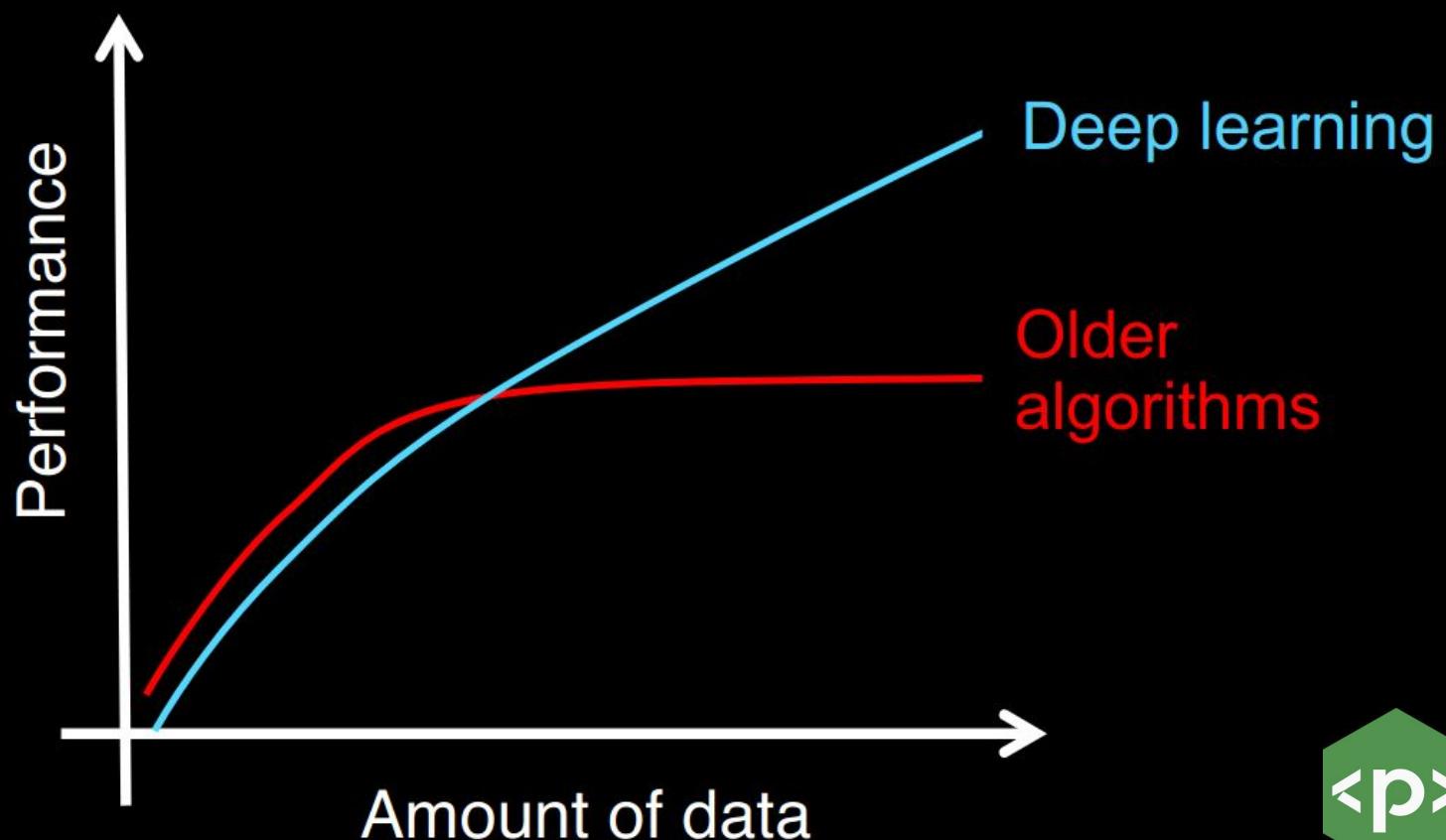
On m Training example



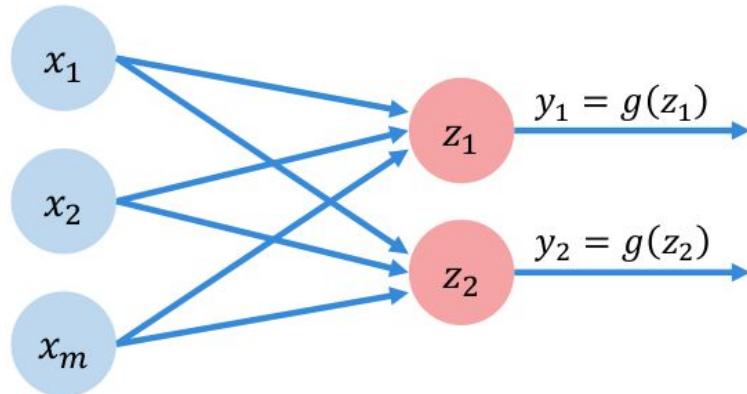
**Now Let's Dive Deeeeeeeeep in?**



## Learning from tagged data



# Multi Output Perceptron

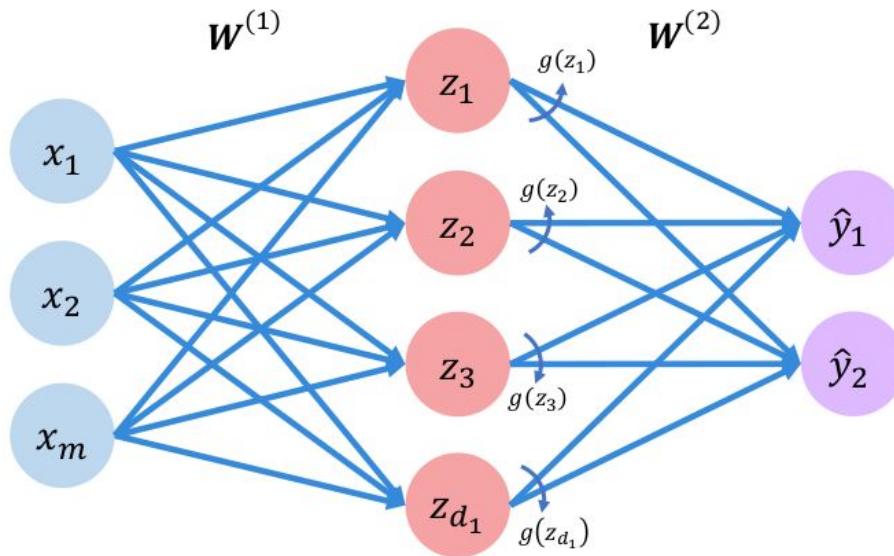


$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)



# Single Layer Neural Network



Inputs

Hidden

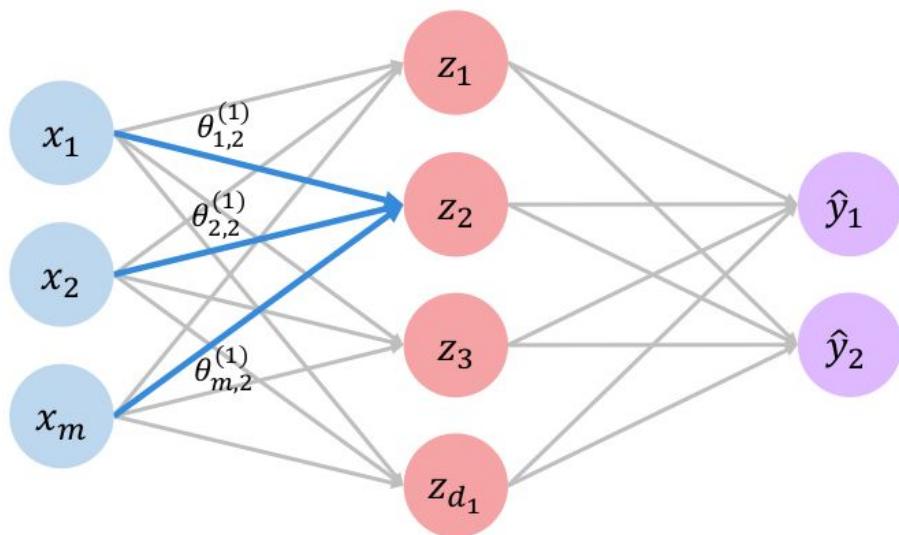
Final Output

$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)} \quad \hat{y}_i = g\left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)}\right)$$

6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)



# Single Layer Neural Network

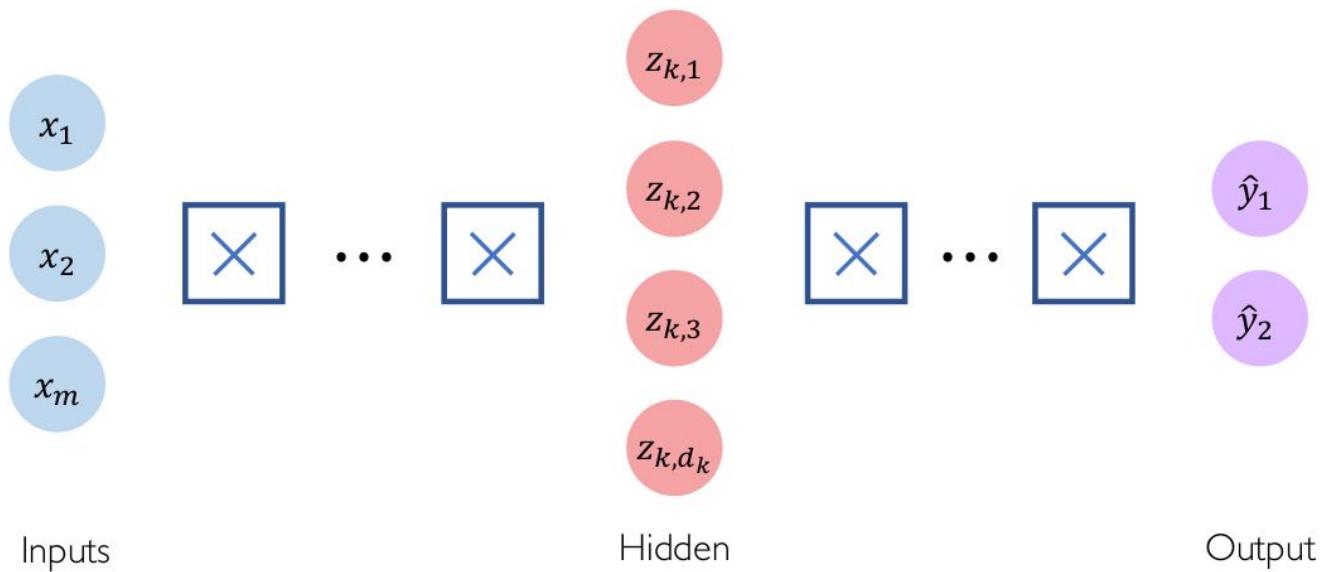


$$\begin{aligned} z_2 &= w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)} \\ &= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)} \end{aligned}$$

6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)

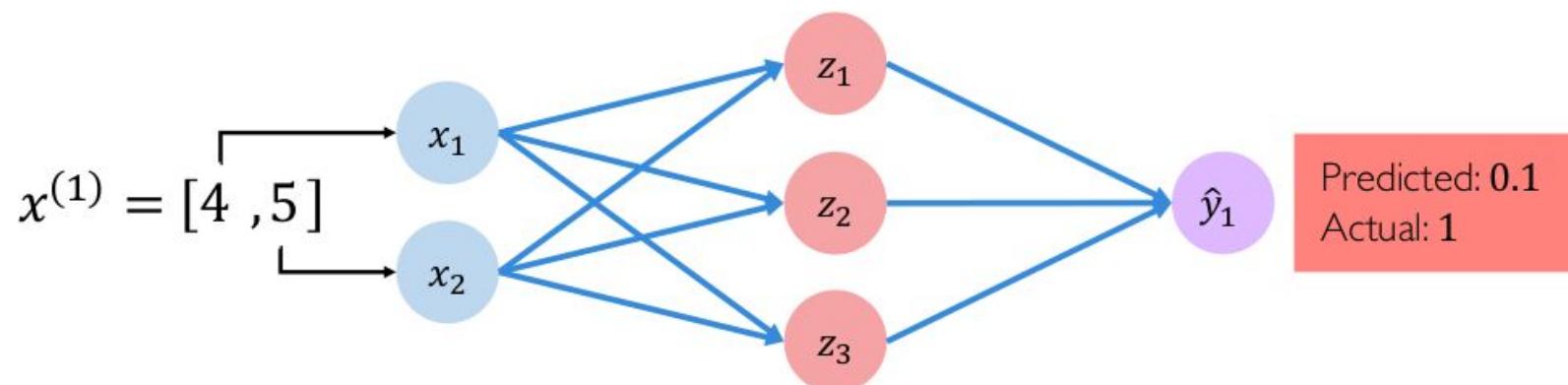


# Deep Neural Network



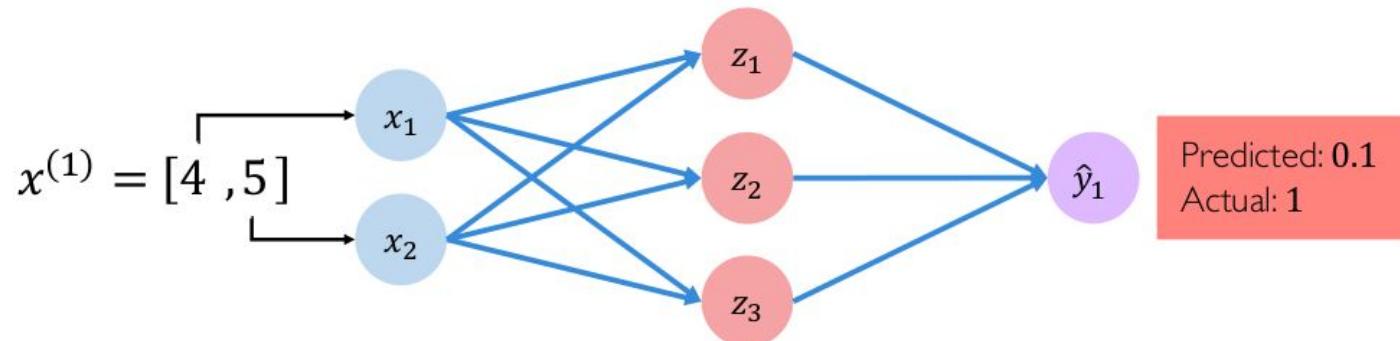
$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

# Example Problem: Will I pass this class?



# Quantifying Loss

The **loss** of our network measures the cost incurred from incorrect predictions



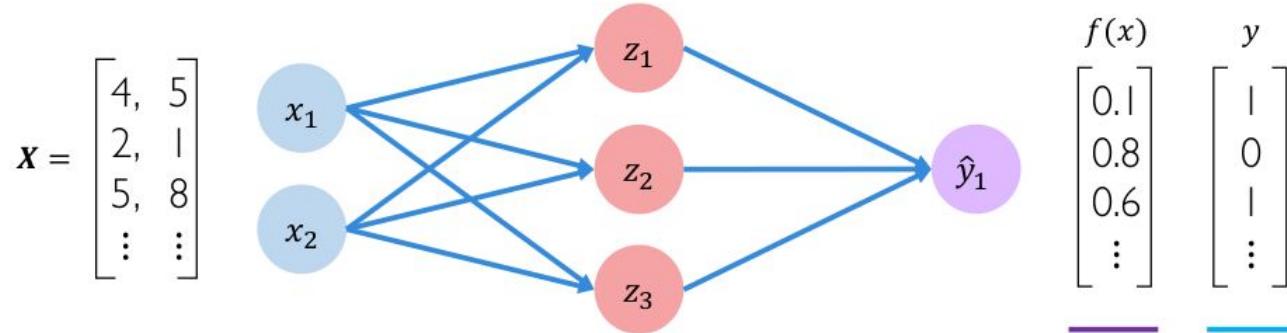
$$\mathcal{L}\left(\underline{f\left(x^{(i)}; \mathbf{W}\right)}, \underline{y^{(i)}}\right)$$

Predicted                      Actual

# Empirical Loss

(Cost Function)

The **empirical loss** measures the total loss over our entire dataset



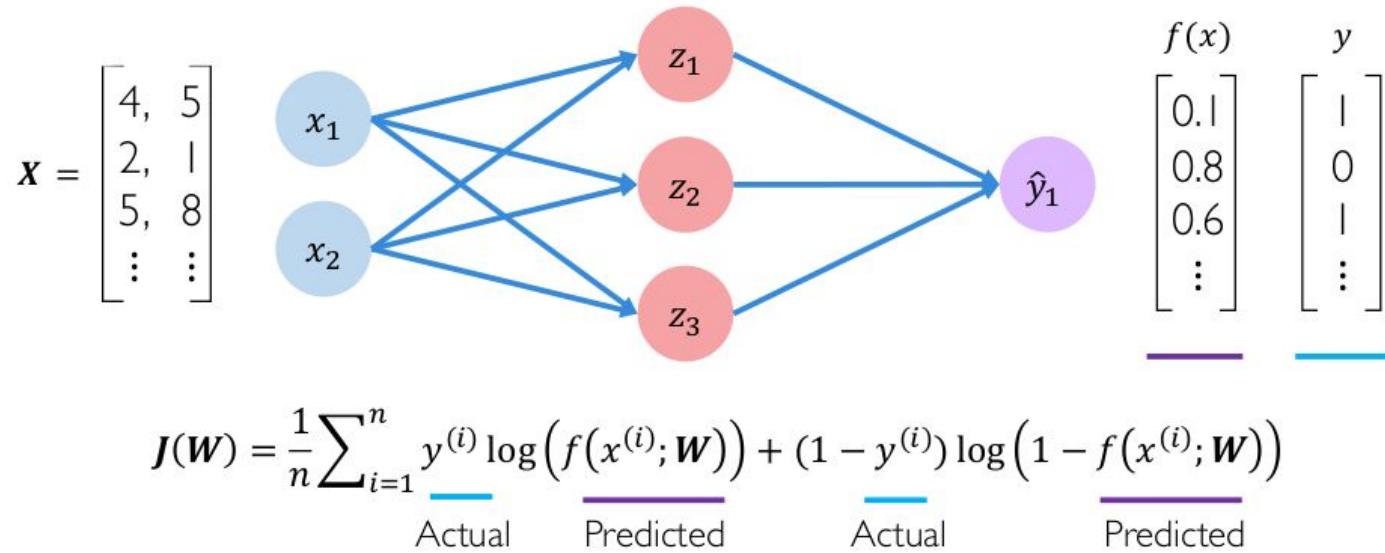
- Also known as:
- Objective function
  - Cost function
  - Empirical Risk

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

Predicted      Actual

# Binary Cross Entropy Loss

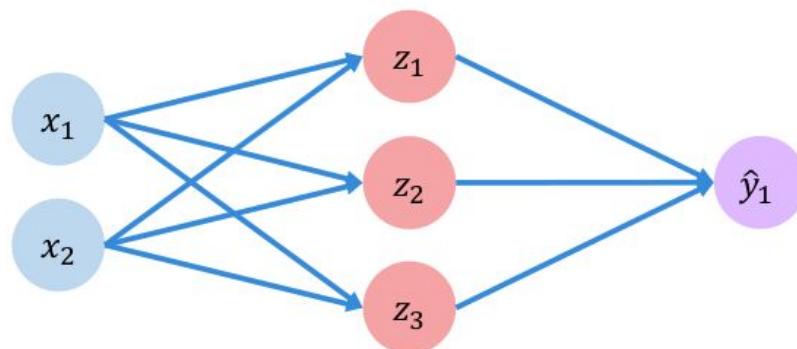
Cross entropy loss can be used with models that output a probability between 0 and 1



# Mean Squared Error Loss

**Mean squared error loss** can be used with regression models that output continuous real numbers

$$\mathbf{x} = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$$J(W) = \frac{1}{n} \sum_{i=1}^n \left( \underline{y^{(i)}} - \underline{f(x^{(i)}; \mathbf{w})} \right)^2$$

Actual      Predicted

$f(x)$	$y$
30	90
80	20
85	95
$\vdots$	$\vdots$

Final Grades  
(percentage)

# Loss Optimization

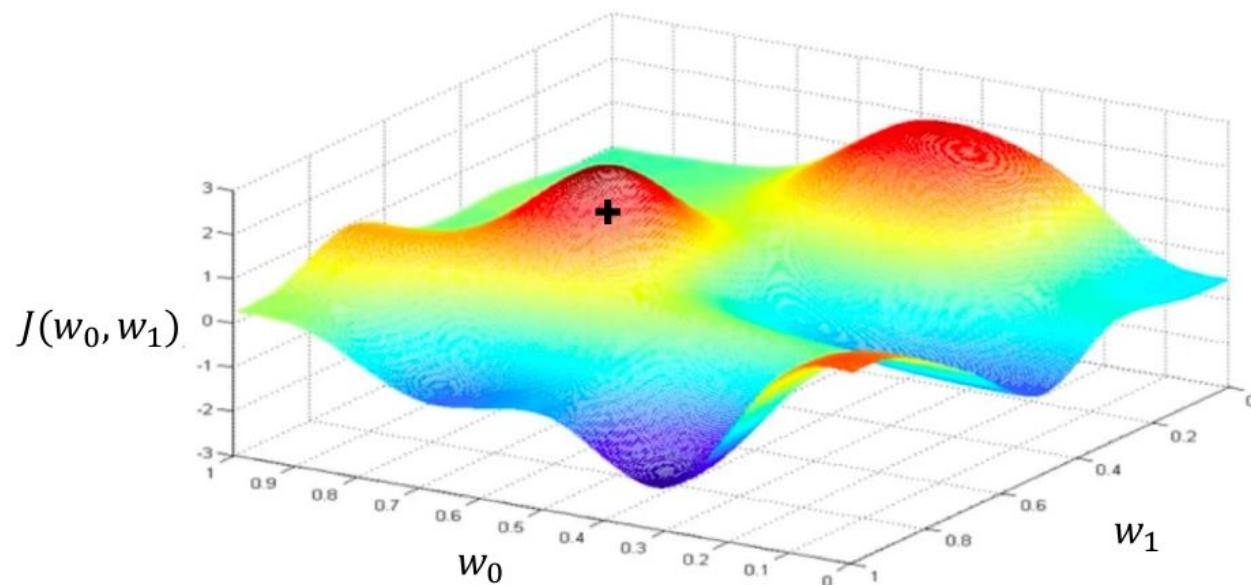
We want to find the network weights that **achieve the lowest loss**

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

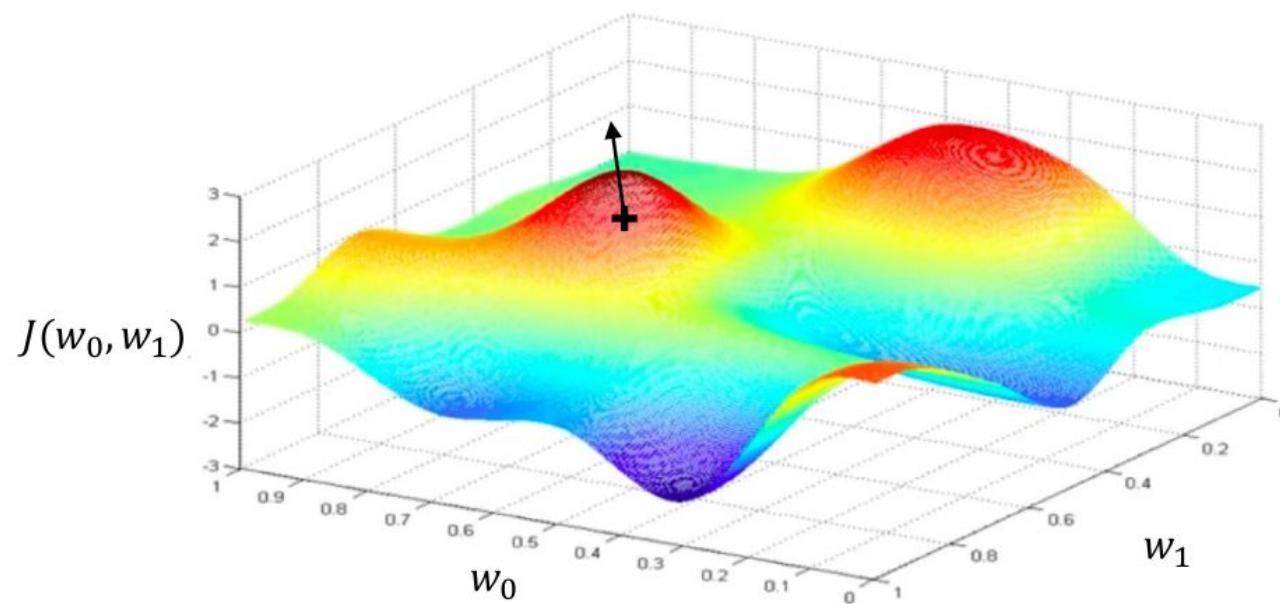
# Loss Optimization

Randomly pick an initial  $(w_0, w_1)$



# Loss Optimization

Compute gradient,  $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$

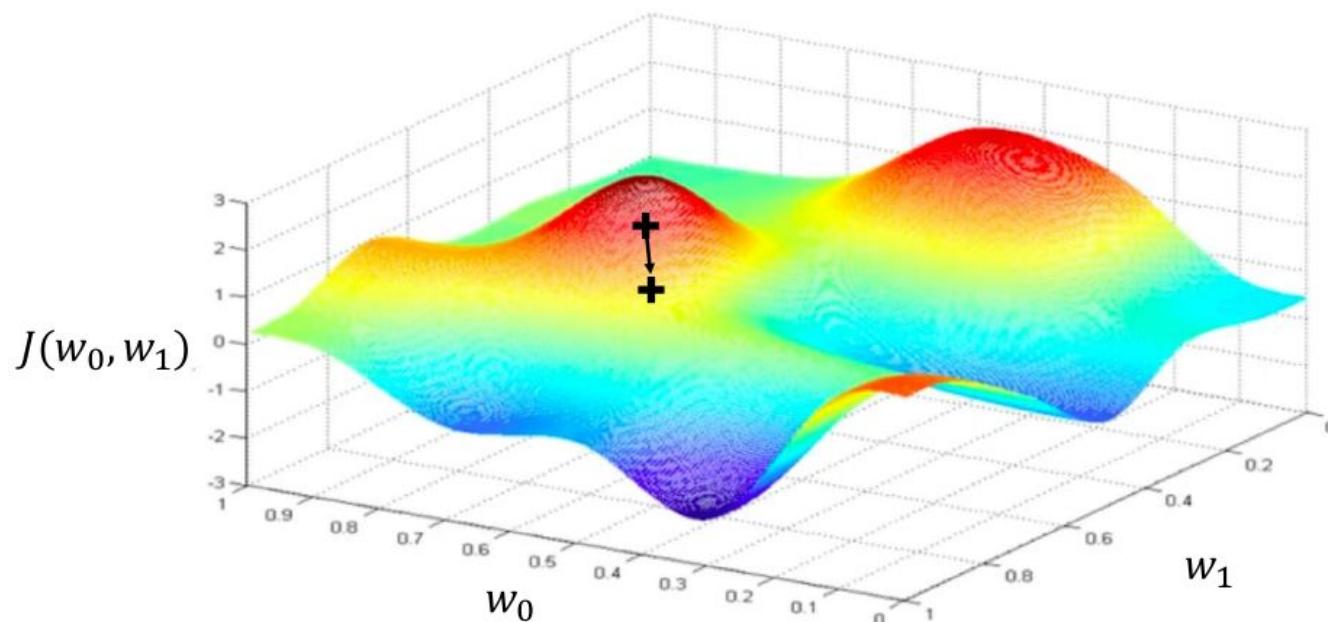


6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)



# Loss Optimization

Take small step in opposite direction of gradient

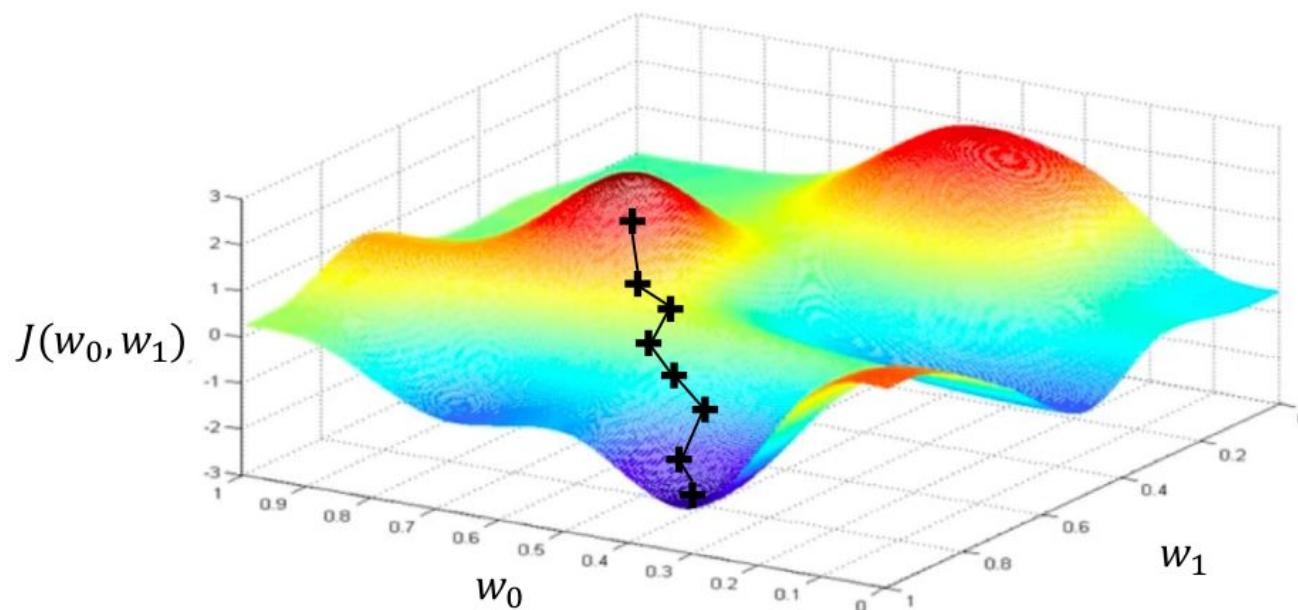


6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)



# Gradient Descent

Repeat until convergence



6.S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com)



# Gradient Descent

## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient,  $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

# Computing Gradients: Backpropagation



$$\frac{\partial J(\mathbf{W})}{\partial w_2} = \underline{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}} * \underline{\frac{\partial \hat{y}}{\partial w_2}}$$

# Computing Gradients: Backpropagation



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

— — — — —

Apply chain rule!      Apply chain rule!

# Computing Gradients: Backpropagation



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

Repeat this for **every weight in the network** using gradients from later layers

# Loss Functions Can Be Difficult to Optimize

**Remember:**

Optimization through gradient descent

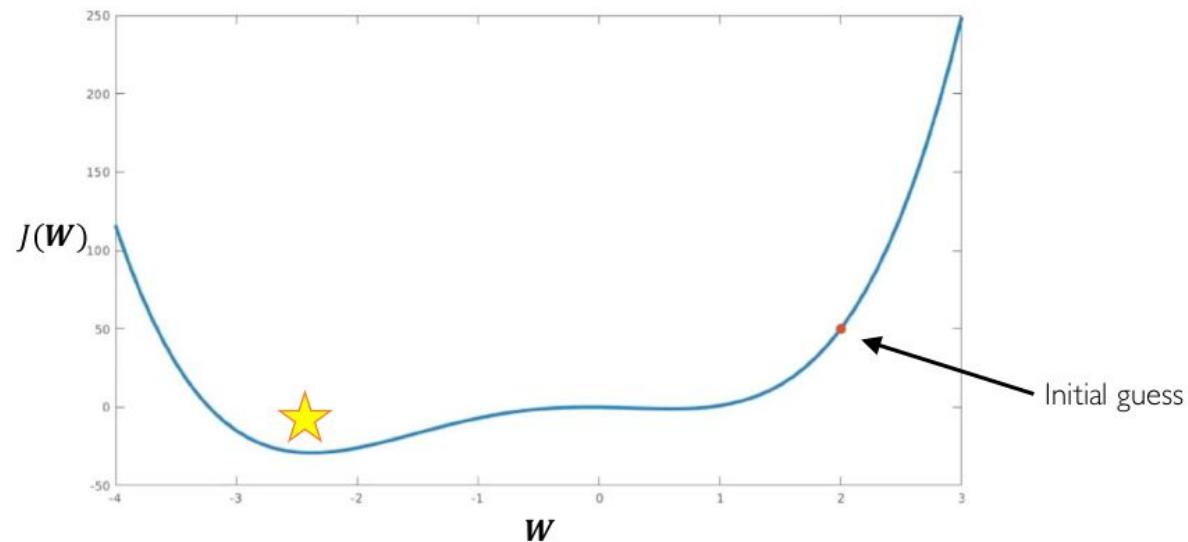
$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$

↑  
How can we set the  
learning rate?

# Setting the Learning Rate

*Large learning rates* overshoot, become unstable and diverge

*Small learning rate* converges slowly and gets stuck in false local minima



*Stable learning rates* converge smoothly and avoid local minima

# How to deal with this?

## Idea 1:

Try lots of different learning rates and see what works “just right”

## Idea 2:

Do something smarter!

Design an adaptive learning rate that “adapts” to the landscape

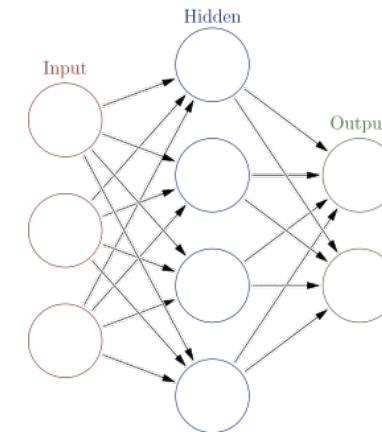
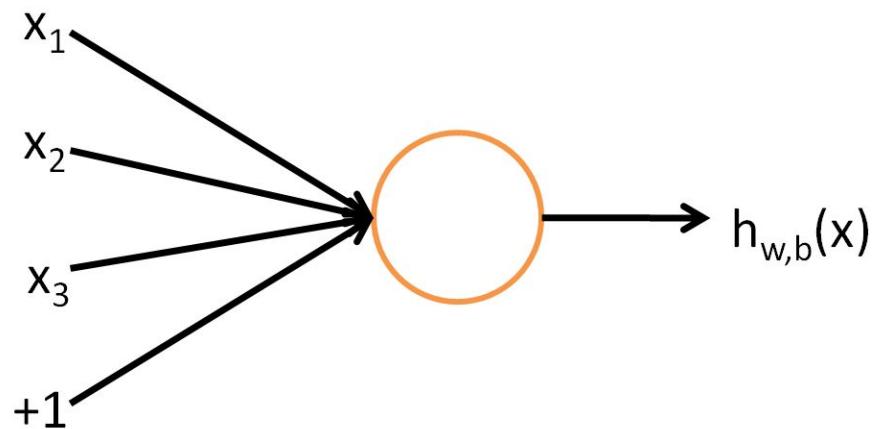
# Adaptive Learning Rates

- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
  - how large gradient is
  - how fast learning is happening
  - size of particular weights
  - etc...

**That's it for today. But just a simple intuition.**



# Why do one need non linear activation function?



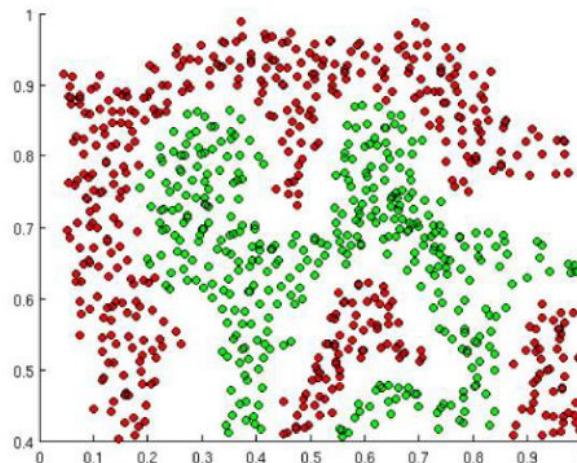
Linear hidden layer  
is useless.

We actually need an AF to include non-linearities.



# Importance of Activation Functions

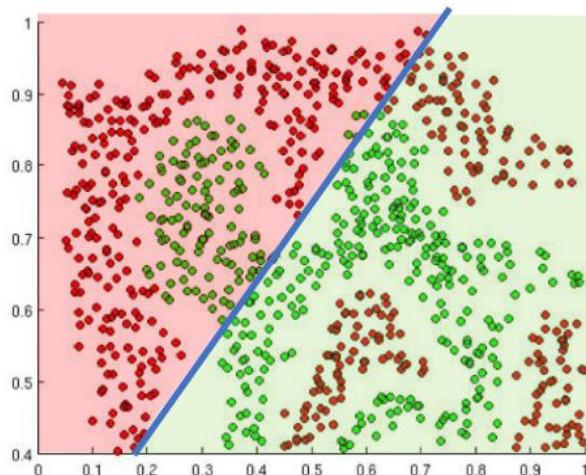
The purpose of activation functions is to **introduce non-linearities** into the network



What if we wanted to build a Neural Network to  
distinguish green vs red points?

# Importance of Activation Functions

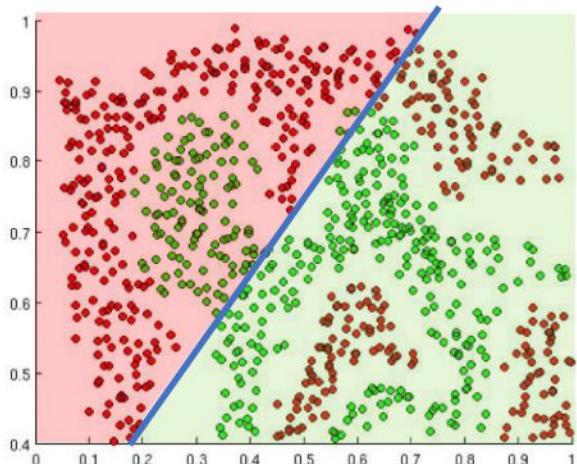
The purpose of activation functions is to **introduce non-linearities** into the network



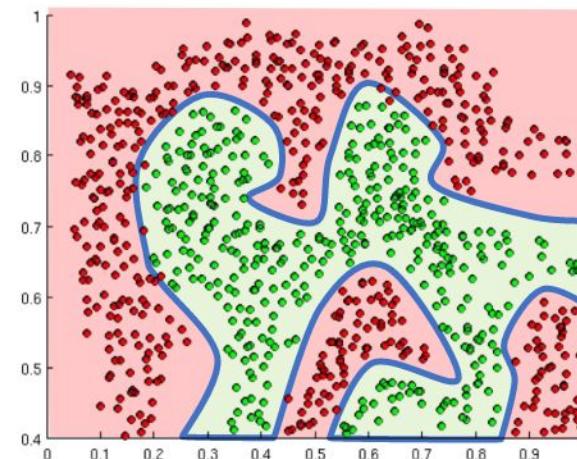
Linear Activation functions produce linear decisions no matter the network size

# Importance of Activation Functions

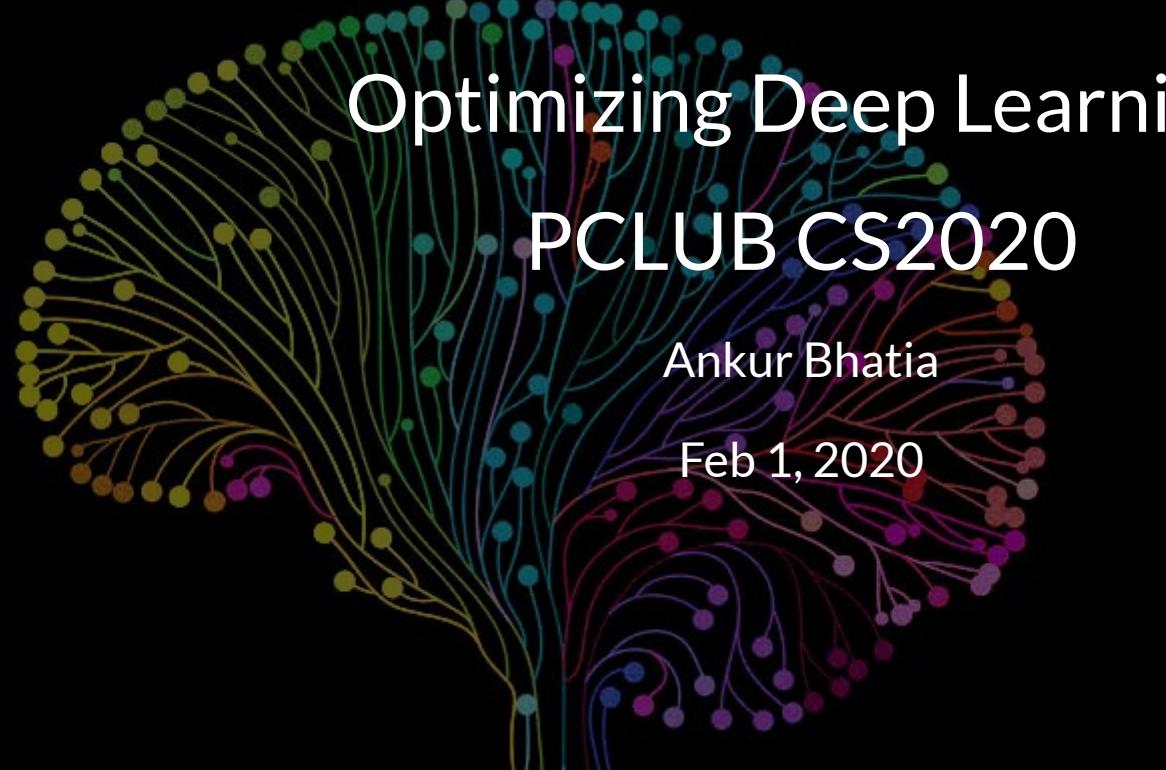
The purpose of activation functions is to **introduce non-linearities** into the network



Linear Activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions



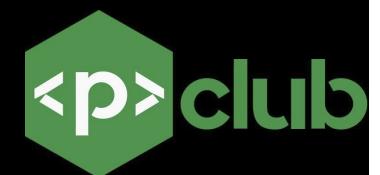
Week 2:

# Optimizing Deep Learning

## PCLUB CS2020

Ankur Bhatia

Feb 1, 2020



# Optimizing NNs.



# Batch (Vanilla) vs Mini Batch Gradient Descent

$x^{(i)}$

1. ith training example.

$z^{[l]}$

2. L - layer

$X^{\{t\}}, Y^{\{t\}}$ .

3. mini-batch

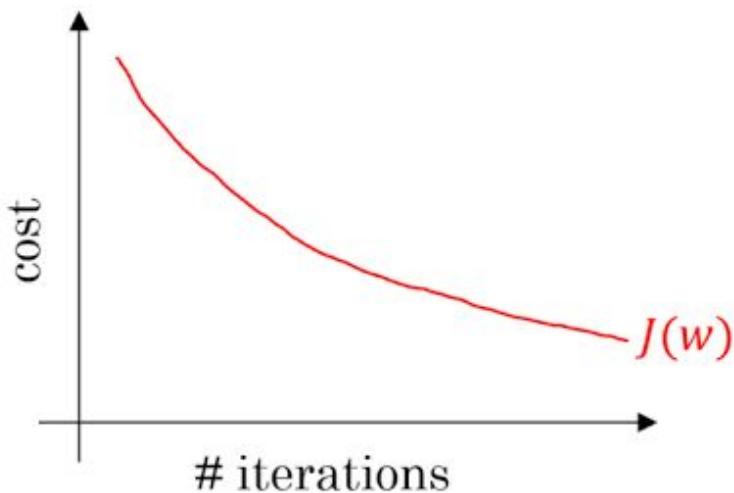


## Mini Batch Gradient Descent.

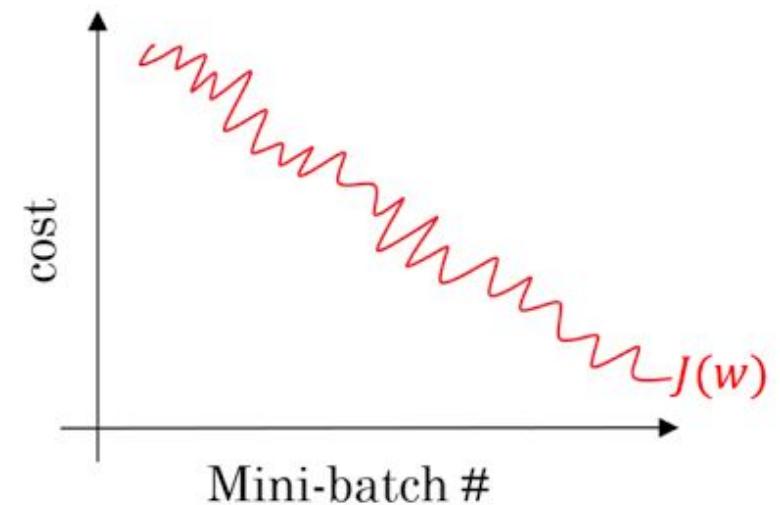
repeat  $t=0$   
for  $t=1, \dots, 5000$  {  
     $m = 5000000$  (say)  
     $\frac{m}{1000}$  (say)  
    }  
    Forward prop on  $X^{[t]}$   
     $Z^{[t]} = W^{[t]} X^{[t]} + b^{[t]}$   
     $A^{[t]} = g^{[t]}(Z^{[t]})$       } vectorized implementation  
     $\vdots$  (1000 ex)  
     $A^{[t]} = g^{[t]}(Z^{[t]})$ .  
    compute cost  $J^{[t]} = \frac{1}{1000} \sum_{n=1}^L L(\hat{y}^{[t]}, y^{[t]}) + \frac{\lambda}{2 \cdot 1000} \|W\|_F^2$   
    Back prop to compute gradients w.r.t  $J^{[t]}$  (using  $(X^{[t]}, y^{[t]})$ )  
     $W^{[t+1]} = W^{[t]} - \alpha dW^{[t]}$ ,  $b^{[t+1]} = b^{[t]} - \alpha db^{[t]}$   
    }  
    1 epoch (after full for loop completed).  
    }

Instead, we're estimating it on a small batch. Which means we're not always going in the optimal direction, because our derivatives are 'noisy'.

Batch gradient descent



Mini-batch gradient descent



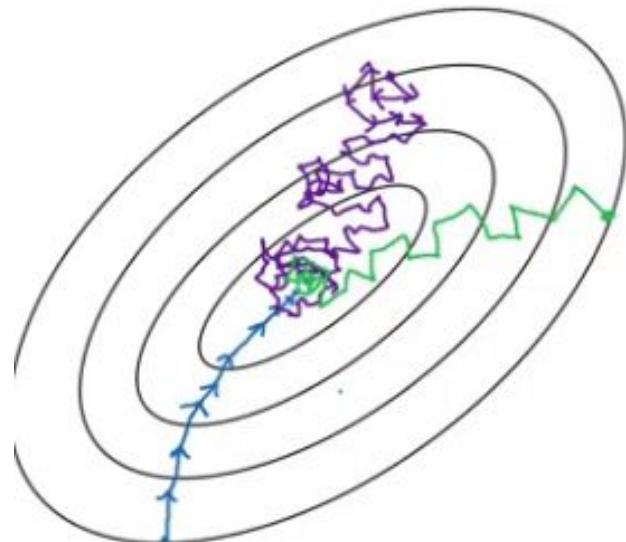
One Mini batch may be easy and other may be difficult to learn.

**Parameter to choose - Size of minibatch**

# Mini Batch Gradient Descents

If  $m = 1$ ; Stochastic Gradient Descent => Easy to compute but very noisy.

If  $m = m$ ; Batch Gradient Descent

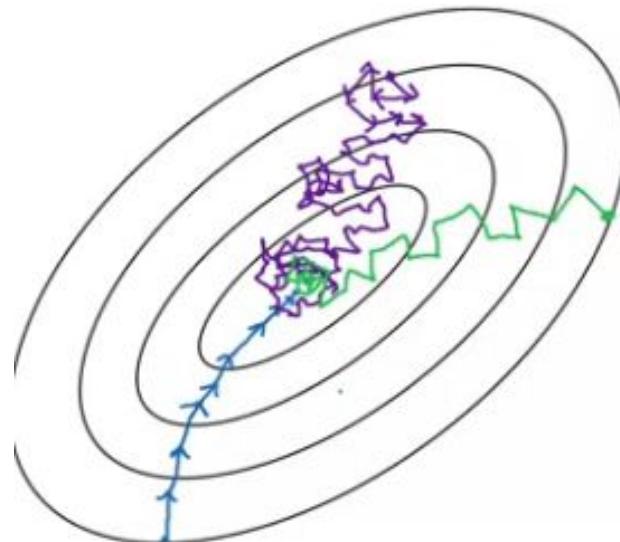


1. So What's the size to choose?..(bet 1 - m)
2. Stochastic GD loses speed of Vectorization.
3. In between mini-batch size gives:
  - a. Fast learning
  - b. Get good vectorization.
  - c. Makes progress without all data to process.
4. Batch-GD takes too long per iteration.

# Mini Batch Gradient Descents

If  $m = 1$ ; Stochastic Gradient Descent

If  $m = m$ ; Batch Gradient Descent



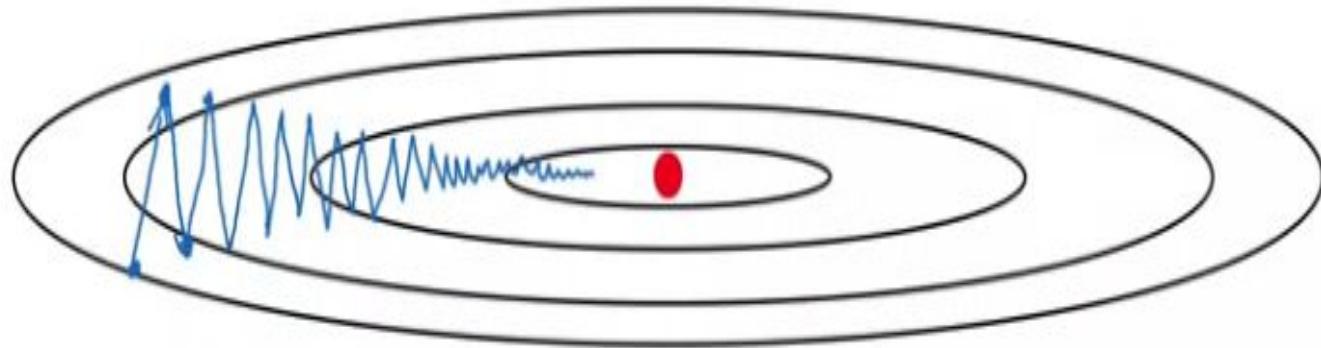
1. For small Trainx -  $m = m$ ;
2. General **mini batch sizes** = 64, 128, 256...
3. mini - batch should fit in CPU - GPU m/m.

Optimization Algorithms faster than GD.

# 1. Gradient Descent with momentum

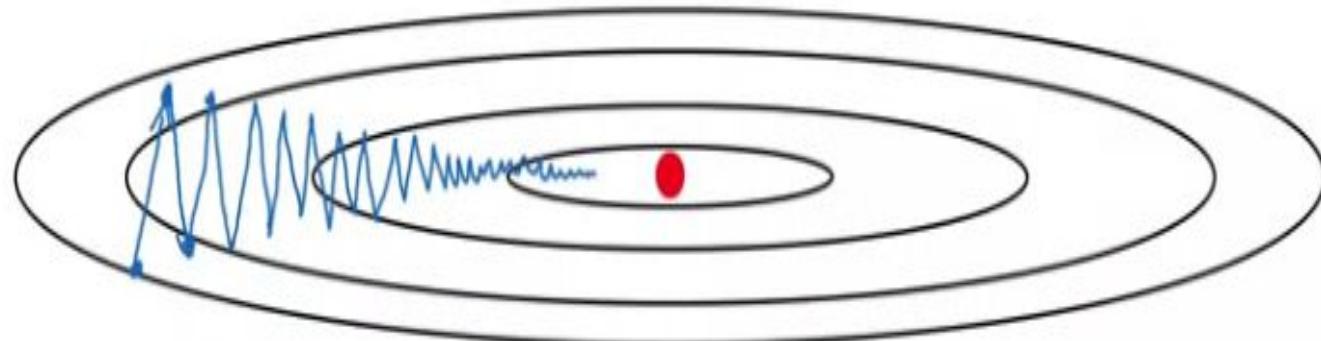
1. GD (either batch-minibatch ) takes a lot of steps to reach minima - oscillates (these osc. Slows GD).
2. That is why can't use larger LR. (overshoot, diverge) (So LR large - X)





1. GD (either batch-minibatch ) takes a lot of steps to reach minima - oscillates (these osc. Slows GD).
2. That is why can't use larger LR. (overshoot, diverge) (So LR large - X)

90 degree steps  
to ellipse.



1. Slower learning in vertical.
2. Faster learning in horizontal.

Optimization Algorithms faster than GD.

# 1. Gradient Descent with momentum

1. GD (either batch-minibatch ) takes a lot of steps to reach minima - oscillates (these osc. Slows GD).
2. That is why can't use larger LR. (overshoot, diverge) (So LR large - X)
3. What we do then?
4. Is running average a good option? Yes maybe!
5. A hyper-parameter **beta** (for the sequence of points) will be controlling the extent of rememberance average.
6. Eg - beta = 0.9 implies averaging over the last 10 gradients to get the next gradient. (it averages over  $1/(1-\text{beta})$ )



## 1. Gradient Descent with momentum : Exponentially weighted averages

Exponentially weighted averages:

Idea is to use exponentially weighted averages across the horizontal and vertical at each time step to average out and make the next steps smooth.



## 1. Gradient Descent with momentum : Exponentially weighted averages

$$V_t = \beta V_{t-1} + (1 - \beta) S_t$$
$$\beta \in [0, 1]$$

Exponentially weighted averages:

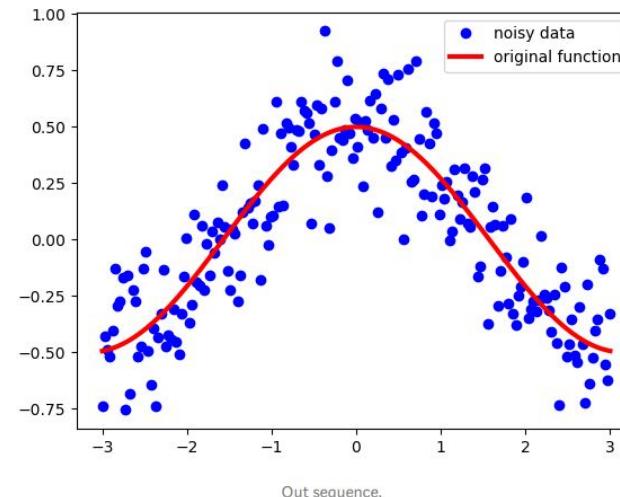
$$v_0 = 0$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$

...



## 1. Gradient Descent with momentum : Exponentially weighted averages

$$V_t = \beta V_{t-1} + (1 - \beta) S_t$$
$$\beta \in [0, 1]$$

Exponentially weighted averages:

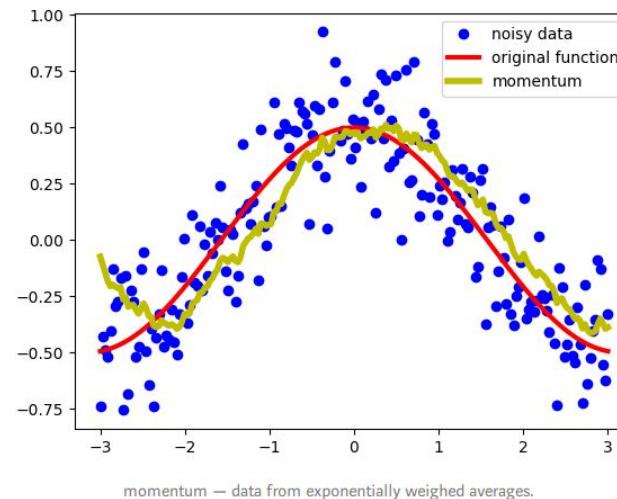
$$v_0 = 0$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$

...



<https://towardsdatascience.com/stochastic-gradient-descent-with-momentum-a84097641a5d>



## 1. Gradient Descent with momentum : Exponentially weighted averages

Implementation  $\Rightarrow$

$v_0 = 0$

Repeat {  
    get current  $\theta_t$ ;  
     $v_t = \beta v_{t-1} + (1-\beta) \theta_t$   
}     ↓  
moving averaged value according to  $\beta$   
    parameter averages over  $\frac{1}{1-\beta}$  steps

Optimization Algorithms faster than GD.

# Gradient Descent with momentum

Initializations:

1.  $V_{dw} = 0$
2.  $V_{db} = 0$

Dimensions are same as  $dW$  and  $db$

Beta = 0.9 is a common choice.

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

$$v_{dw} = \beta v_{dw} + (1 - \beta)dW$$

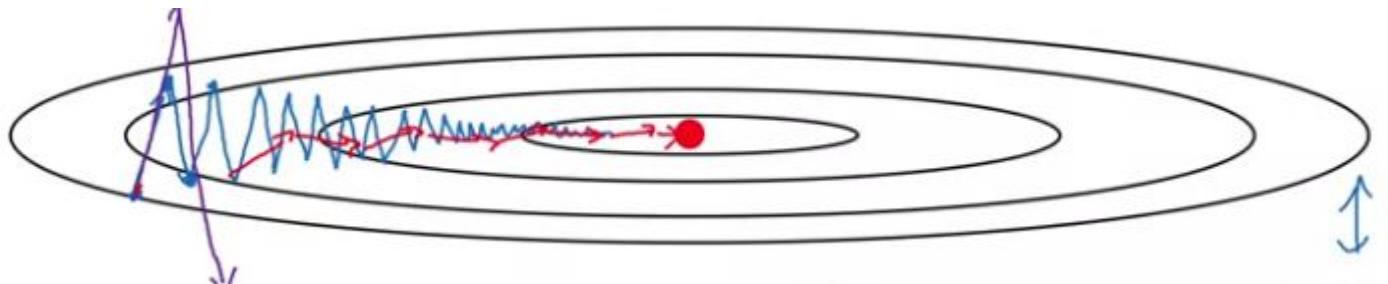
$$v_{db} = \beta v_{db} + (1 - \beta)db$$

$$W = W - \alpha v_{dw}, \quad b = b - \alpha v_{db}$$

Hyperparameters:  $\alpha, \beta$        $\beta = 0.9$



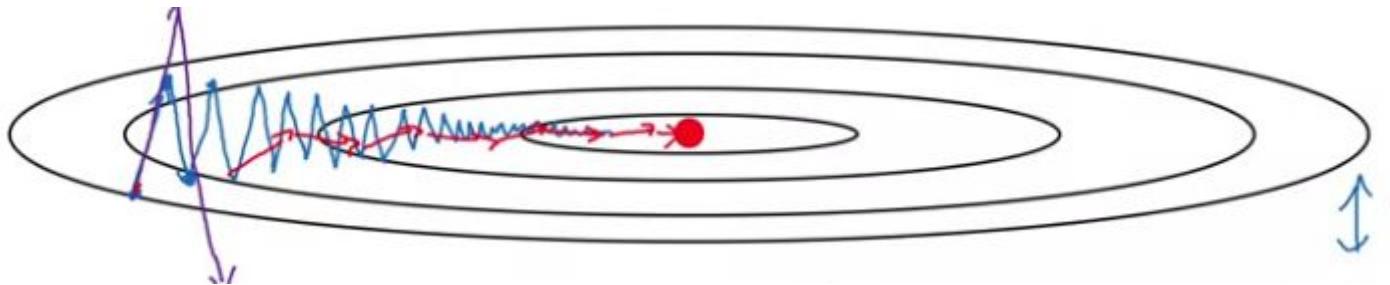
# Gradient Descent with momentum



## After Applying Momentum:

1. Moves more quickly to the horizontal direction.
2. **Oscillations are damped. (on using momentum)**
3. Takes more straight forward path.

1. Vertical average will almost be 0.
2. Horizontal average, all derivatives in horizontal direction, will be still big.



Now, Consider the intuition by: ball rolling

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

Hyperparameters:  $\alpha, \beta$        $\beta = 0.9$

1.  $\beta = 0.9$  implies averaging over the last 10 gradients to get the next gradient. (it averages over  $1/(1-\beta)$ )

# Remember: Dimensions of $\nabla d\mathbf{b}$ , $\nabla d\mathbf{w}$ ?

Now, Consider the intuition by: ball rolling

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

$$v_{dw} = \beta v_{dw} + (1 - \beta) dW \quad \text{momentum}$$
$$v_{db} = \beta v_{db} + (1 - \beta) db \quad \text{velocity}$$

$$W = W - \alpha v_{dw}, \quad b = b - \alpha v_{db}$$

Hyperparameters:  $\alpha, \beta \quad \beta = 0.9$

1. beta = 0.9 implies averaging over the last 10 gradients to get the next gradient. (it averages over  $1/(1-\beta)$ )

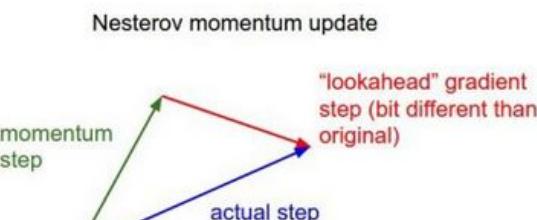
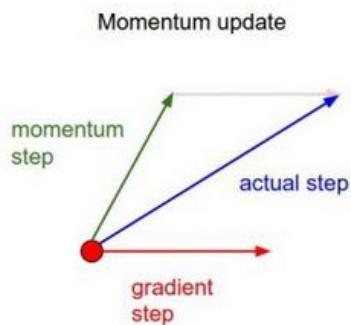


**So after using Momentum, we can also increase the learning rate a bit because the algorithm automatically tries to maintain the updates.**



# Nesterov accelerated gradient (NAG)

**Nesterov Momentum** is a slightly different version of the momentum update that has recently been gaining popularity. In this version we're first looking at a point where current momentum is pointing to and computing gradients from that point. It becomes much clearer when you look at the picture.

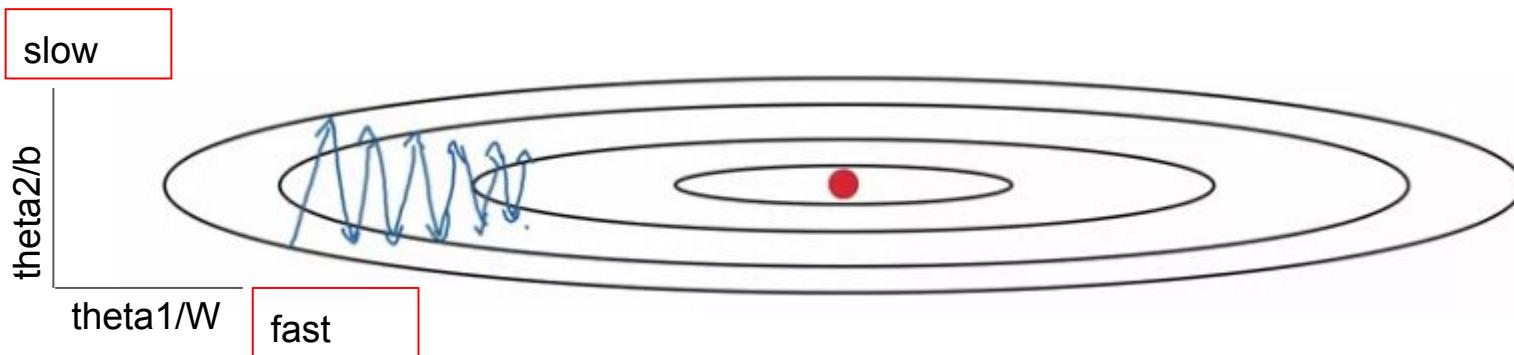


Source ([Stanford CS231n class](#))

Optimization Algorithms faster than GD.

## 2. Understanding Gradient Descent with RMS Prop in for 2D.

### Root Mean Square Propagation



## 2. Gradient Descent with Root Mean Square Propagation:

Or mentioned Sdw  
in place of Vdw

$$v_{dw} = \beta \cdot v_{dw} + (1 - \beta) \cdot dw^2$$

Exponential weighted average of squares.

$$v_{db} = \beta \cdot v_{db} + (1 - \beta) \cdot db^2$$

$$W = W - \alpha \cdot \frac{dw}{\sqrt{v_{dw}} + \epsilon}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{v_{db}} + \epsilon}$$



## 2. Gradient Descent with Root Mean Square Propagation:

Or mentioned Sdw  
in place of Vdw

$$v_{dw} = \beta \cdot v_{dw} + (1 - \beta) \cdot dw^2 \quad \text{Elementwise}$$

$$v_{db} = \beta \cdot v_{db} + (1 - \beta) \cdot db^2$$

$$W = W - \alpha \cdot \frac{dw}{\sqrt{v_{dw}} + \epsilon}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{v_{db}} + \epsilon}$$



## 2. Gradient Descent with Root Mean Square Propagation:

Or mentioned Sdw  
in place of Vdw

$$v_{dw} = \beta \cdot v_{dw} + (1 - \beta) \cdot dw^2$$

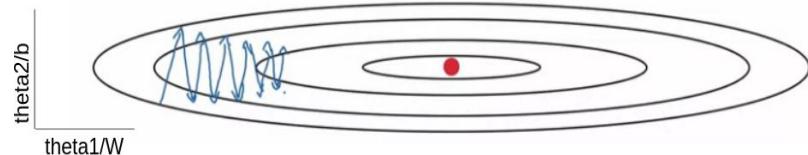
Sdw to be small

$$v_{db} = \beta \cdot v_{db} + (1 - \beta) \cdot db^2$$

Sdb to be large

$$W = W - \alpha \cdot \frac{dw}{\sqrt{v_{dw}} + \epsilon}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{v_{db}} + \epsilon}$$



## 2. Gradient Descent with Root Mean Square Propagation:

$$v_{dw} = \beta \cdot v_{dw} + (1 - \beta) \cdot dw^2$$

Vdw to be small

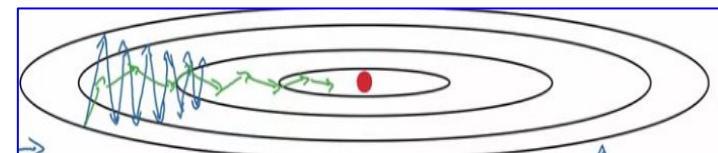
$$v_{db} = \beta \cdot v_{dw} + (1 - \beta) \cdot db^2$$

Vdb to be large

Or mentioned Sdw  
in place of Vdw

$$W = W - \alpha \cdot \frac{dw}{\sqrt{v_{dw}} + \epsilon}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{v_{db}} + \epsilon}$$



- Slopes in b direction is very large compared to w direction. => large db and small dw => Vdb large and Vdw small
- Also can use larger value of learning rate.

## 2. Gradient Descent with Root Mean Square Propagation:

Or mentioned Sdw  
in place of Vdw

$$v_{dw} = \beta \cdot v_{dw} + (1 - \beta) \cdot dw^2$$

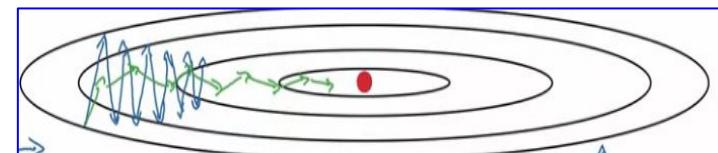
Vdw to be small

$$v_{db} = \beta \cdot v_{db} + (1 - \beta) \cdot db^2$$

Vdb to be large

$$W = W - \alpha \cdot \frac{dw}{\sqrt{v_{dw}} + \epsilon}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{v_{db}} + \epsilon}$$



- Why called RMS?
- Why is epsilon used? - In adam paper - epsilon mentioned  $10^{-8}$  so that algo doesn't divide by 0.

## 2. Gradient Descent with Root Mean Square Propagation:

Dimensions will be large.

W1, W2, W3, ..... , b1, b2, b3, .....

So damping out in vertical direction than in horizontal direction would actually mean - say damping out (or slowing the updates) in w1, w2, b5, w4 ..... direction and fastening the updates in say w3, w5, b1, b2,.... direction.



Optimization Algorithms faster than GD.

## **3. Momentum + RMSProp = Adam**

### **Adaptive Moment Estimation**

It is seen that this algorithm generalized over lots of deep learning architectures.  
Worked well on many problems.



ADAM OPTIMIZATION

### Adaptive moment estimation

$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$

On iteration  $t$ :

Compute  $dw, db$  using minibatch  $c_{t,b}$ .

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

Bias correction:

$$V_{dw}^{corrected} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{corrected} = V_{db} / (1 - \beta_2^t)$$

$$S_{dw}^{corrected} = \frac{S_{dw}}{(1 - \beta_2^t)}, \quad S_{db}^{corrected} = \frac{S_{db}}{(1 - \beta_2^t)}$$

$$w = w - \frac{\alpha V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}}$$

$$b = b - \frac{\alpha V_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$$

Hyperparameters used

$\alpha$  needs to be tuned - (Try own)

$\beta_1 = 0.9$  ( $dw$ )

$\beta_2 = 0.999$  ( $dw^2$ ) (Adam paper)

$\epsilon = 10^{-8}$  (Adam n)

Beta1 - 1st moment  
Beta2 - 2nd moment

Exponential weighted average of squares.

# Adaptive Moment Estimation

ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION - ICLR 2015

Diederik P. Kingma\*, Jimmy Lei Ba\*

<https://arxiv.org/pdf/1412.6980.pdf>

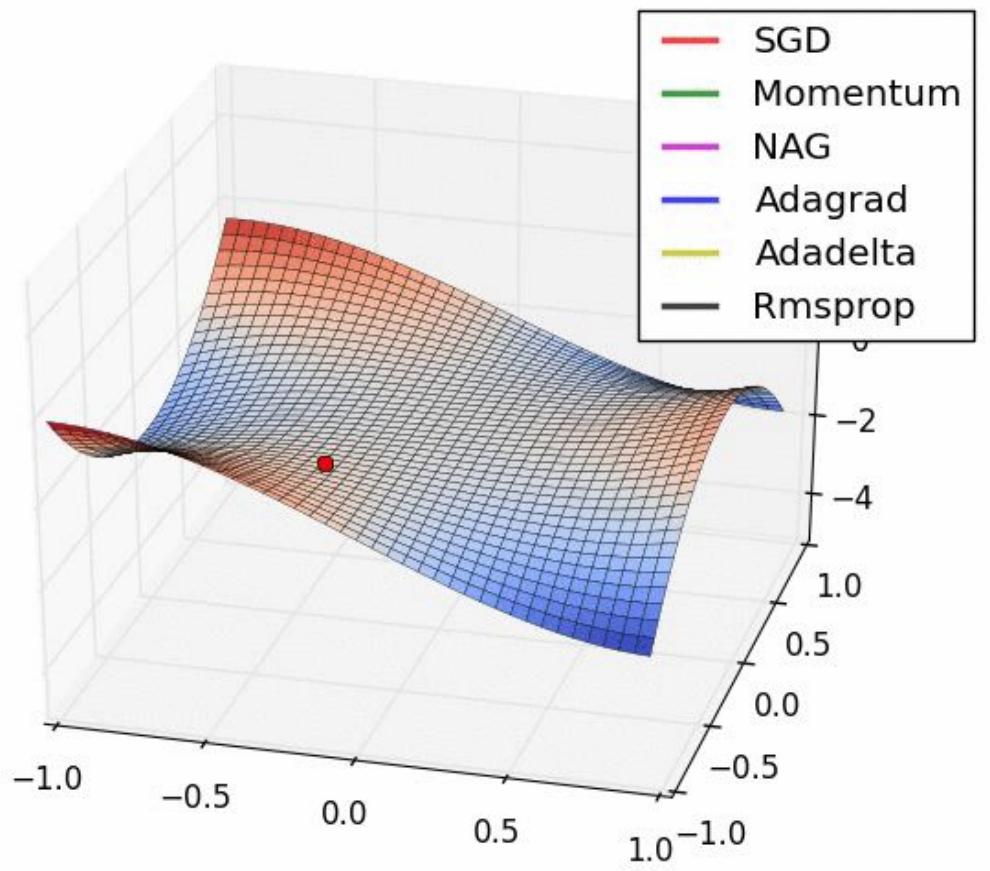
**Algorithm 1:** Adam, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation.  $g_t^2$  indicates the elementwise square  $g_t \odot g_t$ . Good default settings for the tested machine learning problems are  $\alpha = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . All operations on vectors are element-wise. With  $\beta_1^t$  and  $\beta_2^t$  we denote  $\beta_1$  and  $\beta_2$  to the power  $t$ .

```
Require:  $\alpha$ : Stepsize
Require:  $\beta_1, \beta_2 \in [0, 1]$ : Exponential decay rates for the moment estimates
Require:  $f(\theta)$ : Stochastic objective function with parameters  $\theta$ 
Require:  $\theta_0$ : Initial parameter vector
 $m_0 \leftarrow 0$  (Initialize 1st moment vector)
 $v_0 \leftarrow 0$  (Initialize 2nd moment vector)
 $t \leftarrow 0$  (Initialize timestep)
while  $\theta_t$  not converged do
     $t \leftarrow t + 1$ 
     $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$  (Get gradients w.r.t. stochastic objective at timestep  $t$ )
     $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$  (Update biased first moment estimate)
     $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$  (Update biased second raw moment estimate)
     $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  (Compute bias-corrected first moment estimate)
     $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  (Compute bias-corrected second raw moment estimate)
     $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$  (Update parameters)
end while
return  $\theta_t$  (Resulting parameters)
```

## Hyperparameters



# Various Descent's:



**Nesterov accelerated gradient**



# Applied Deep learning in one domain do not applies in other

## ITERATIVE PROCESS

1. Depends on amount of data.
2. GPU/CPU
3. No. of input features
4. Type of data

Still we cannot choose the hyperparameters  
well.



# Applied Deep learning in one domain do not applies in other

## ITERATIVE PROCESS

1. Setting up data well in train, test and dev.
2. Splitting percentages.



# Splitting up of Train/Test/Dev sets

1. Smaller sets. - 100,  $10^3$ ,  $10^4$
2. Larger sets (Big Data) -  $> 10^6$



# Distribution of Train and Test Sets.

## Examples

1. Effects of online vs offline data.
2. Effects of user mischief.

Dev and Test sets should be from same distributions.

Dev set / cross validation set / development set

Test set not necessary.  
/Call test set dev set.



# Bias and Variance

## BIAS:

**Bias** is disproportionate weight *in favor of* or *against* an idea or thing, usually in a way that is **closed-minded**, **prejudicial**, or unfair. Biases can be innate or learned. People may develop biases for or against an individual, a group, or a belief.<sup>[1]</sup> In science and engineering, a bias is a **systematic error**. **Statistical bias** results from an unfair **sampling** of a population, or from an **estimation** process that does not give accurate results on average.<sup>[2]</sup>



# Bias and Variance

## VARIANCE:

In probability theory and statistics, **variance** is the expectation of the squared deviation of a random variable from its mean. Informally, it measures how far a set of (random) numbers are spread out from their average value. Variance has a central role in statistics, where some ideas that use it include descriptive statistics, statistical inference, hypothesis testing, **goodness of fit**, and Monte

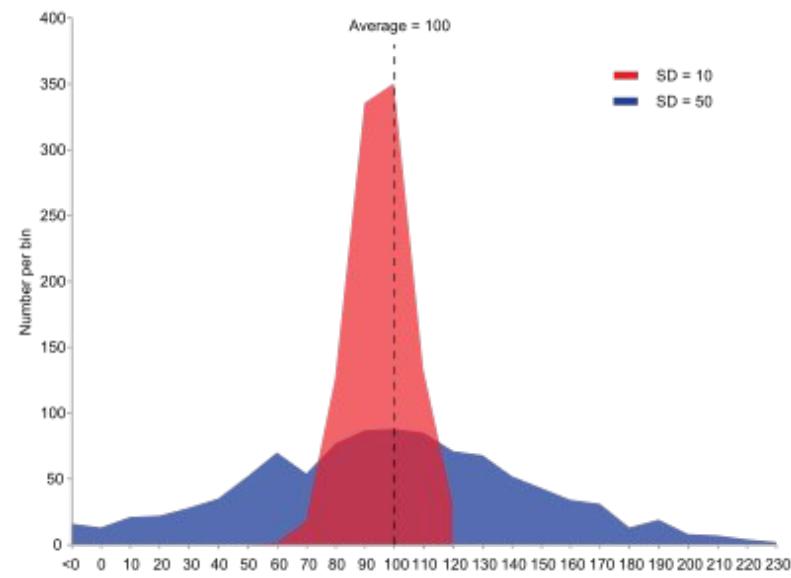


# Bias and Variance

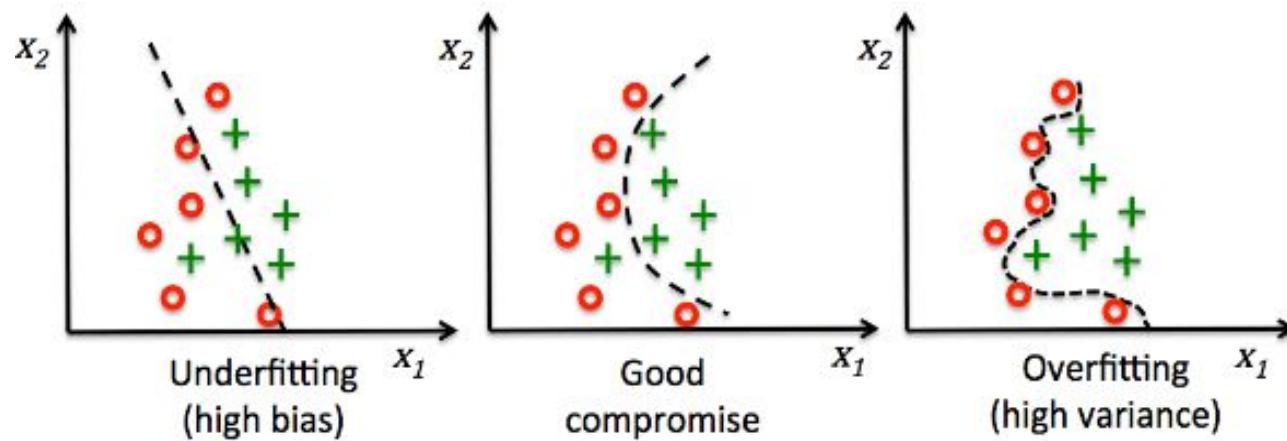
## VARIANCE:

The variance of a random variable  $X$  is the [expected value](#) of the squared deviation from the [mean](#) of  $X$ ,  $\mu = \text{E}[X]$ :

$$\text{Var}(X) = \text{E}[(X - \mu)^2].$$



# Bias and Variance



Logistic Reg. is not  
a good fit.

# Bias and Variance

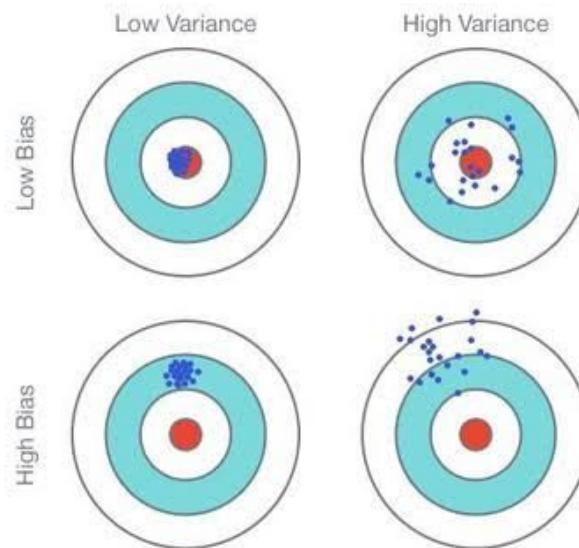
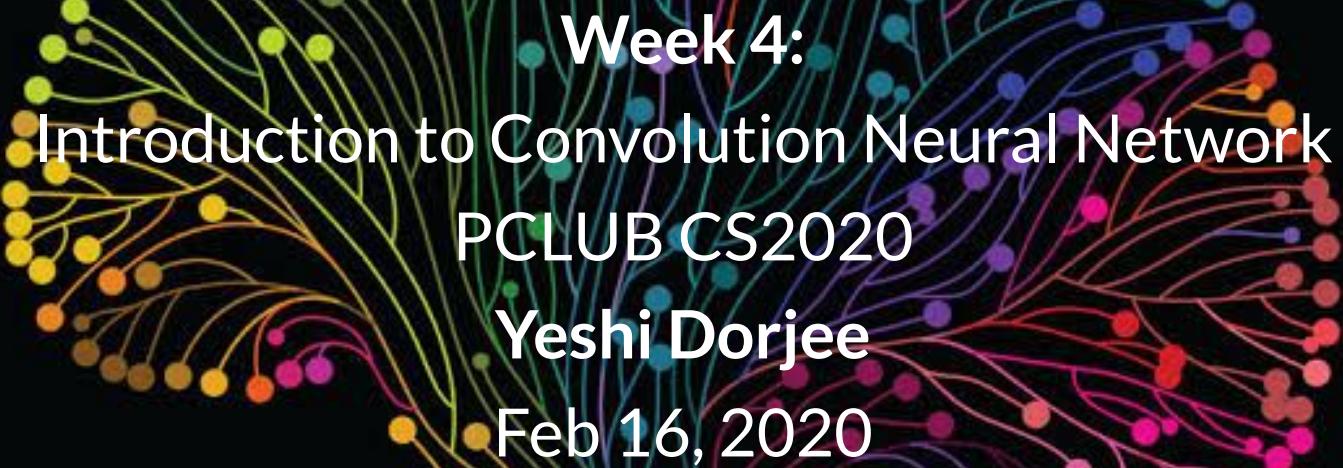


Fig. 1: Graphical Illustration of bias-variance trade-off , Source: Scott Fortmann-Roe., Understanding Bias-Variance Trade-off



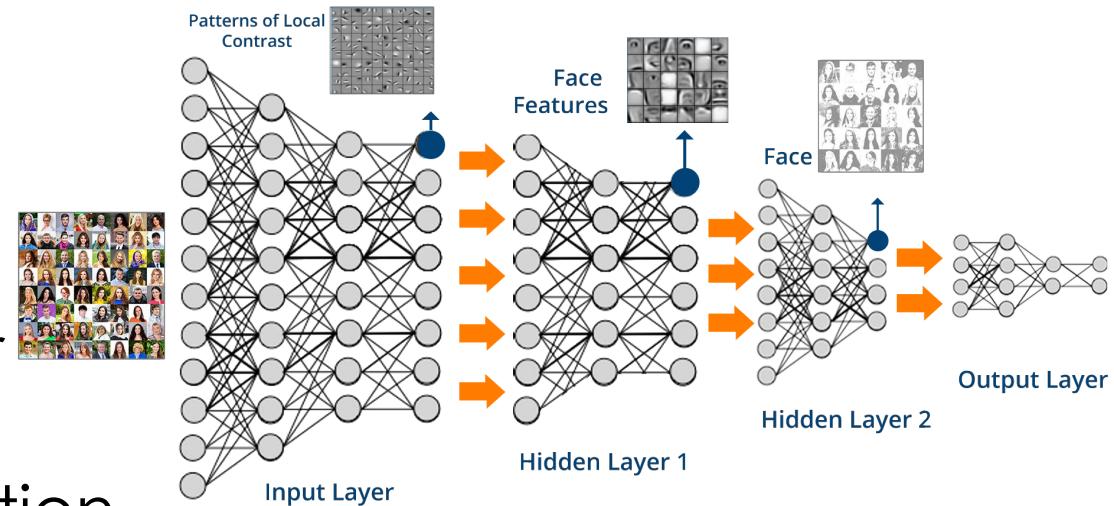


**Week 4:**  
Introduction to Convolution Neural Network  
PCLUB CS2020  
**Yeshi Dorjee**  
Feb 16, 2020

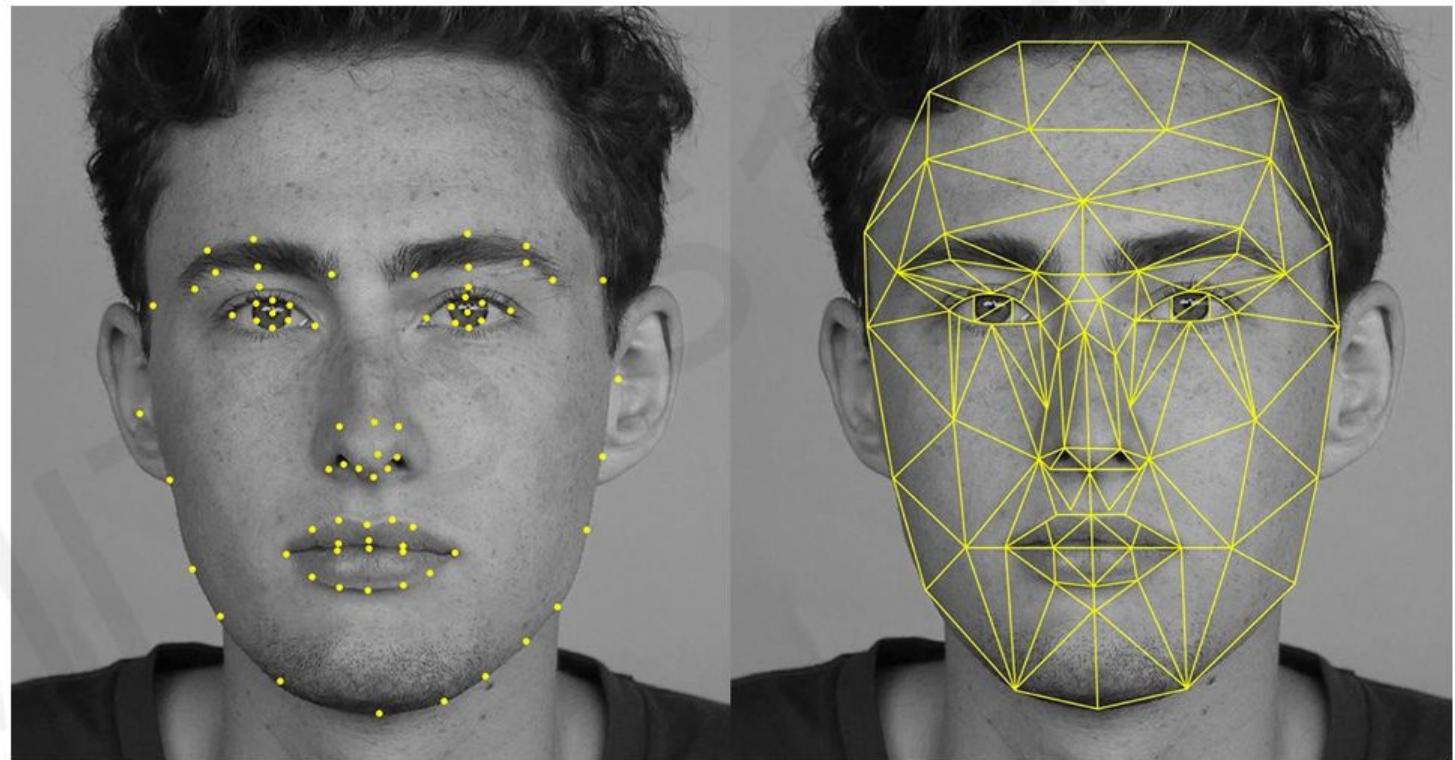
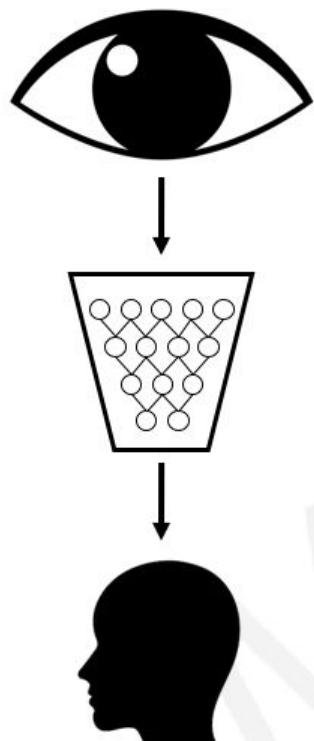


# Overview

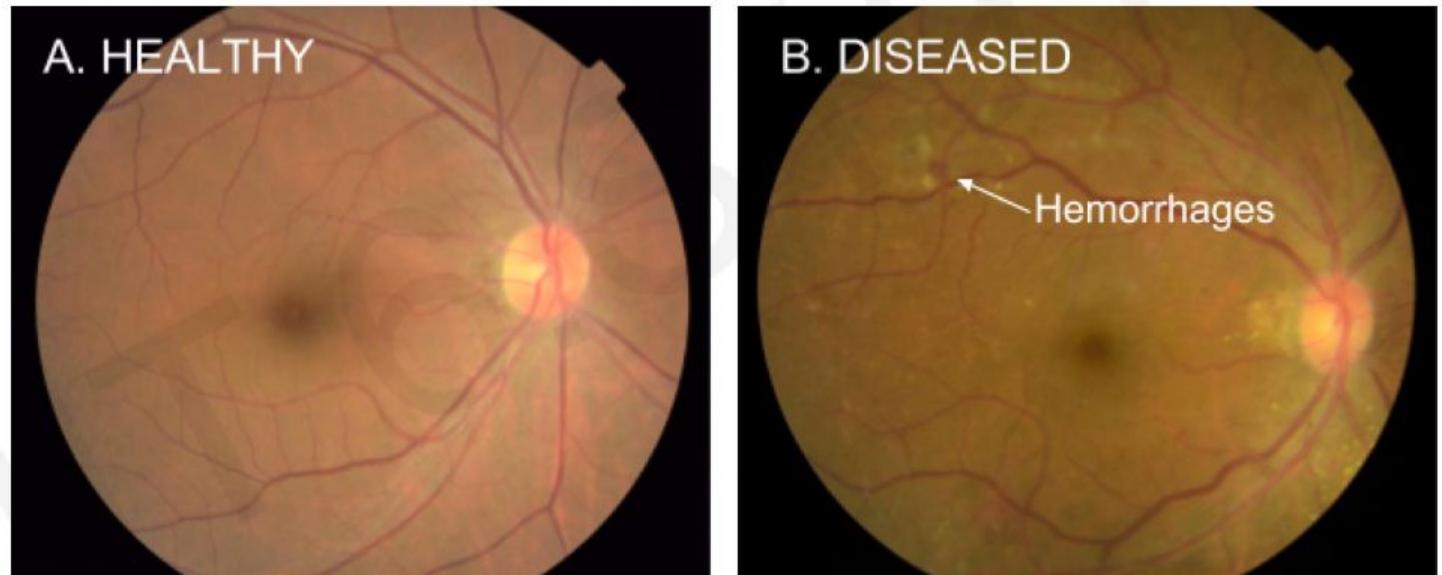
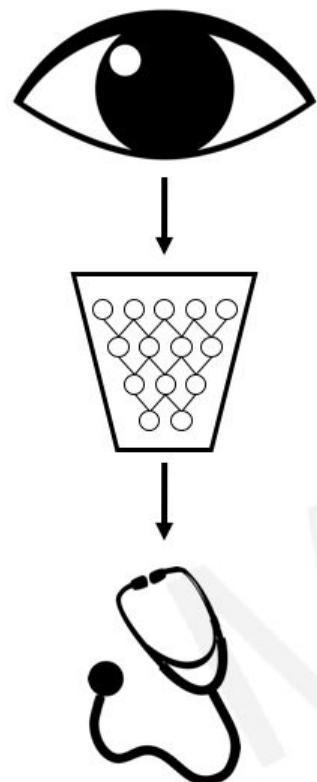
- Introduction
- Convolution Layer
- Pooling
- Image Augmentation
- Classical Network(Case Study)
- ResNets



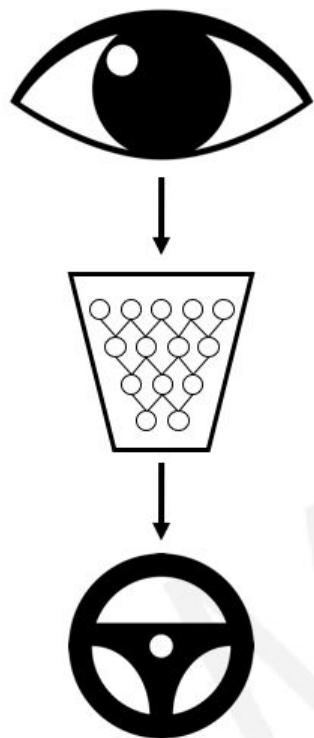
# Impact: Facial Detection & Recognition

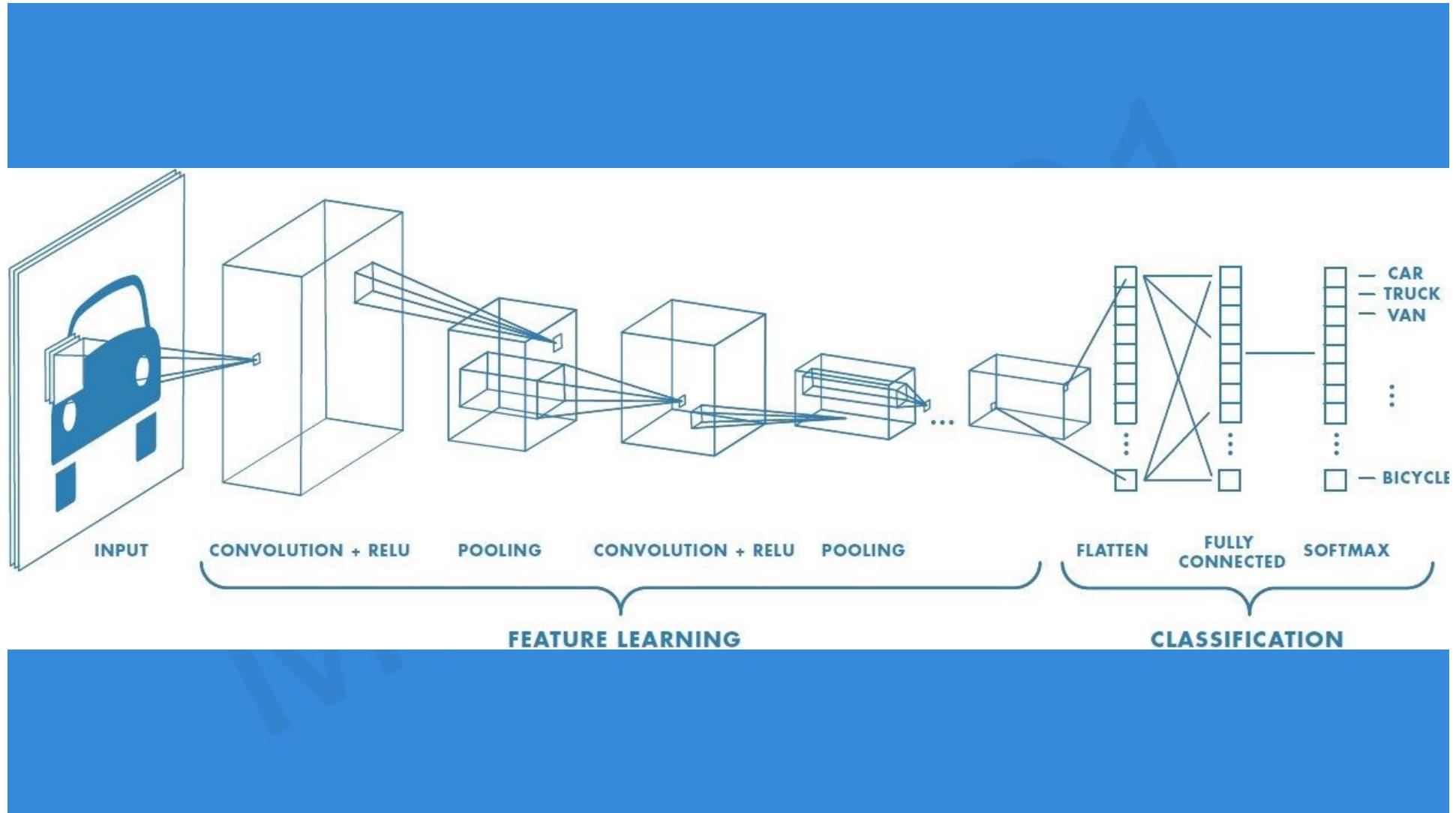


# Impact: Medicine, Biology, Healthcare

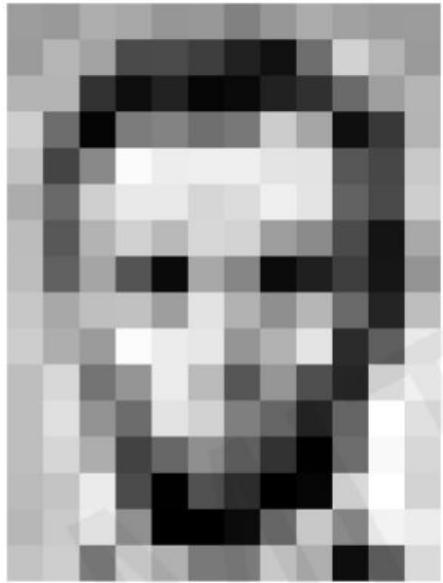


# Impact: Self-Driving Cars

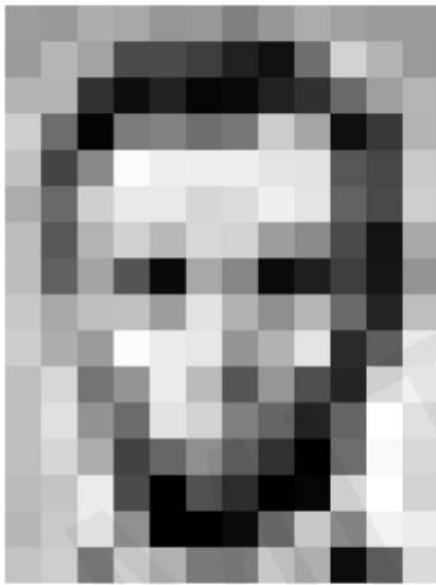




# Images are Numbers

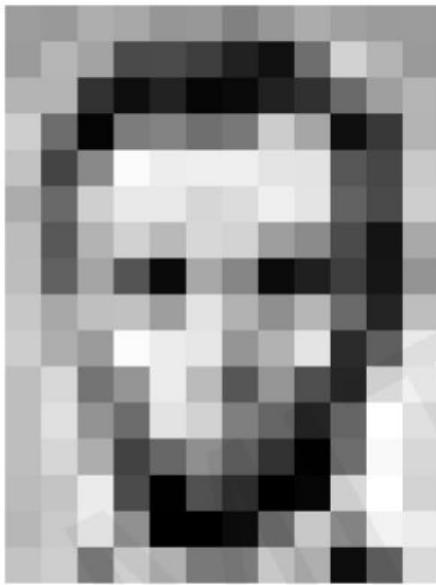


# Images are Numbers



157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	116	210	180	154
180	180	50	14	84	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	197	251	257	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	86	150	79	38	218	241
190	224	147	108	227	210	127	102	36	105	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	238	75	1	81	47	0	6	217	255	211
183	202	237	145	6	9	12	109	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

# Images are Numbers



157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	116	210	180	154
180	180	50	14	84	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	105	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	109	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

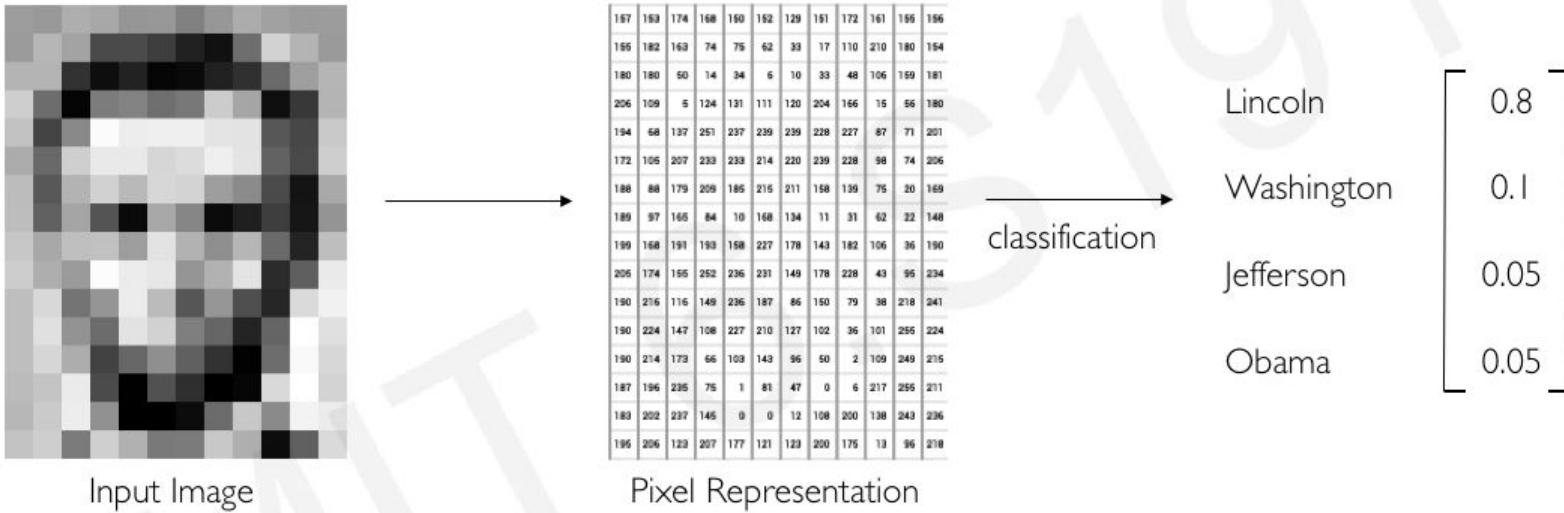
What the computer sees

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	84	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	105	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	109	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

An image is just a matrix of numbers [0,255]!

i.e., 1080x1080x3 for an RGB image

# Tasks in Computer Vision



- **Regression:** output variable takes continuous value
- **Classification:** output variable takes class label. Can produce probability of belonging to a particular class

# High Level Feature Detection

Let's identify key features in each image category



Nose,  
Eyes,  
Mouth



Wheels,  
License Plate,  
Headlights



Door,  
Windows,  
Steps

# Manual Feature Extraction



Problems?

# Manual Feature Extraction



# Manual Feature Extraction

Domain knowledge

Define features

Detect features  
to classify

Viewpoint variation



Scale variation



Deformation



Occlusion



Illumination conditions



Background clutter



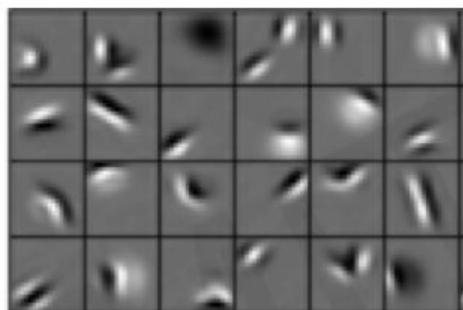
Intra-class variation



# Learning Feature Representations

Can we learn a **hierarchy of features** directly from the data instead of hand engineering?

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

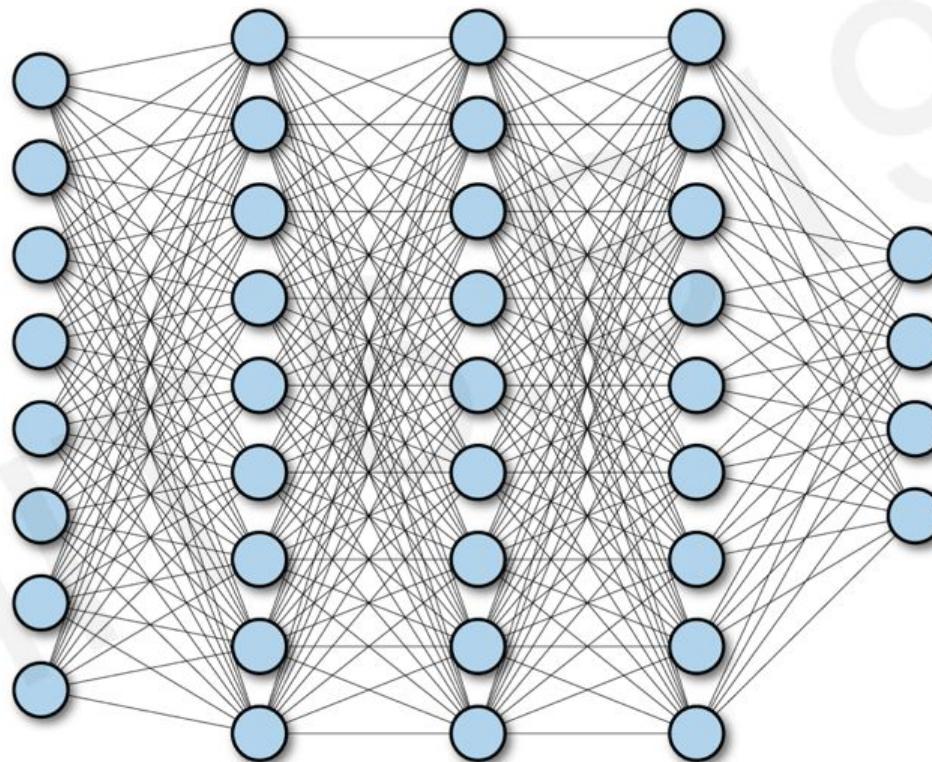
High level features



Facial structure

# Learning Visual Features

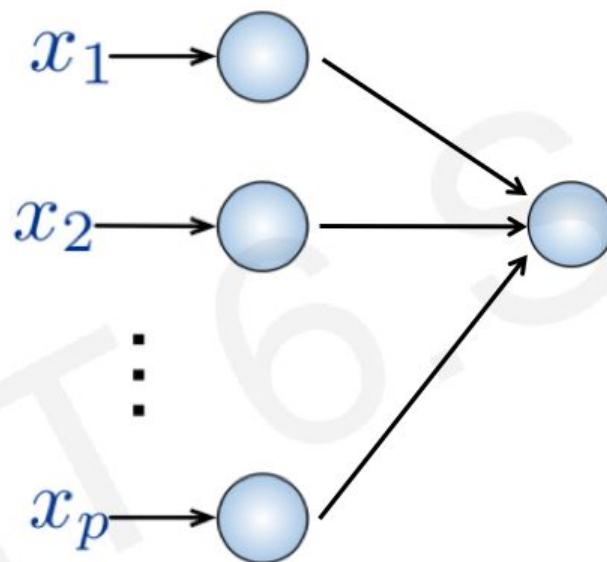
# Fully Connected Neural Network



# Fully Connected Neural Network

## Input:

- 2D image
- Vector of pixel values



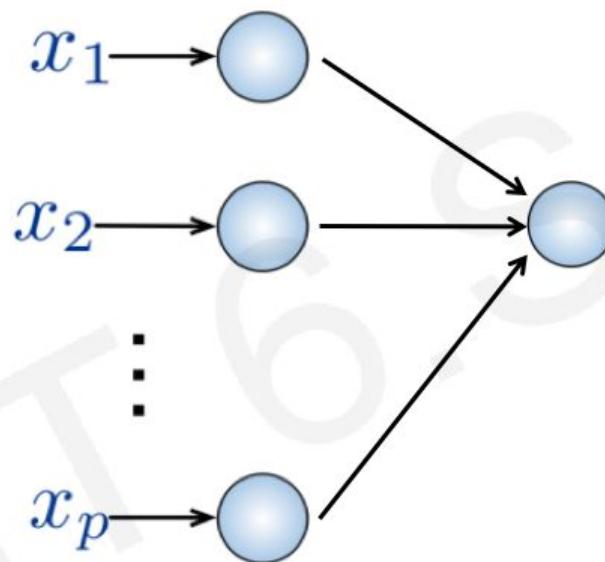
## Fully Connected:

- Connect neuron in hidden layer to all neurons in input layer
- No spatial information!
- And many, many parameters!

# Fully Connected Neural Network

## Input:

- 2D image
- Vector of pixel values



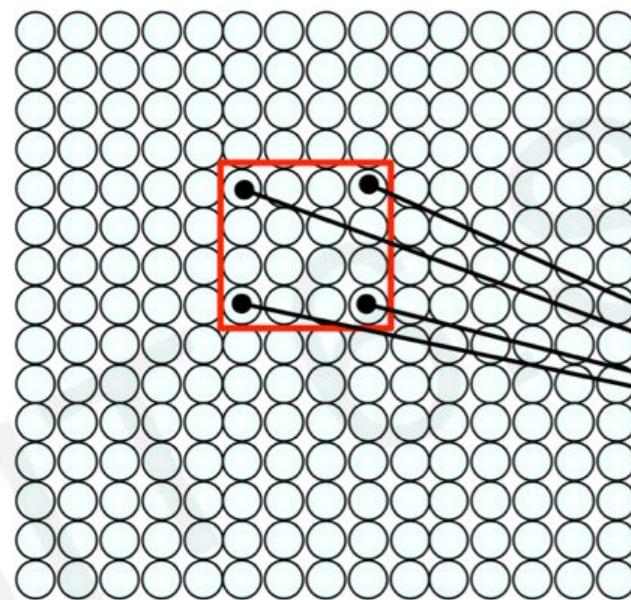
## Fully Connected:

- Connect neuron in hidden layer to all neurons in input layer
- No spatial information!
- And many, many parameters!

How can we use **spatial structure** in the input to inform the architecture of the network?

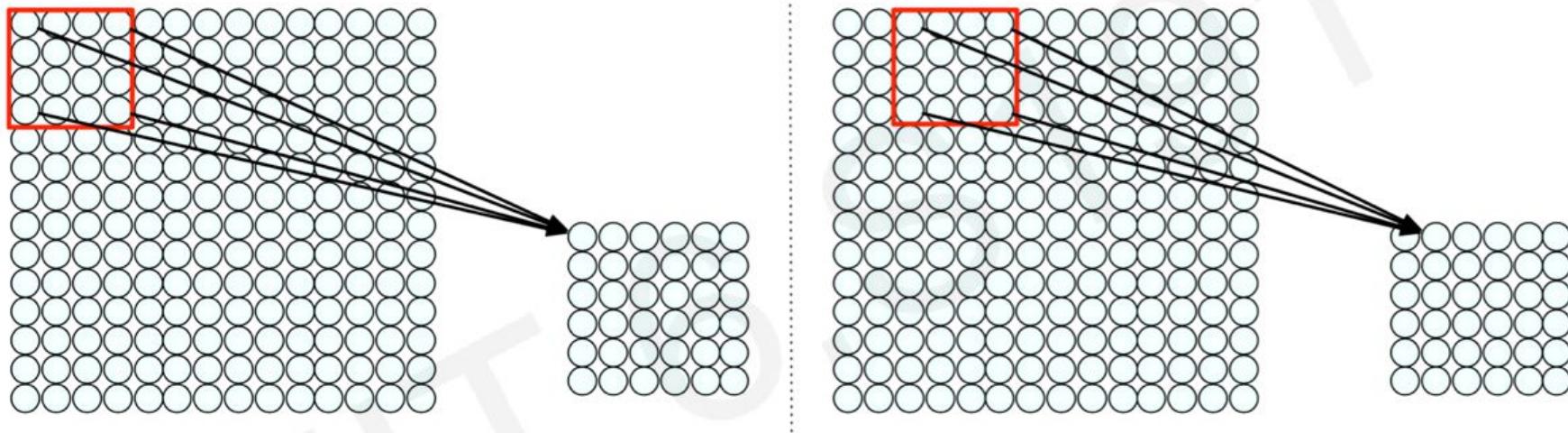
# Using Spatial Structure

**Input:** 2D image.  
Array of pixel values



**Idea:** connect patches of input  
to neurons in hidden layer:  
Neuron connected to region of  
input. Only "sees" these values.

# Using Spatial Structure



Connect patch in input layer to a single neuron in subsequent layer.

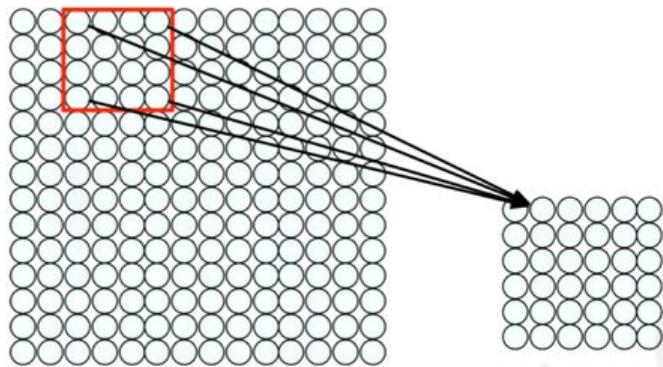
Use a sliding window to define connections.

How can we **weight** the patch to detect particular features?

# Applying Filters to Extract Features

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) Spatially **share** parameters of each filter  
(features that matter in one part of the input should matter elsewhere)

# Feature Extraction with Convolution



- Filter of size  $4 \times 4$  : 16 different weights
- Apply this same filter to  $4 \times 4$  patches in input
- Shift by 2 pixels for next patch

This “patchy” operation is **convolution**

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

# Feature Extraction and Convolution

## A Case Study

# X or X?

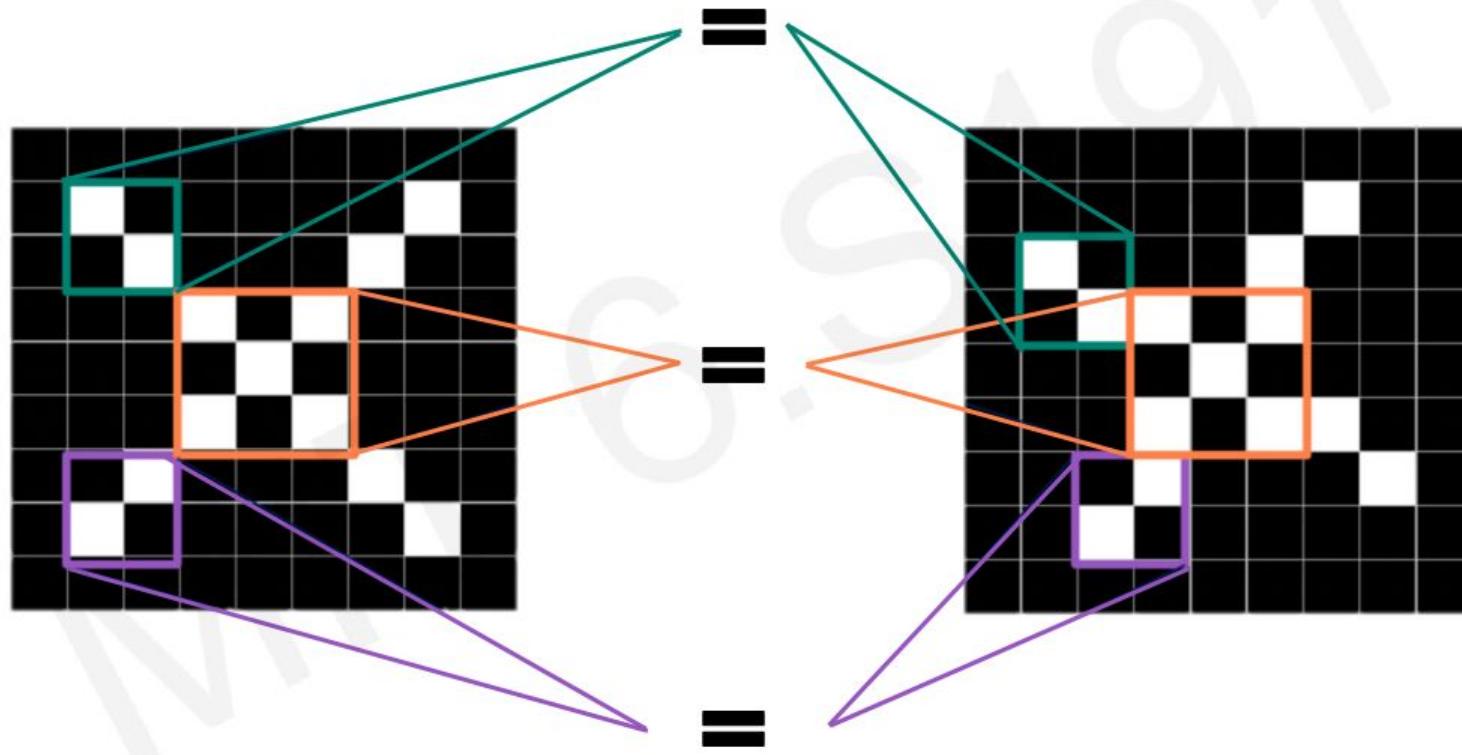


-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

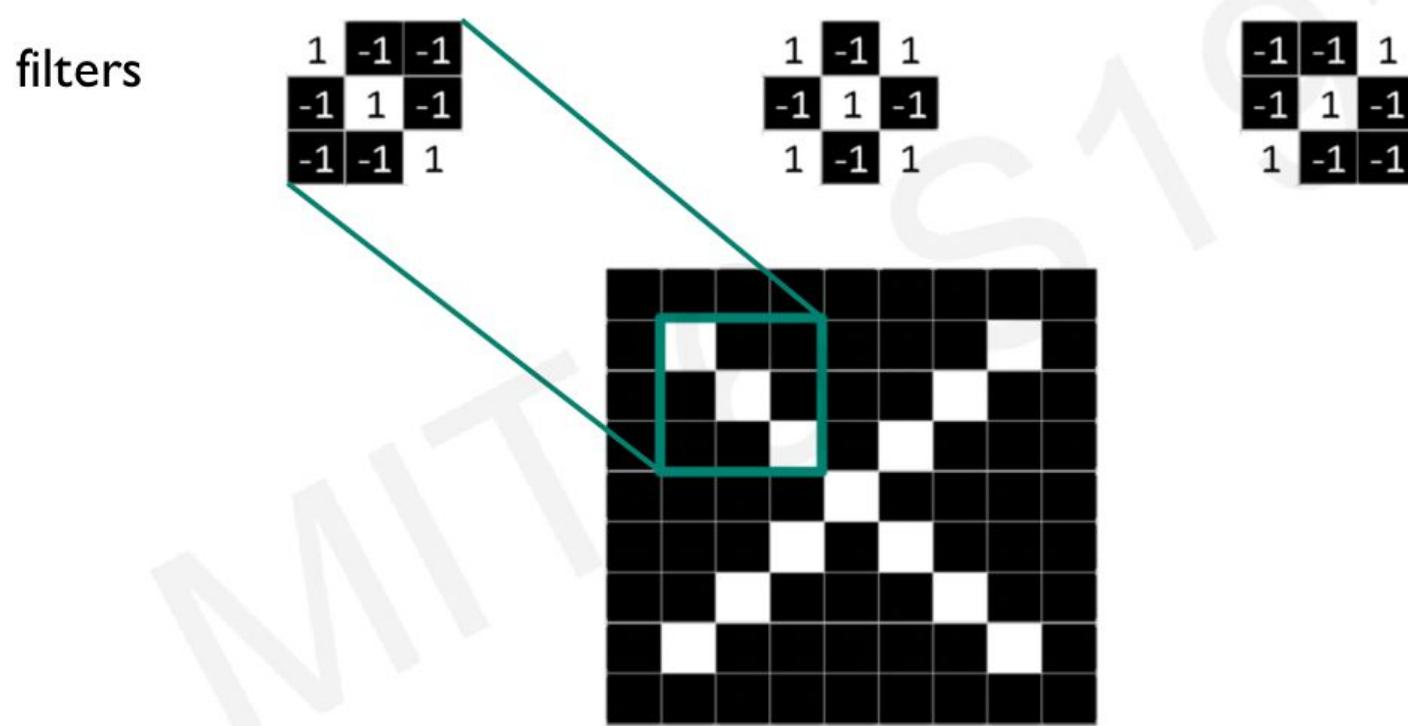
-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	1	-1
-1	1	-1	-1	-1	1	-1	-1	-1
-1	-1	1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	1	-1	-1
-1	-1	-1	1	-1	1	1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	1
-1	-1	1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

Image is represented as matrix of pixel values... and computers are literal!  
We want to be able to classify an X as an X even if it's shifted, shrunk, rotated, deformed.

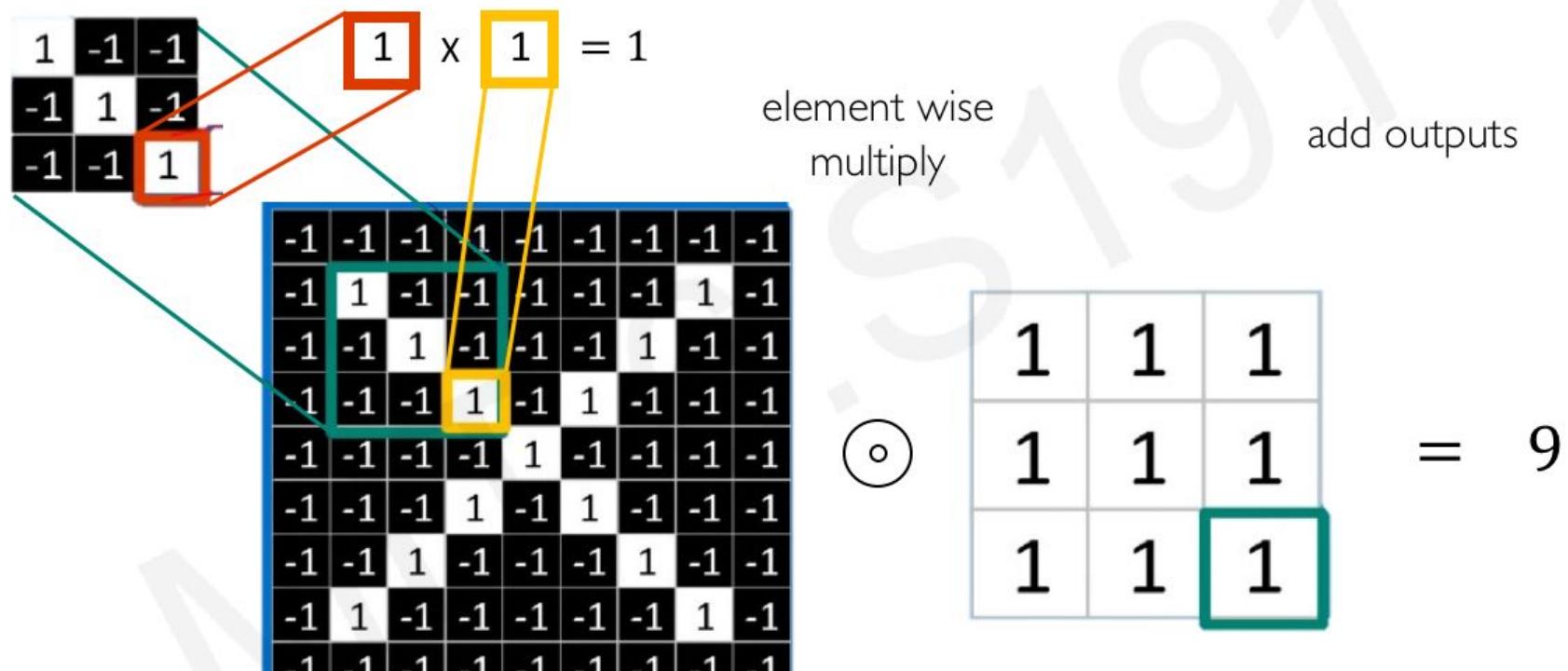
# Features of X



# Filters to Detect X Features

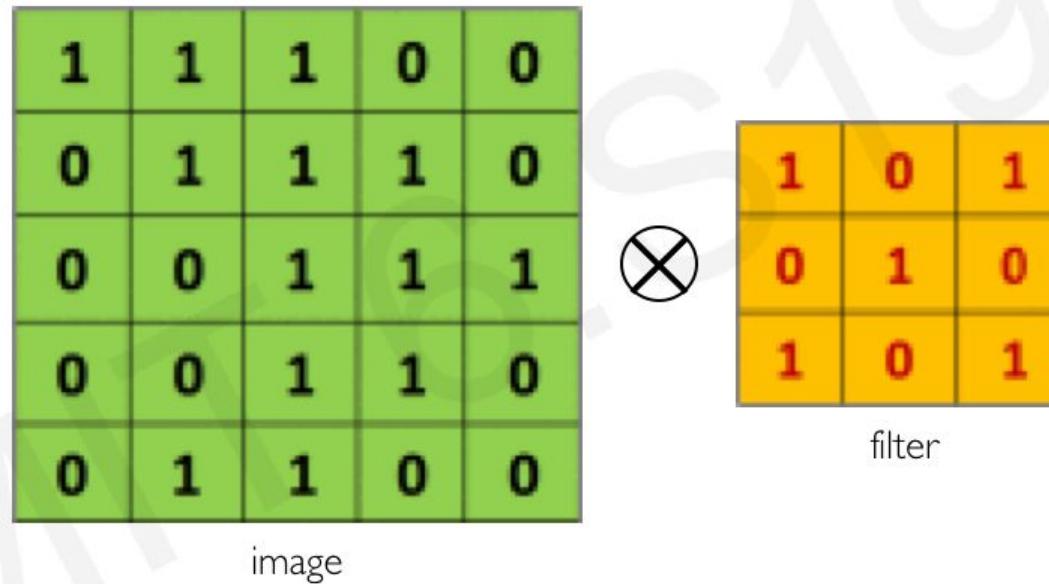


# The Convolution Operation



# The Convolution Operation

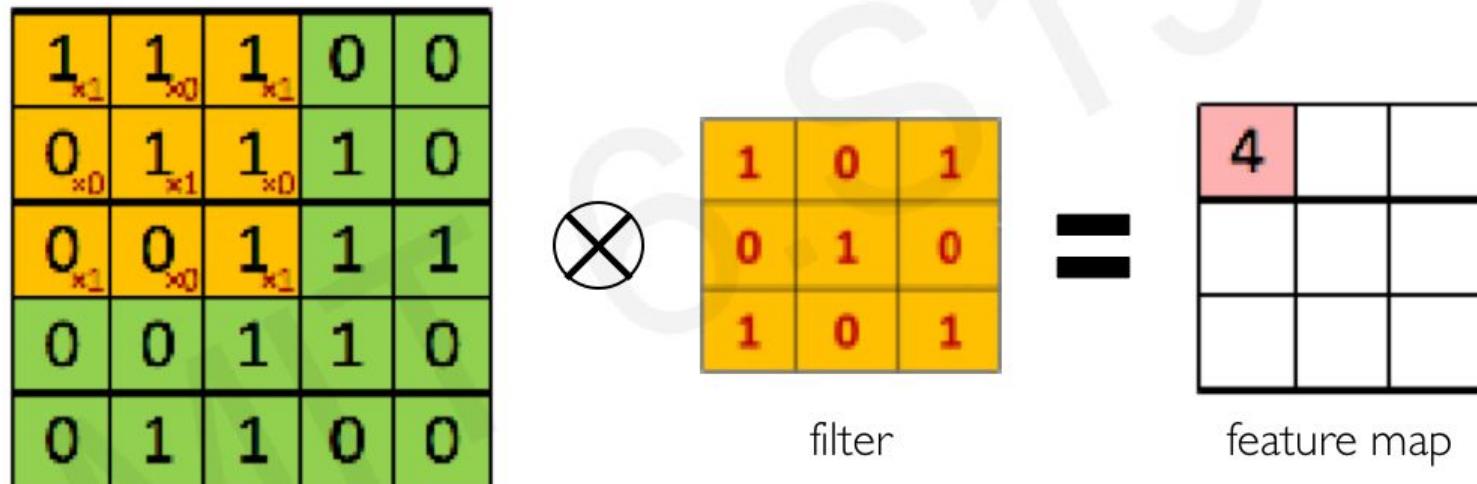
Suppose we want to compute the convolution of a 5x5 image and a 3x3 filter:



We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs...

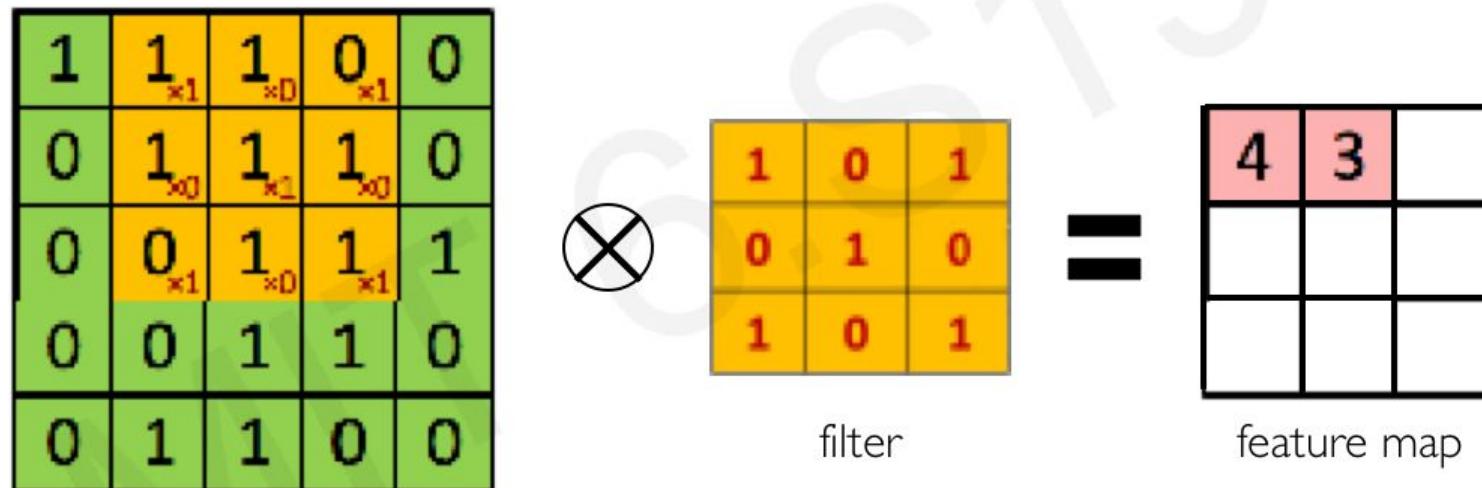
# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



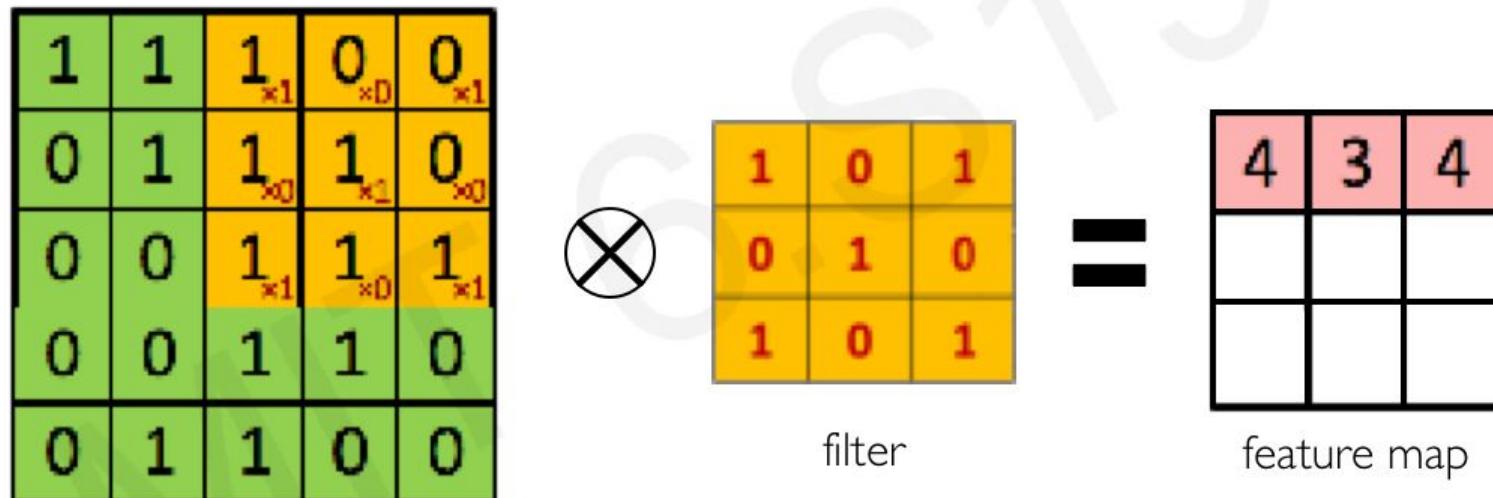
# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



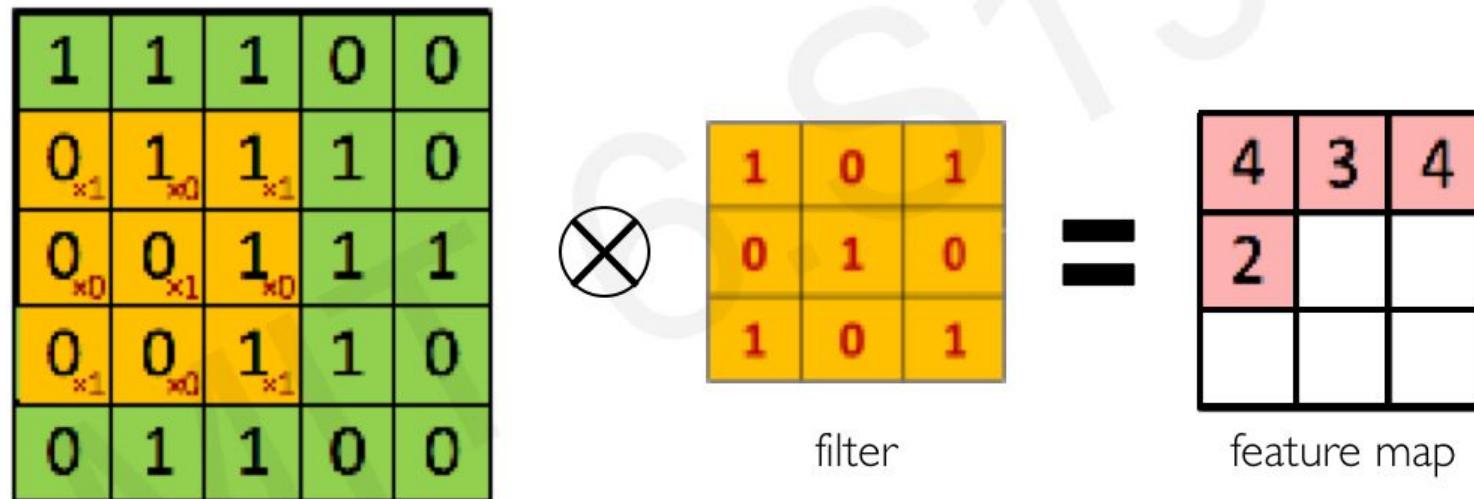
# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



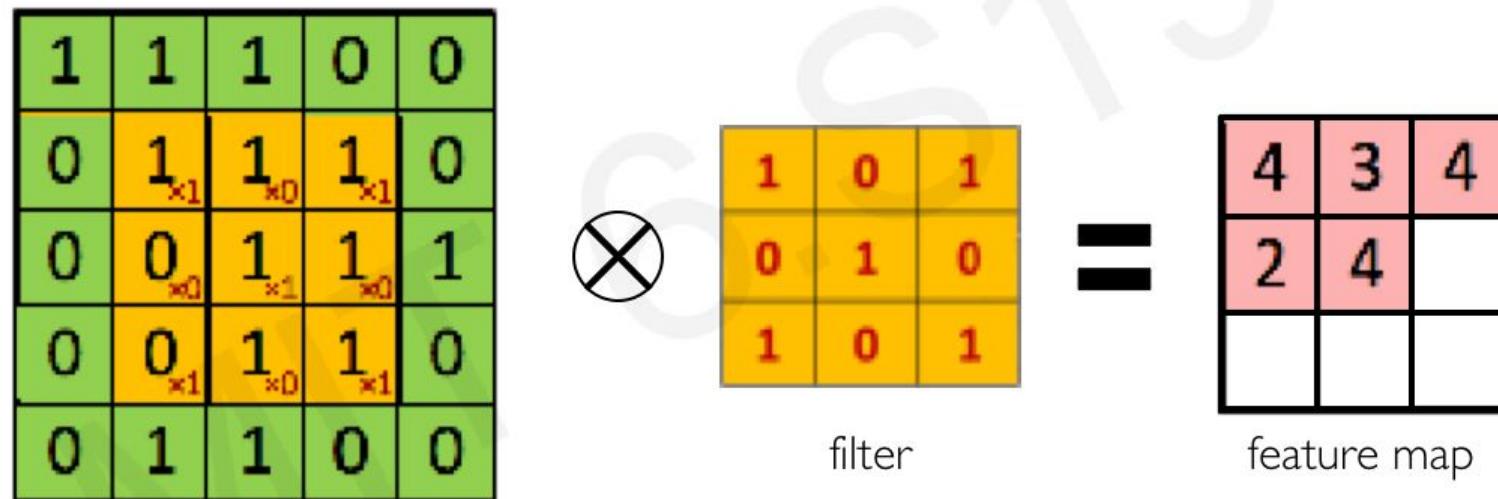
# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



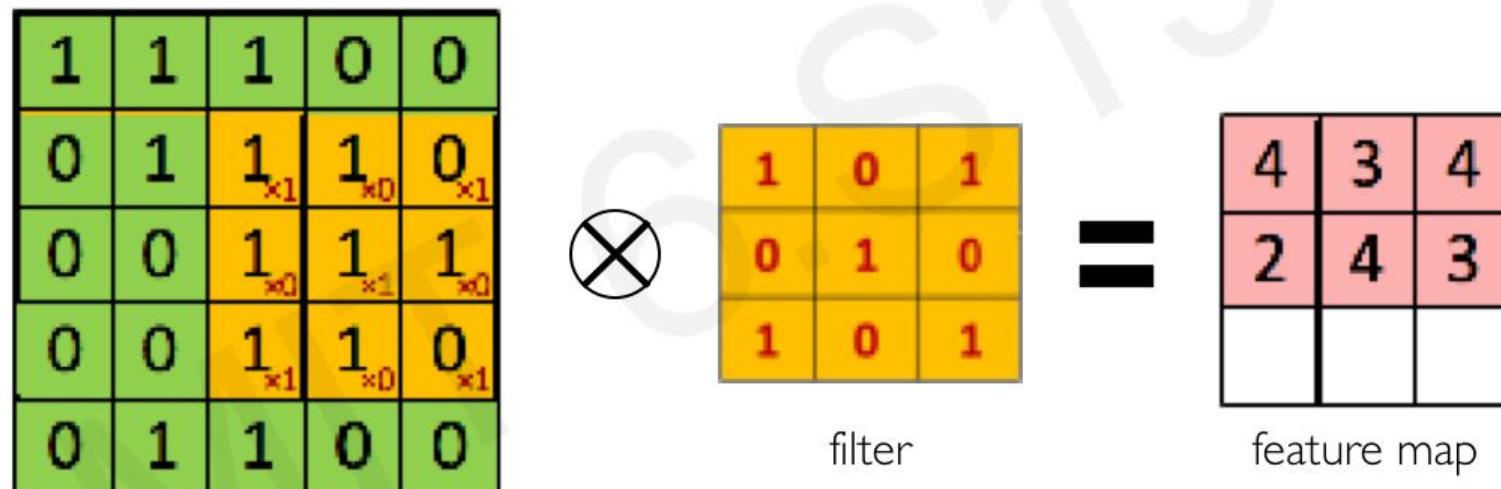
# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



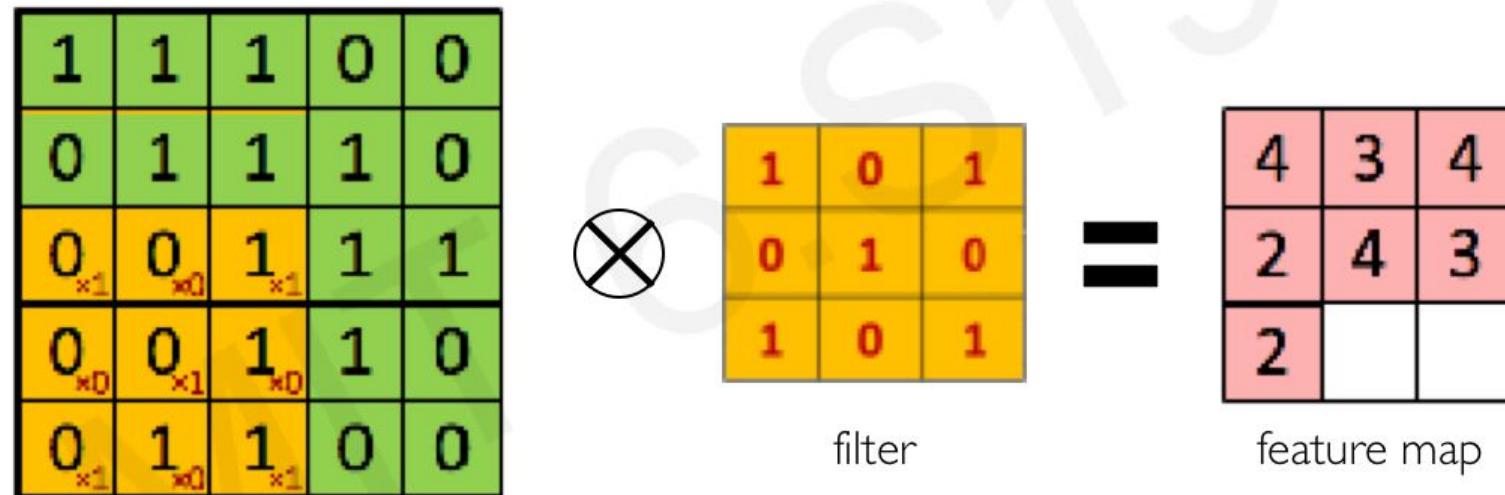
# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



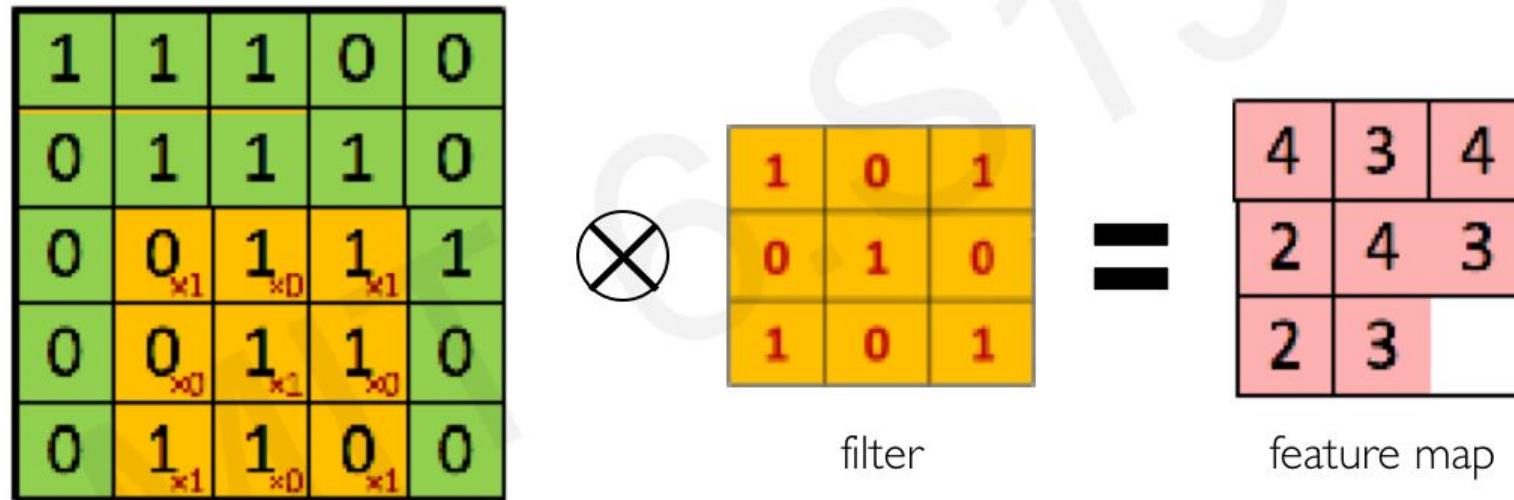
# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



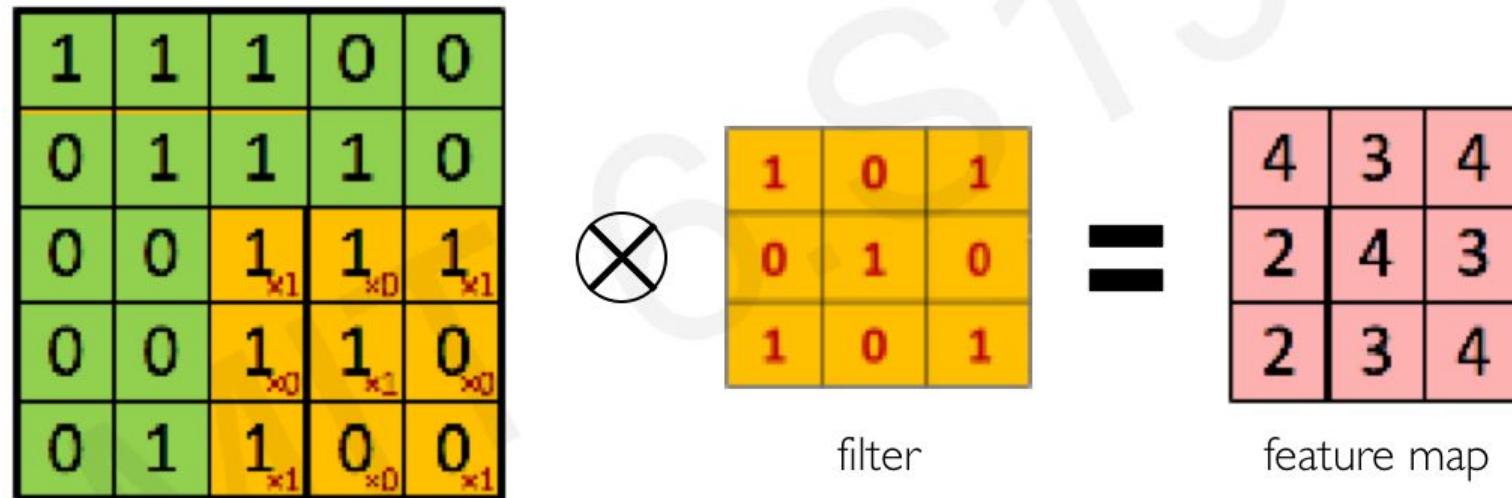
# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



# The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



# Producing Feature Maps



Original



Sharpen

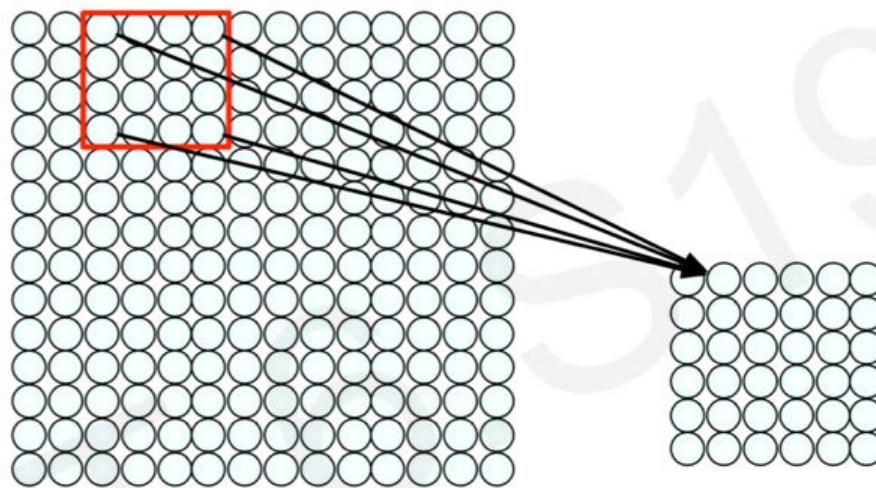


Edge Detect



“Strong” Edge  
Detect

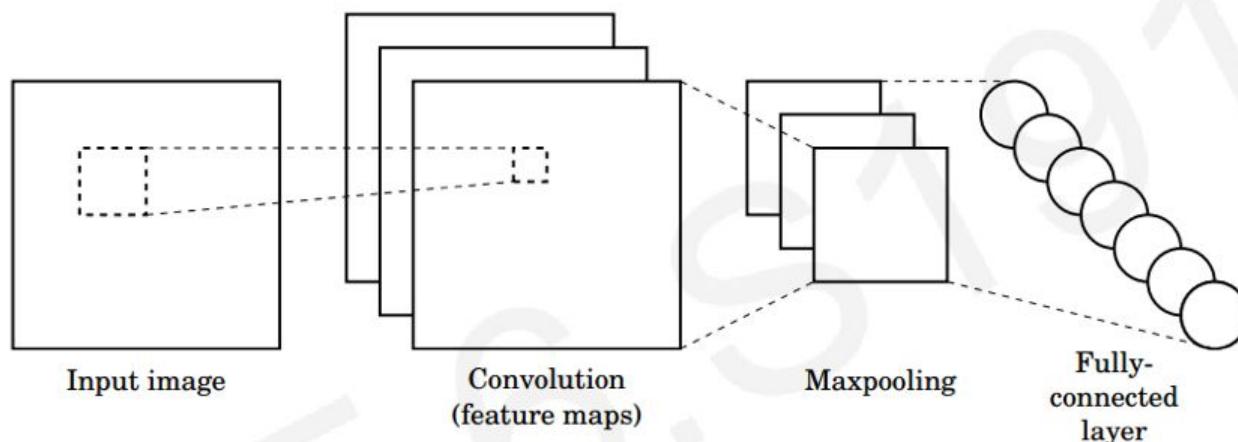
# Feature Extraction with Convolution



- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

# Convolutional Neural Networks (CNNs)

# CNNs for Classification



- 1. Convolution:** Apply filters to generate feature maps.
- 2. Non-linearity:** Often ReLU.
- 3. Pooling:** Downsampling operation on each feature map.

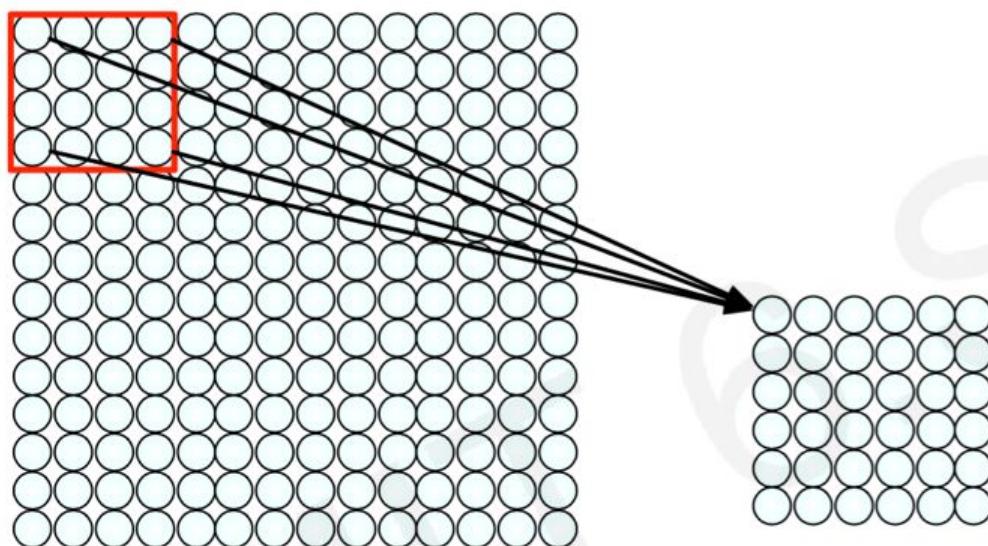
 `tf.keras.layers.Conv2D`

 `tf.keras.activations.*`

 `tf.keras.layers.MaxPool2D`

**Train model with image data.  
Learn weights of filters in convolutional layers.**

# Convolutional Layers: Local Connectivity

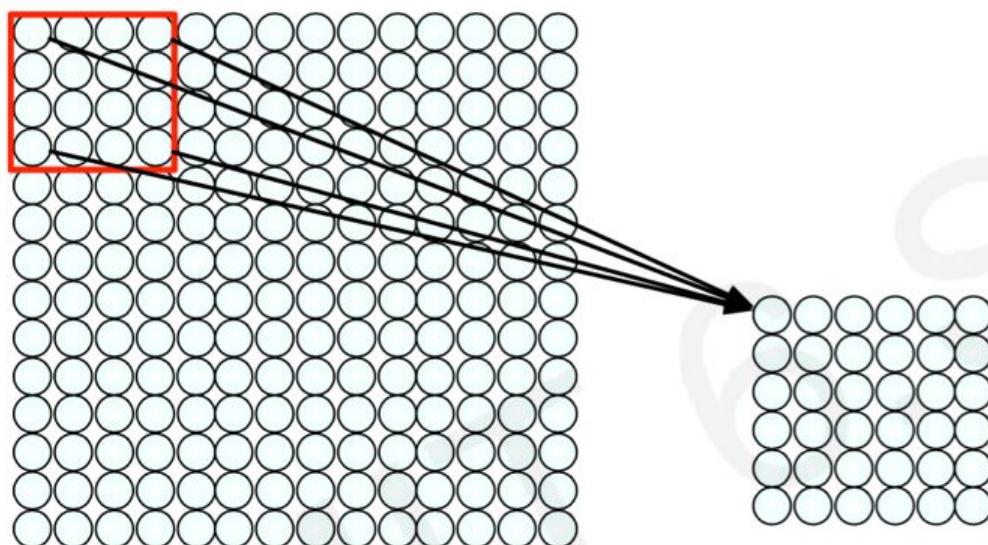


 `tf.keras.layers.Conv2D`

For a neuron in hidden layer:

- Take inputs from patch
- Compute weighted sum
- Apply bias

# Convolutional Layers: Local Connectivity



$$\sum_{i=1}^4 \sum_{j=1}^4 w_{ij} x_{i+p,j+q} + b$$

for neuron  $(p,q)$  in hidden layer

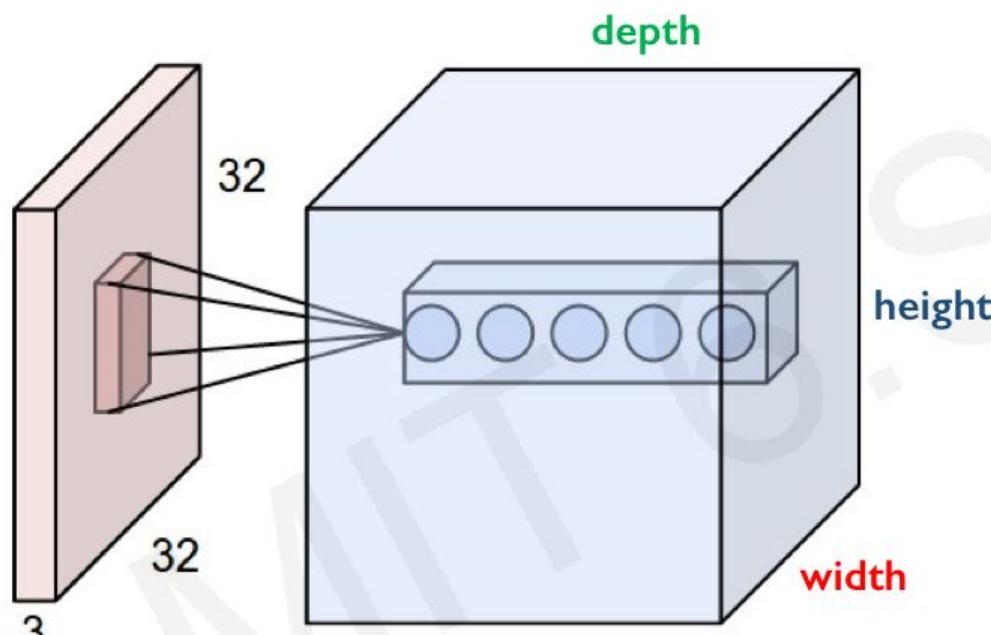
 `tf.keras.layers.Conv2D`

For a neuron in hidden layer:

- Take inputs from patch
- Compute weighted sum
- Apply bias

- 1) applying a window of weights
- 2) computing linear combinations
- 3) activating with non-linear function

# CNNs: Spatial Arrangement of Output Volume



**Layer Dimensions:**

$$h \times w \times d$$

where h and w are spatial dimensions  
d (depth) = number of filters

**Stride:**

Filter step size

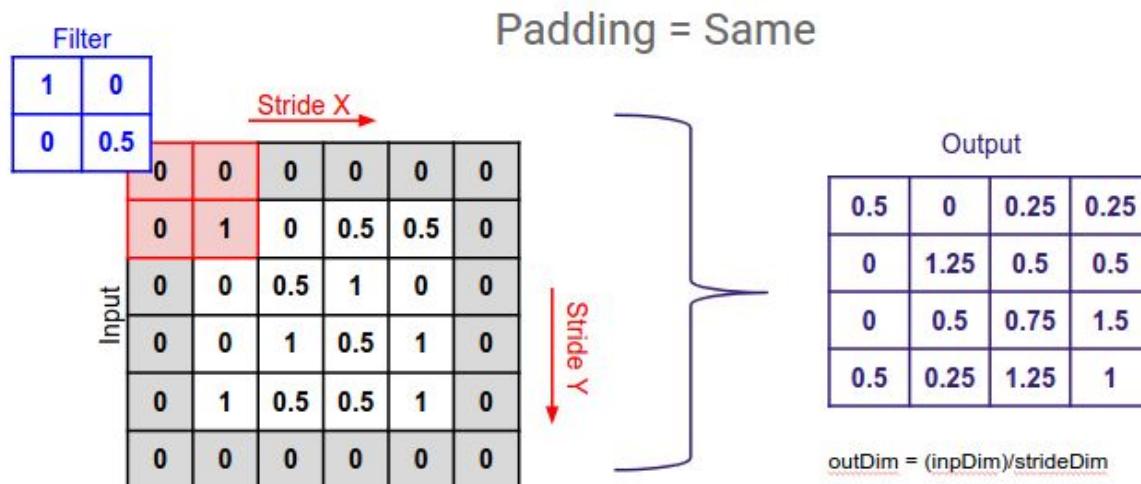
**Receptive Field:**

Locations in input image that  
a node is path connected to



`tf.keras.layers.Conv2D( filters=d, kernel_size=(h,w), strides=s )`

# Padding



$$n_{out} = \left\lfloor \frac{n_{in} + 2p - k}{s} \right\rfloor + 1$$

$n_{in}$ : number of input features

$n_{out}$ : number of output features

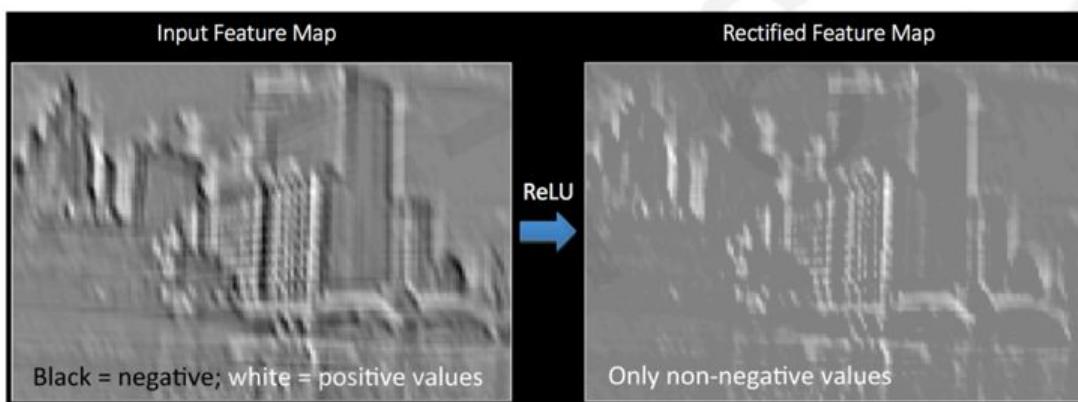
$k$ : convolution kernel size

$p$ : convolution padding size

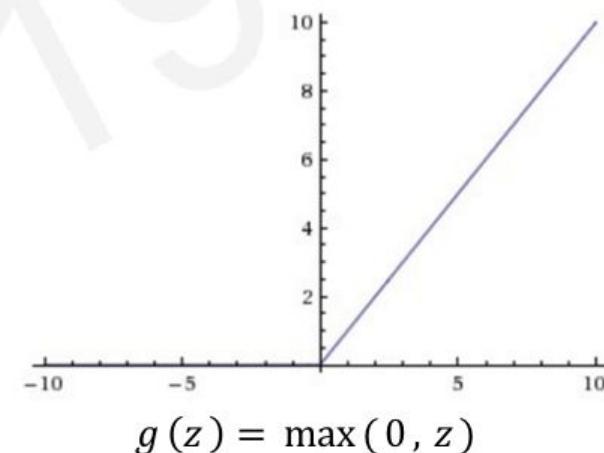
$s$ : convolution stride size

# Introducing Non-Linearity

- Apply after every convolution operation (i.e., after convolutional layers)
- ReLU: pixel-by-pixel operation that replaces all negative values by zero. **Non-linear operation!**

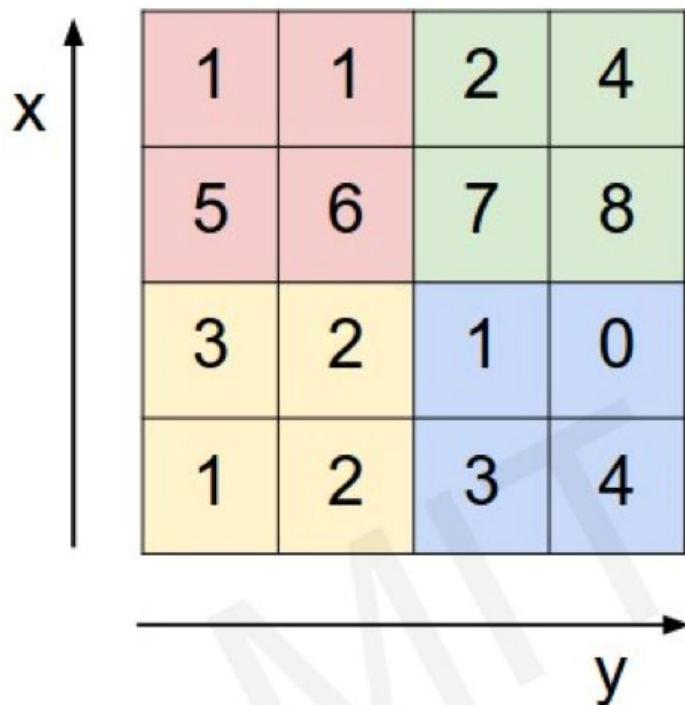


Rectified Linear Unit (ReLU)



 `tf.keras.layers.ReLU`

# Pooling



max pool with 2x2 filters  
and stride 2

```
tf.keras.layers.MaxPool2D(  
    pool_size=(2,2),  
    strides=2  
)
```



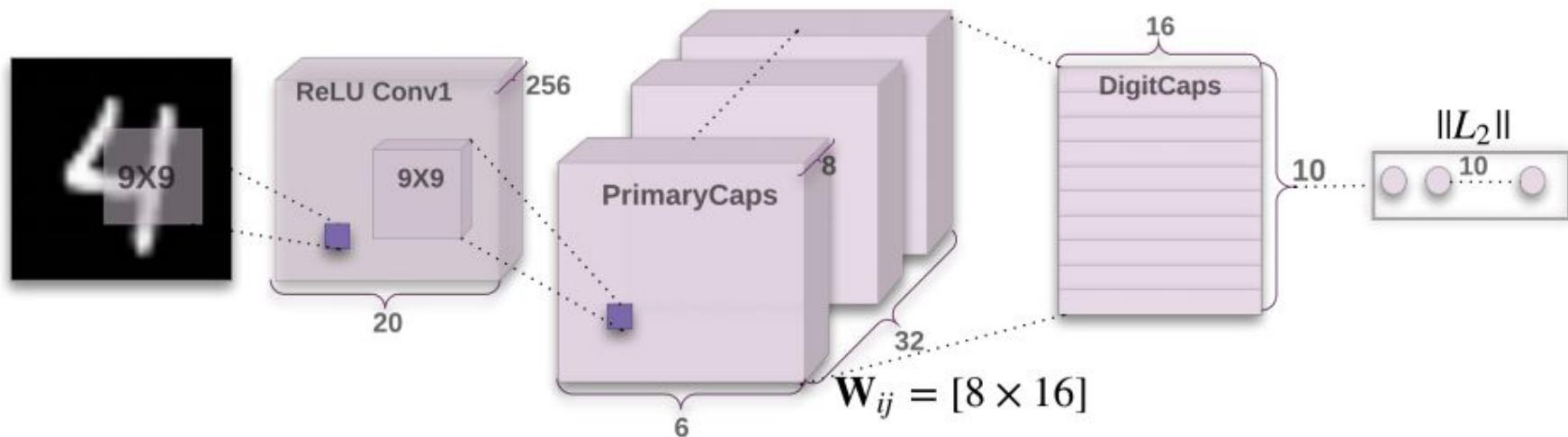
6	8
3	4

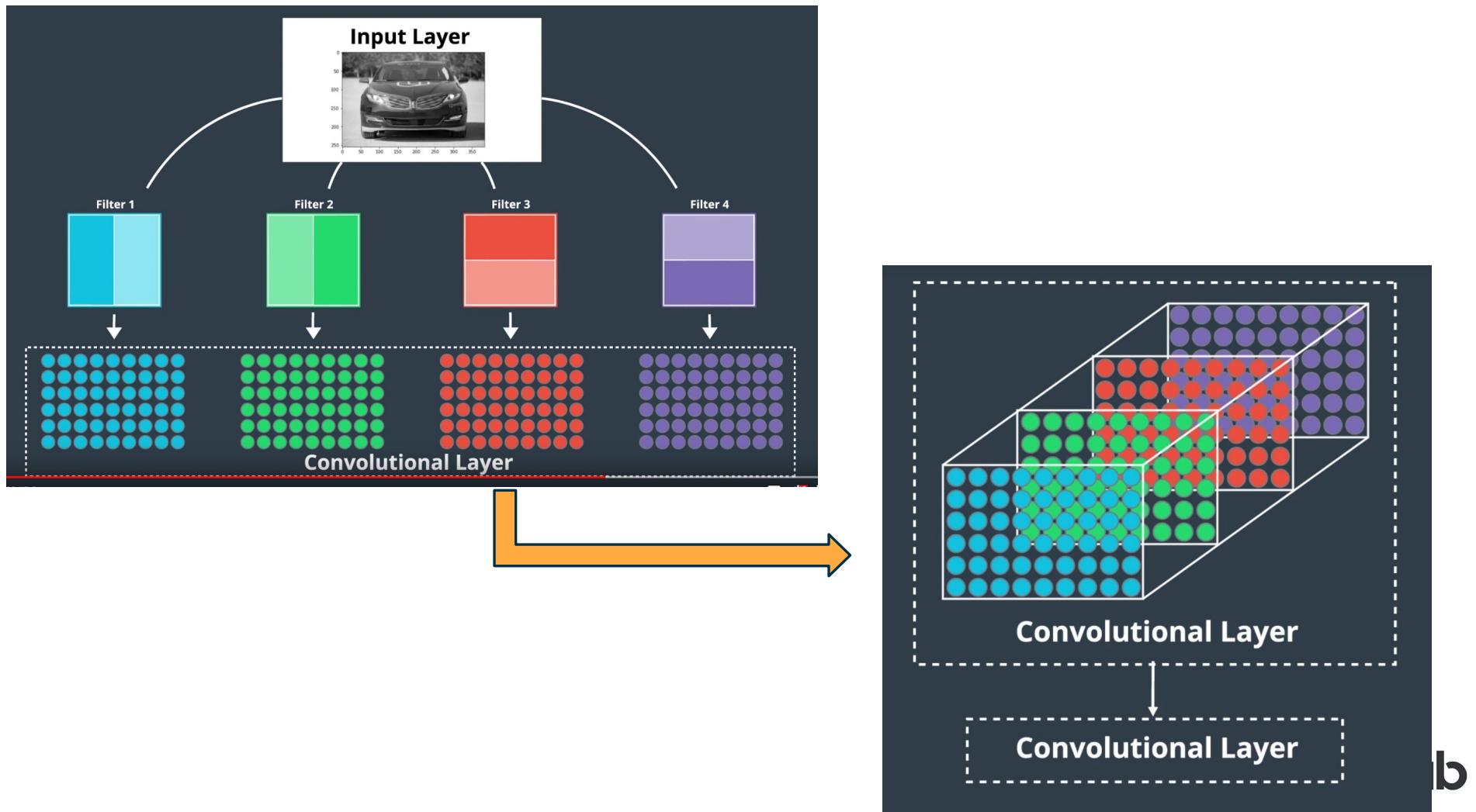
- 1) Reduced dimensionality
- 2) Spatial invariance

How else can we downsample and preserve spatial invariance?

Paper by  
**(Dynamic Routing Between Capsules)**  
<https://arxiv.org/abs/1710.09829>

# The Capsule Network





## Layer $l$ which is a convolution layer

$$\text{filter size} = f^{[l]} \quad \text{input} = n_h^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]}$$

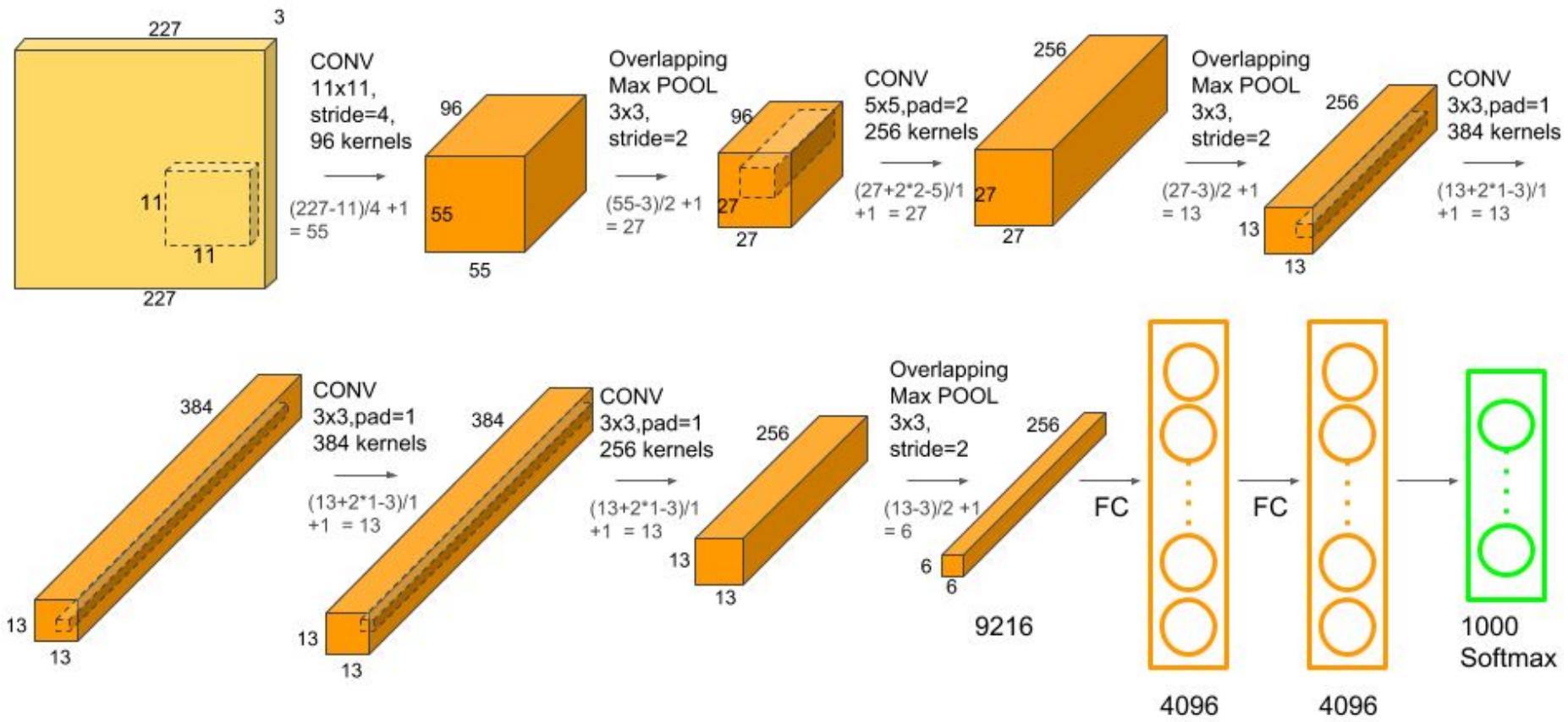
$$\text{padding} = p^{[l]} \quad \text{each filter} = f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$$

$$\text{stride} = s^{[l]} \quad \text{weights} = f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$$

$$\text{output} = n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$$

$$\text{bias} = 1 \times 1 \times 1 \times n_c^{[l]}$$



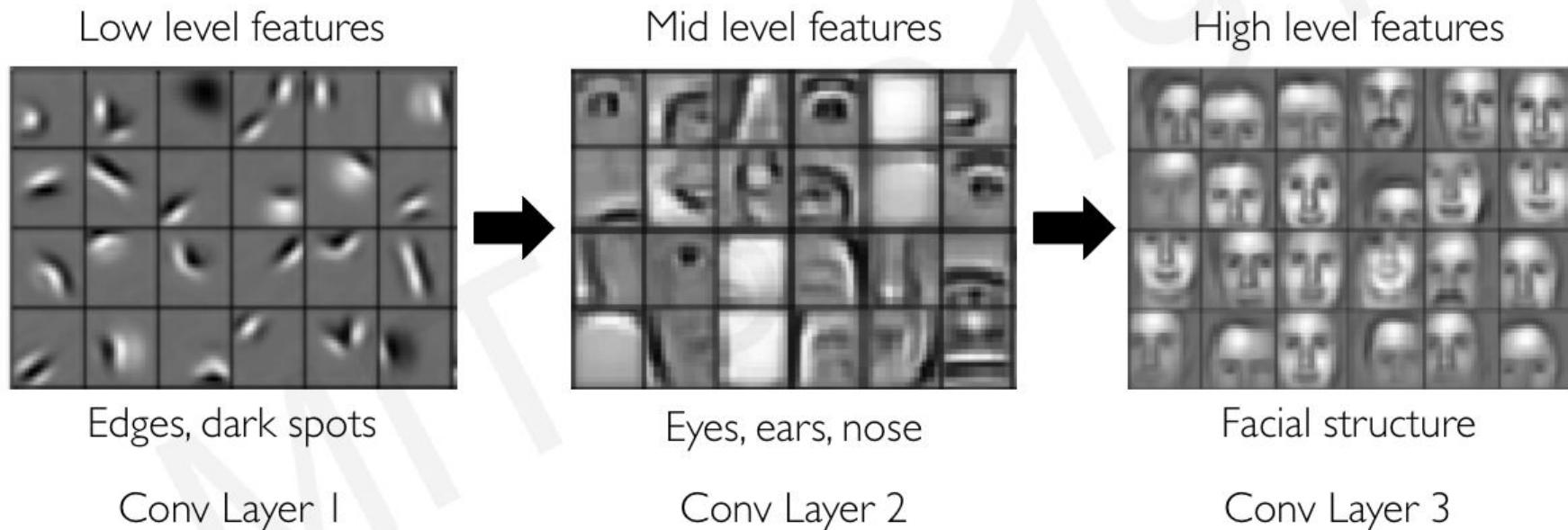


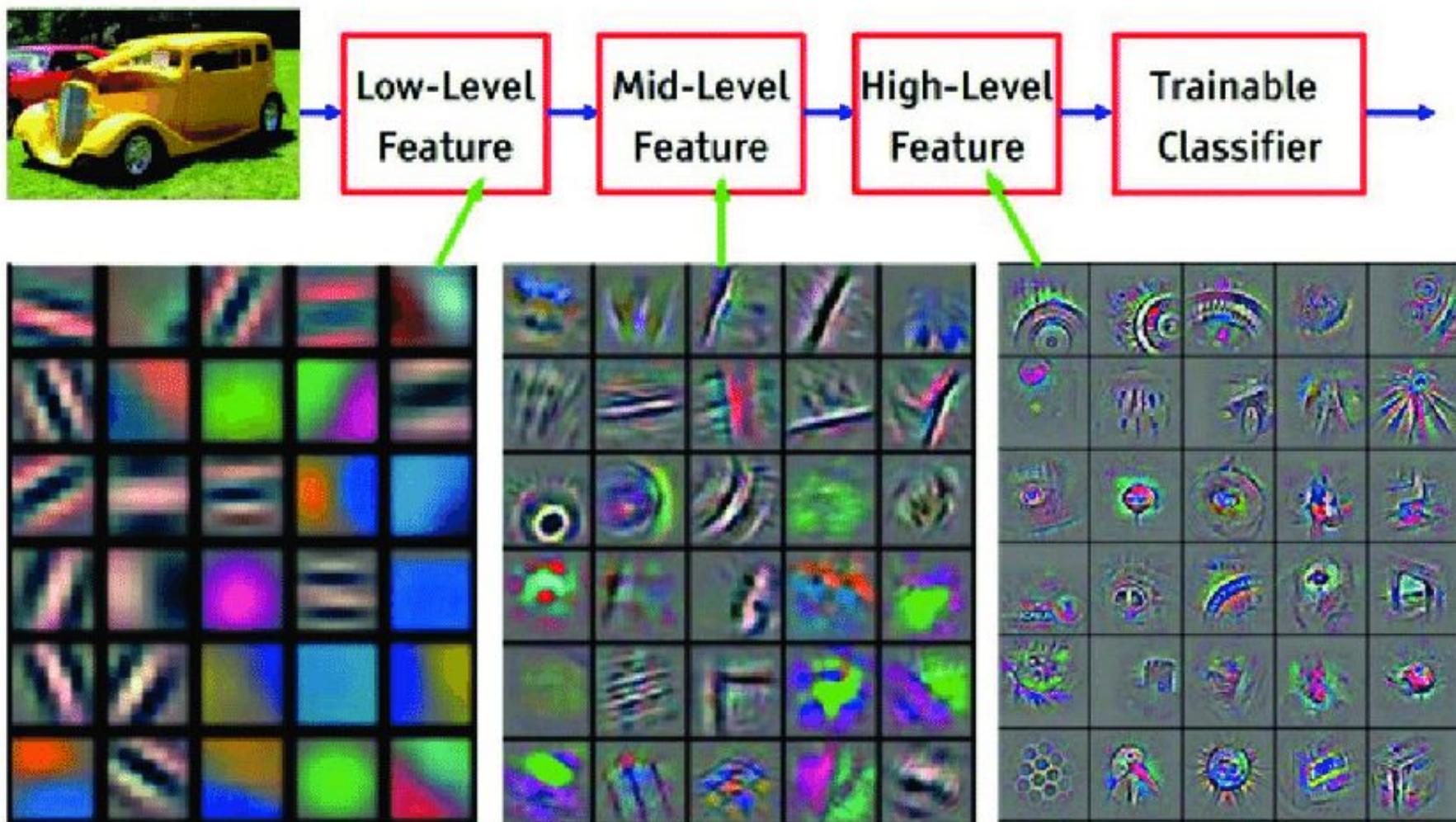
Height and Weight Decrease Where Channel

.

)

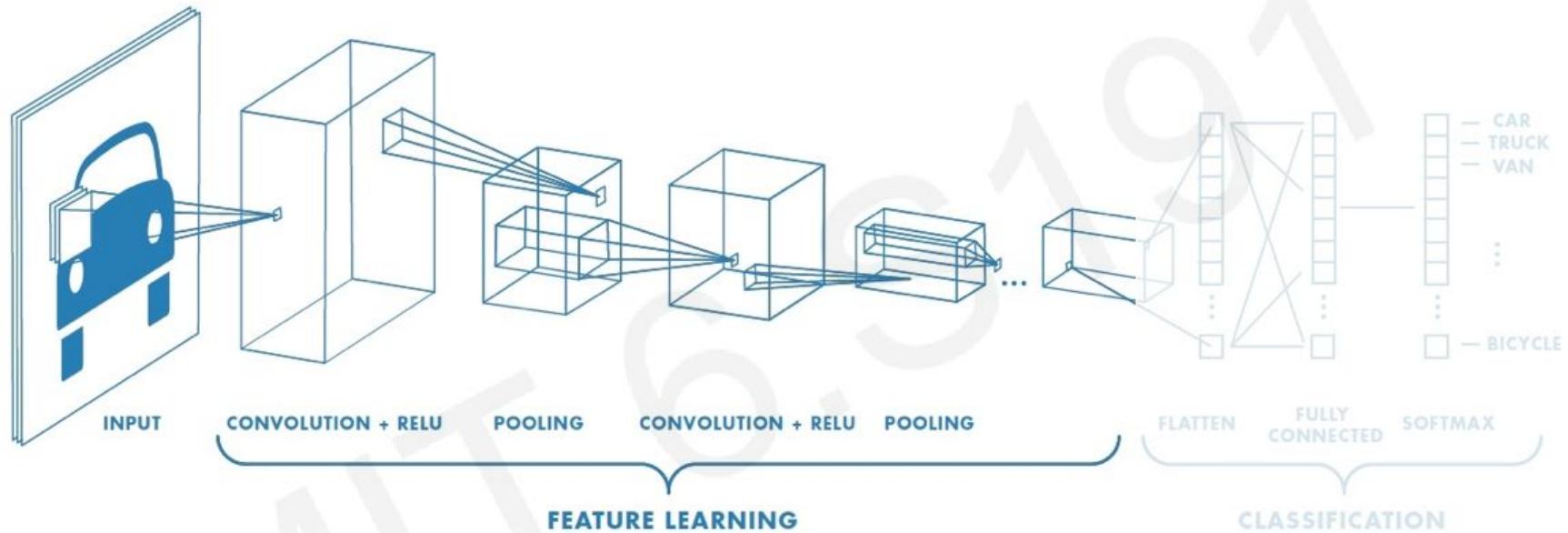
# Representation Learning in Deep CNNs





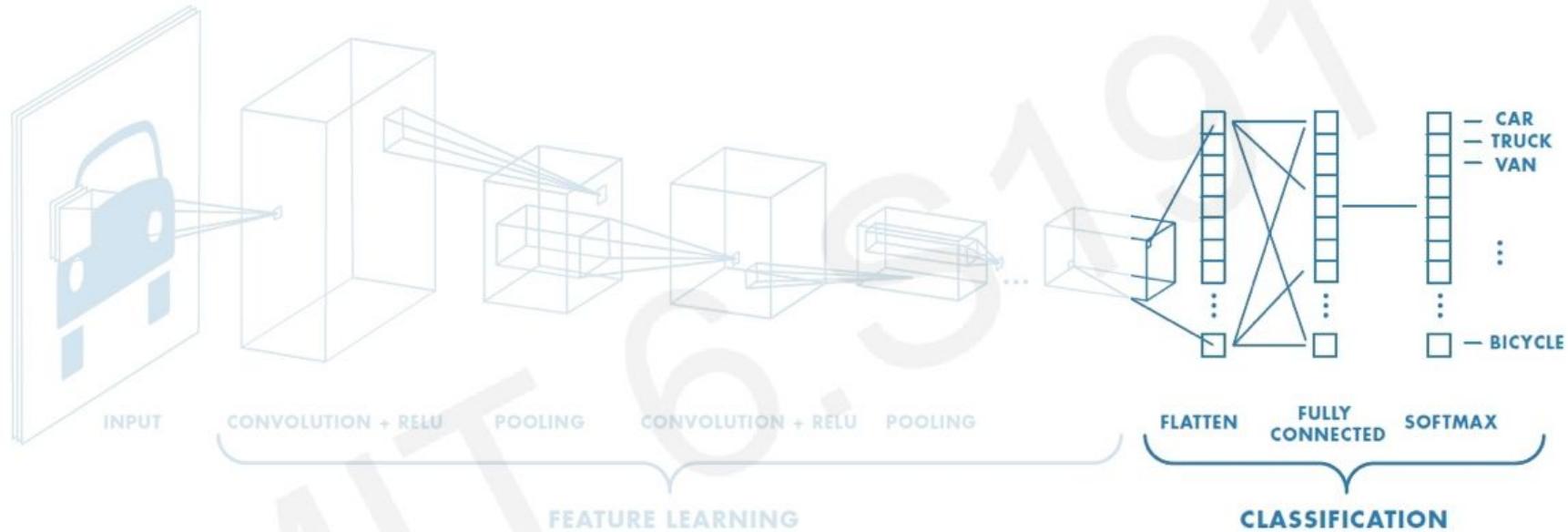
b

# CNNs for Classification: Feature Learning



1. Learn features in input image through **convolution**
2. Introduce **non-linearity** through activation function (real-world data is non-linear!)
3. Reduce dimensionality and preserve spatial invariance with **pooling**

# CNNs for Classification: Class Probabilities



- CONV and POOL layers output high-level features of input
- Fully connected layer uses these features for classifying input image
- Express output as **probability** of image belonging to a particular class

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

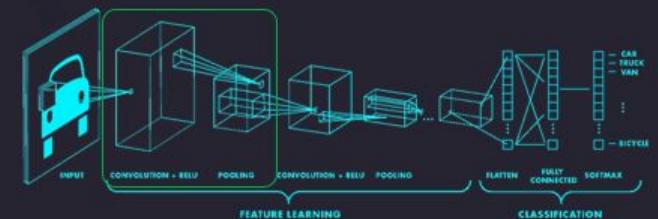
# Putting it all together

```
import tensorflow as tf

def generate_model():
    model = tf.keras.Sequential([
        # first convolutional layer
        tf.keras.layers.Conv2D(32, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),  
  

        # second convolutional layer
        tf.keras.layers.Conv2D(64, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),  
  

        # fully connected classifier
        tf.keras.layers.Flatten(),
        tf.keras.layers.Dense(1024, activation='relu'),
        tf.keras.layers.Dense(10, activation='softmax') # 10 outputs
    ])
    return model
```

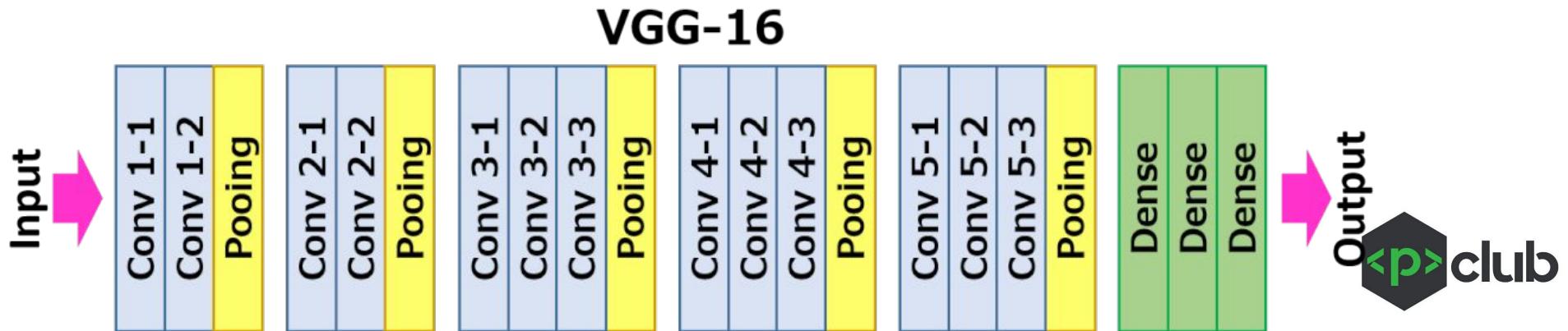
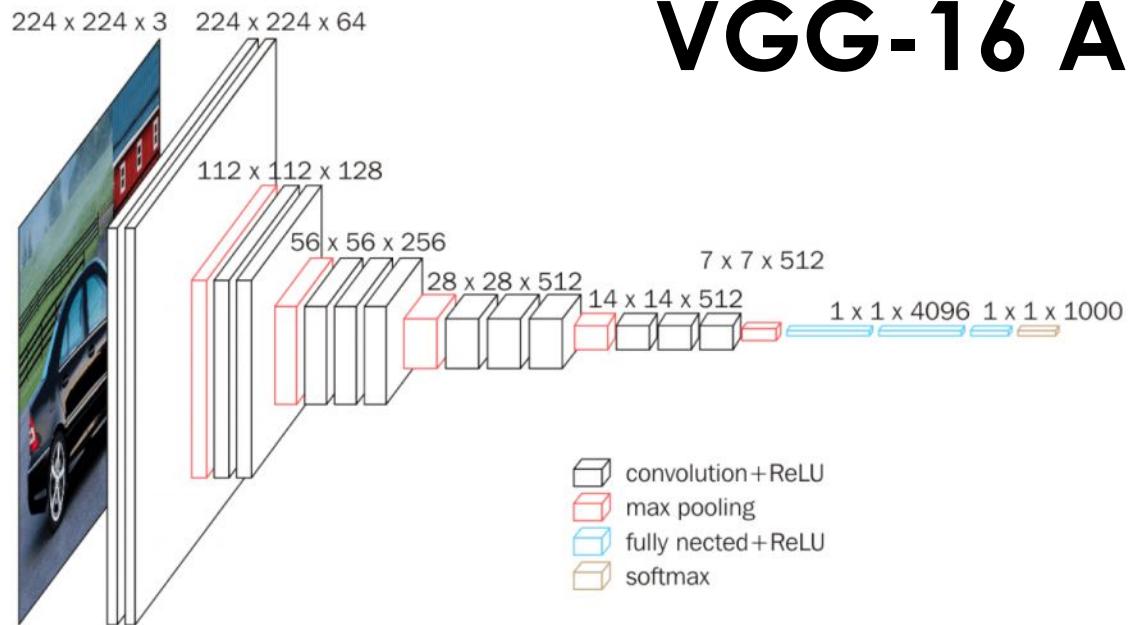


# Classical Networks

- LeNet-5
- Alex Net
- VGG-16/19



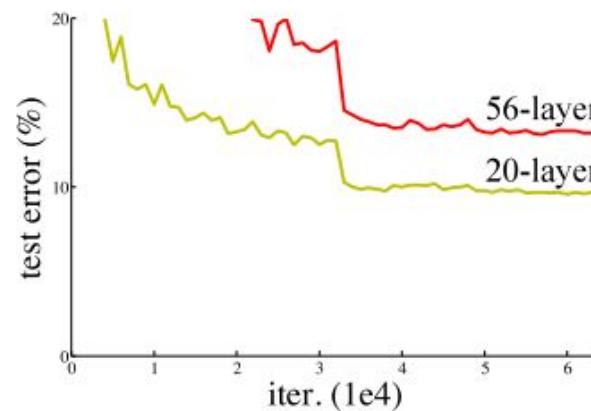
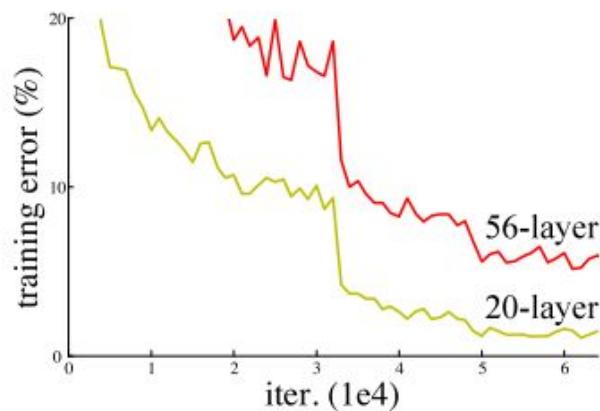
# VGG-16 Architecture



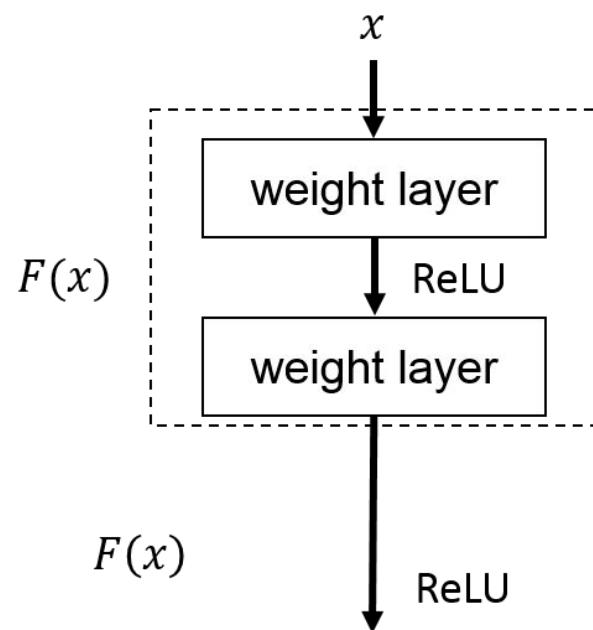
**Researchers observed that it makes sense to affirm that “*the deeper the better*” when it comes to convolutional neural networks. This makes sense, since the models should be more capable ( their flexibility to adapt to any space increase because they have a bigger parameter space to explore ).**

**However, it has been noticed that after some depth, the performance degrades.**

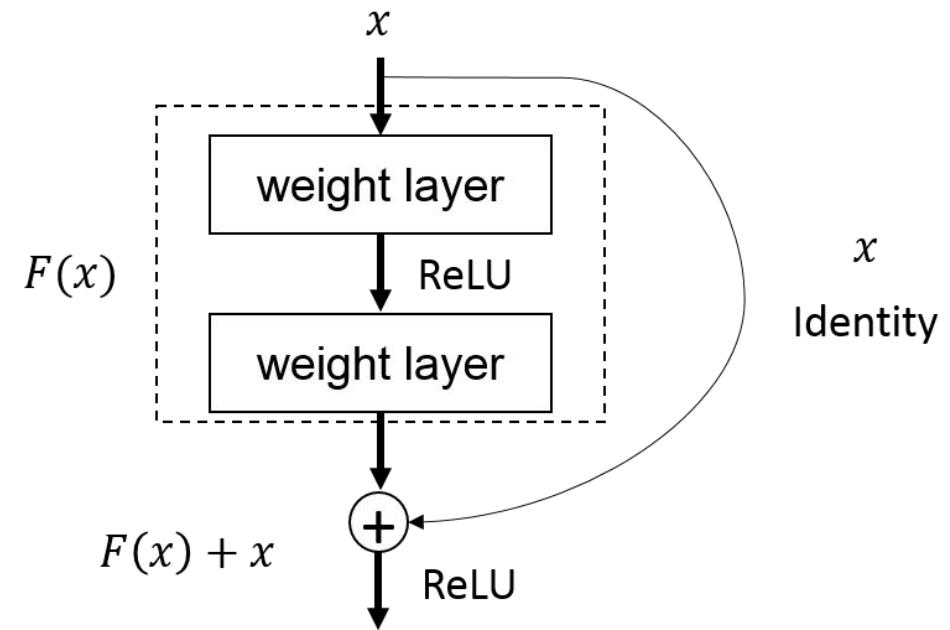
**This was one of the bottlenecks of VGG. They couldn’t go as deep as wanted, because they started to lose generalization capability.**

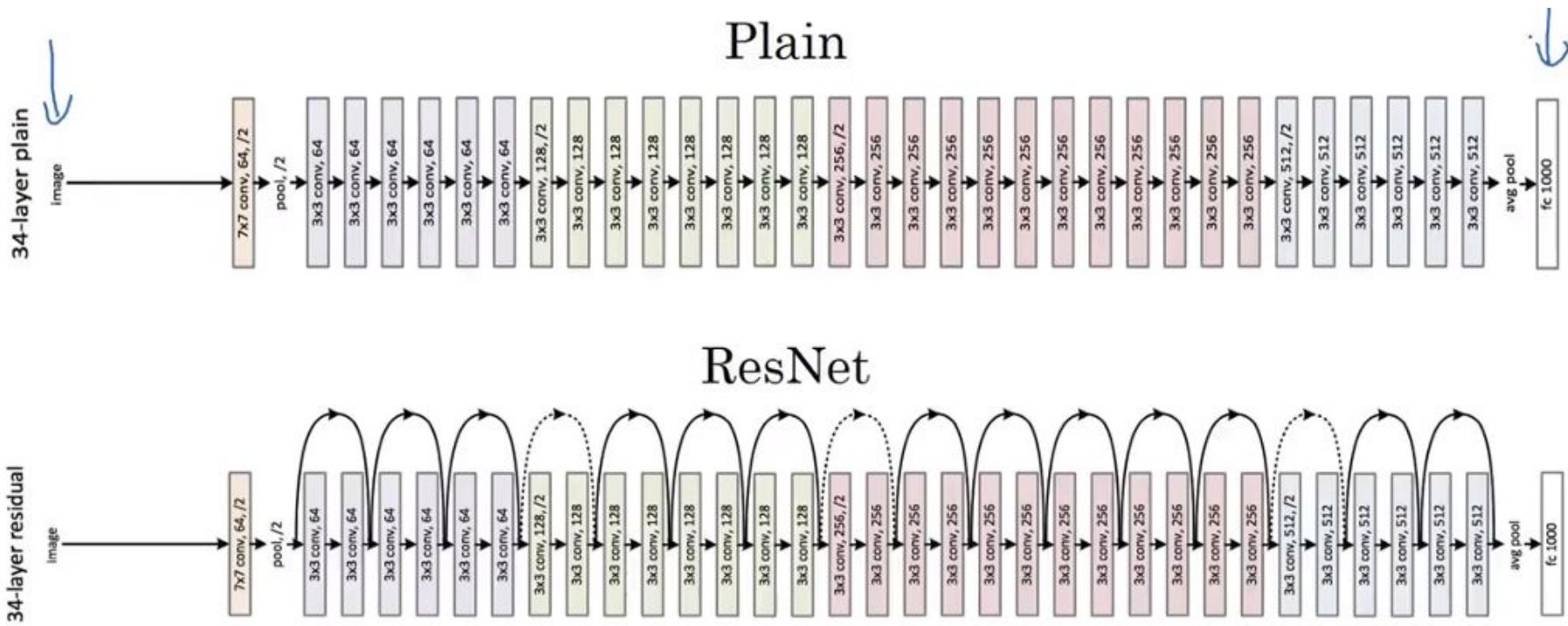


Plain network



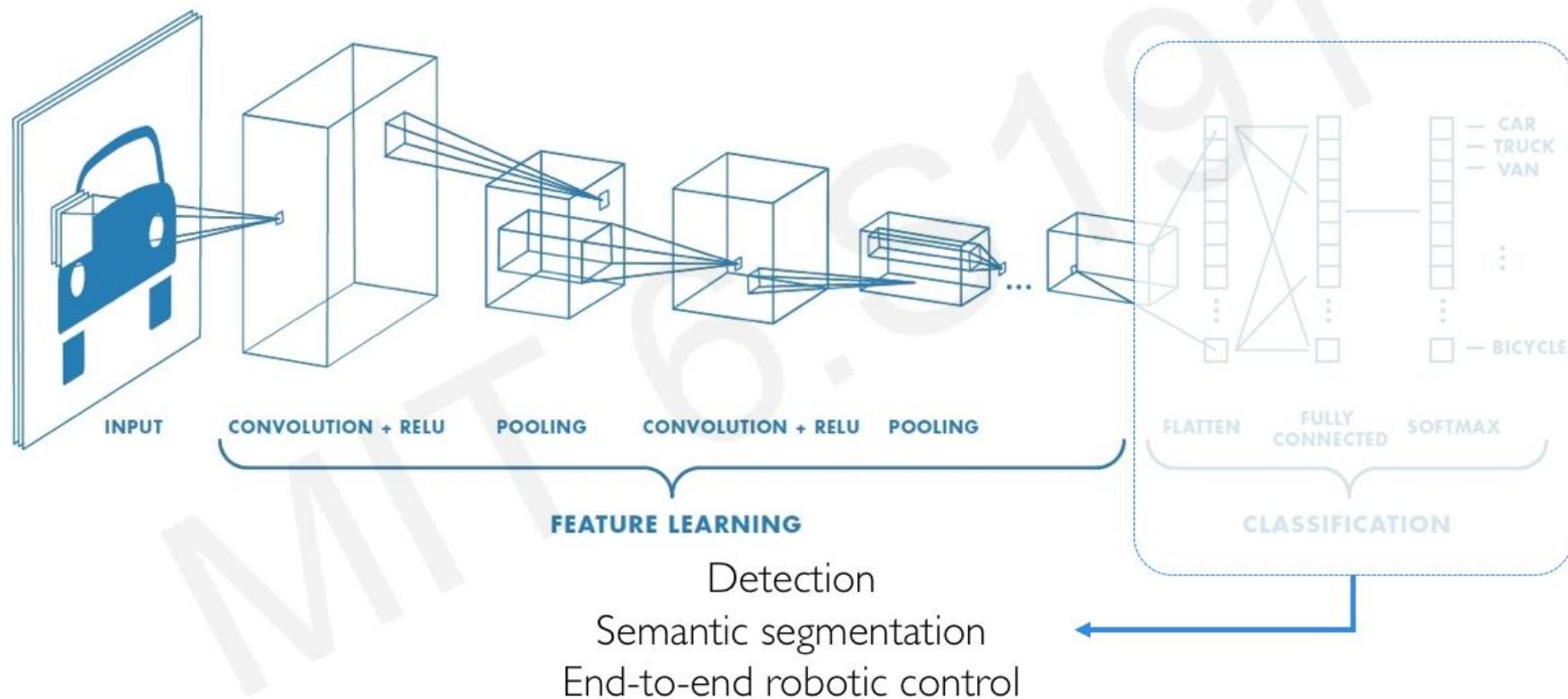
Residual network



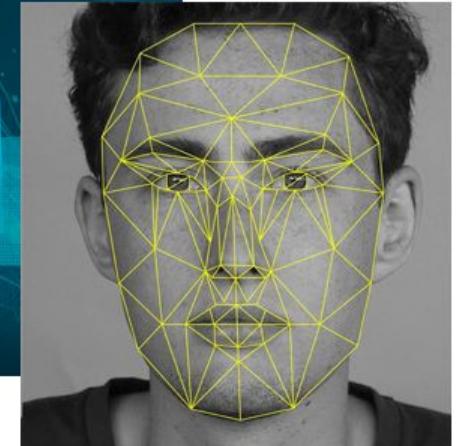
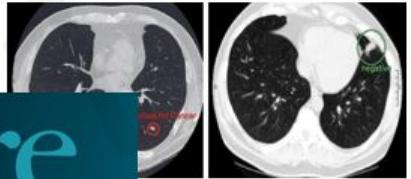


# An Architecture for Many Applications

# An Architecture for Many Applications



# Deep Learning for Computer Vision: Impact



# Deep Learning for Computer Vision: Summary

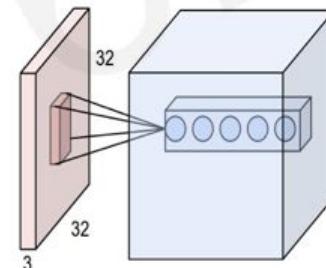
## Foundations

- Why computer vision?
- Representing images
- Convolutions for feature extraction



## CNNs

- CNN architecture
- Application to classification
- ImageNet



## Applications

- Segmentation, image captioning, control
- Security, medicine, robotics





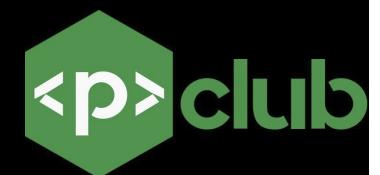
Week 5:

# Optimizing NN and Paper Presentation

## PCLUB CS2020

Ankur Bhatia

Feb 22, 2020



# Today's Agenda:

1. Normalization and Batch Normalization
2. PAPER/ IDEA PRESENTATION



# Normal Distribution

$$\sigma = \sqrt{\frac{\sum(X - \mu)^2}{n}}$$

where,

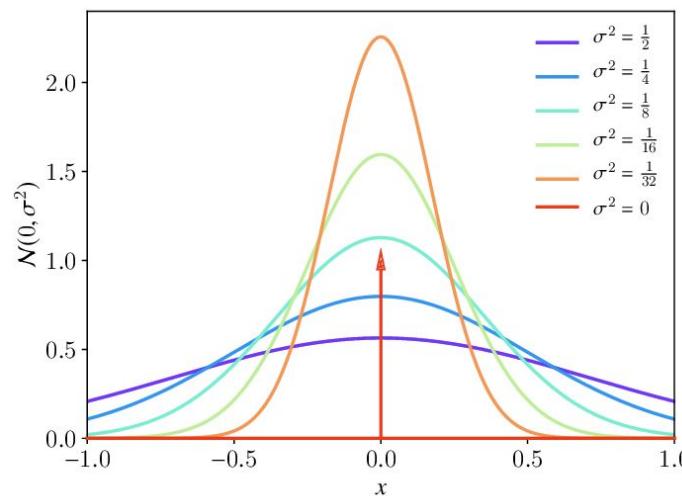
$\sigma$  = population standard deviation

$\sum$  = sum of...

$\mu$  = population mean

$n$  = number of scores in sample.

**Normal distribution**, also known as the **Gaussian distribution**, is a probability **distribution** that is symmetric about the mean, showing that data near the mean are more frequent in occurrence than data far from the mean.



## Skewness and Kurtosis

<https://www.investopedia.com/terms/n/normaldistribution.asp>



# Uniform Distribution

## Uniform Distribution

 DOWNLOAD  
Wolfram Notebook

 EXPLORE THIS TOPIC IN  
The MathWorld Classroom

A uniform distribution, sometimes also known as a rectangular distribution, is a distribution that has constant probability.



The probability density function and cumulative distribution function for a continuous uniform distribution on the interval  $[a, b]$  are

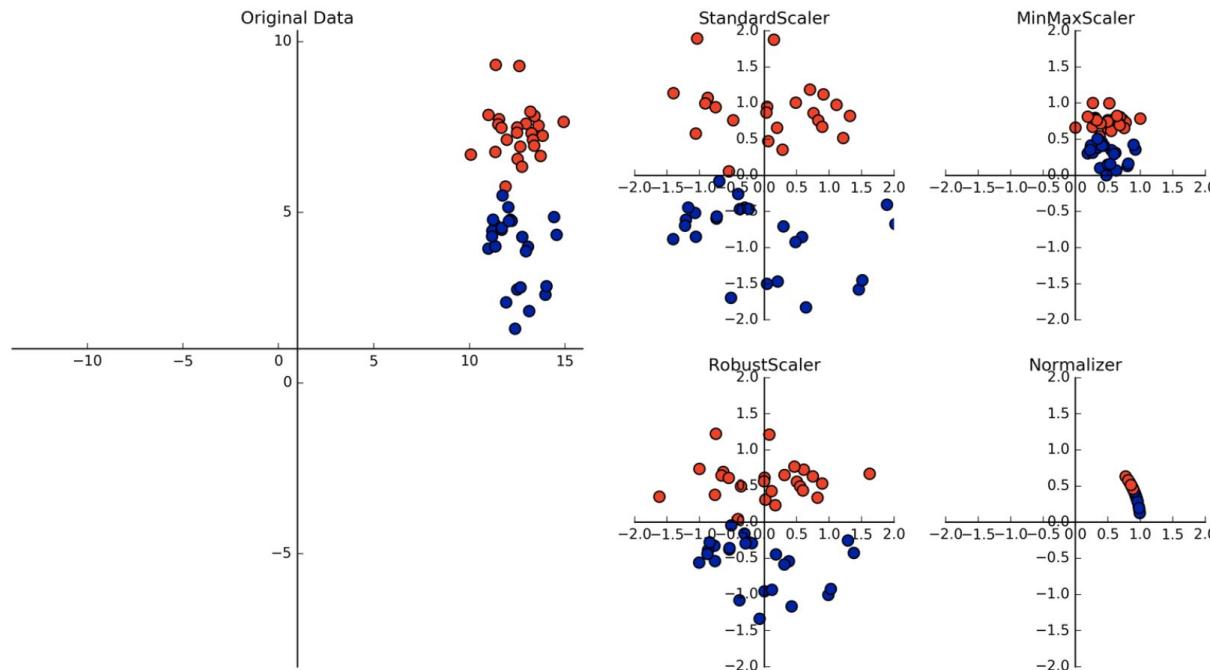
$$P(x) = \begin{cases} 0 & \text{for } x < a \\ \frac{1}{b-a} & \text{for } a \leq x \leq b \\ 0 & \text{for } x > b \end{cases}$$

$$D(x) = \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } a \leq x \leq b \\ 1 & \text{for } x > b. \end{cases}$$

<http://mathworld.wolfram.com/UniformDistribution.html>



# Normalizing Training Set

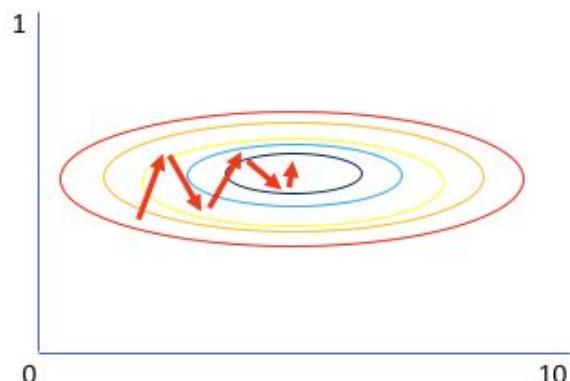


$$\begin{aligned}\mu_{\mathcal{B}} &\leftarrow \frac{1}{m} \sum_{i=1}^m x_i \\ \sigma_{\mathcal{B}}^2 &\leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \\ \hat{x}_i &\leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}\end{aligned}$$

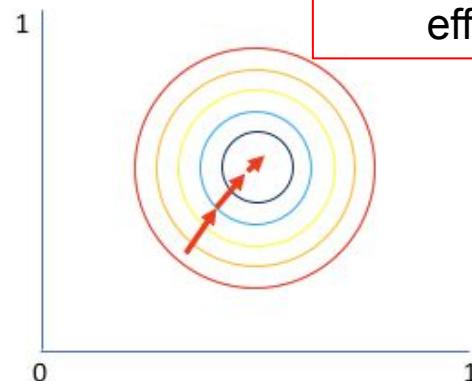
1. Use the same mean sd variance for train and test data.



## Why normalize?



Gradient of larger parameter  
dominates the update



Both parameters can be  
updated in equal proportions

- 1. Cos func is round and easier to optimize.
- 2. Features on similar scale(similar variances)
- 3. Normalizing  $x_1, x_2, x_3 \dots$  helps you train  $w$  and  $b$  more efficiently.

If not normalized, then the parameter ratios can take very different values.  
Cost func. will be an elongated bowl like.

- 1. Easier and faster to optimize.



# Random Weight Initialization

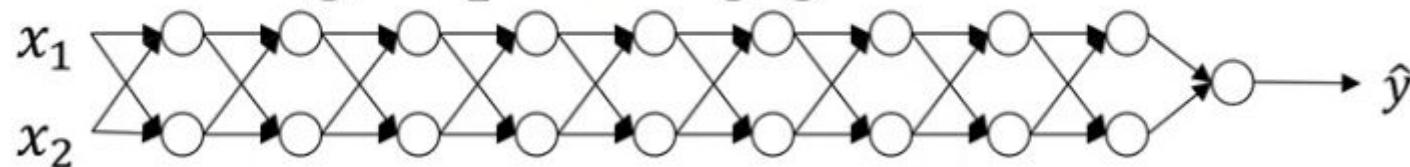
1. Vanishing and Exploding Gradients.
2. How to reduce the problem of V and E gradients?



# Random Weight Initialization

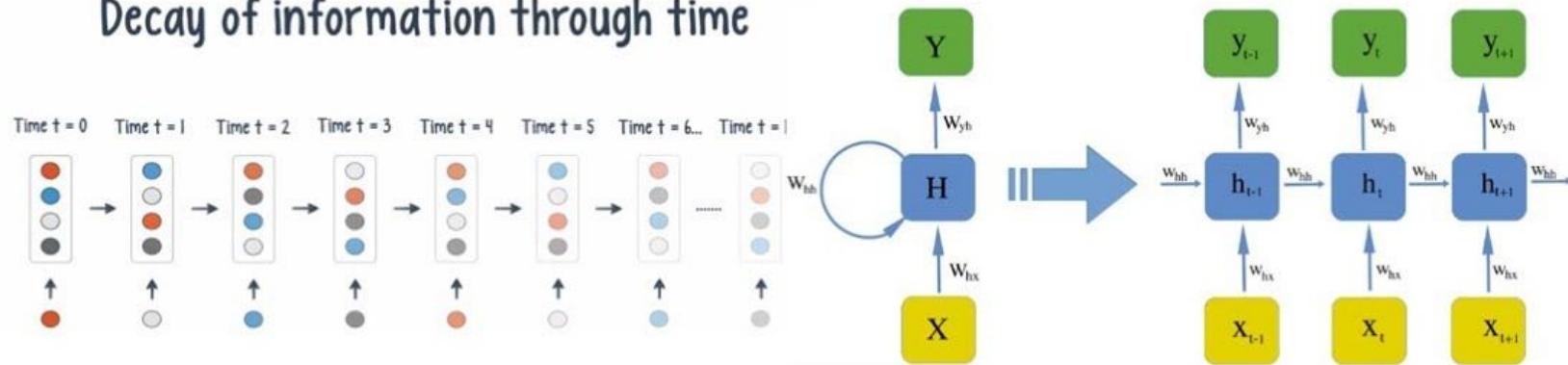
1. In Deep NN's , the derivatives become extremely small or large can cause problem in training or converging of the Neural Network.

## Vanishing/exploding gradients

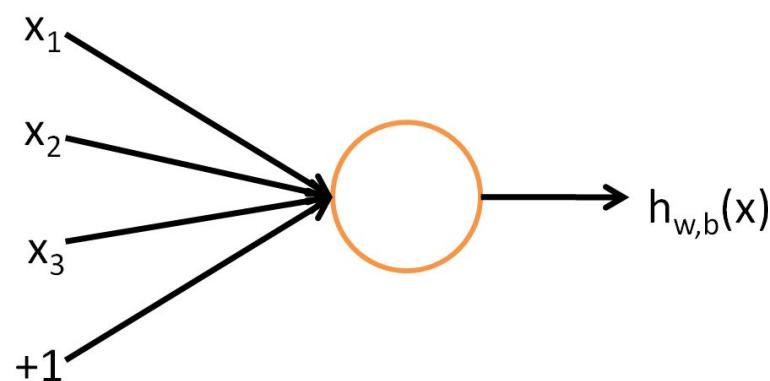


# Exploding Gradient Problem

Decay of information through time



# Weight Initialization



1.  $Z = W_1X_1 + W_2X_2 + W_3X_3 \dots$
2. Larger the value of  $n$ , smaller the value of  $W_i$  required so that  $z$  do not explode.
3.  $\text{Variance}(W_i) = 1/n$

# Weight Initialization

1. Relu =>  $2/(n^{[l-1]})$ .
2. Tanh =>  $1/(n^{[l-1]})$  (Xavier Initialization)
3. =>  $2/(n^{[l-1]} + n^{[l]})$

```
In [25]: def tanh(x): return torch.tanh(x)
```

```
In [26]: x = torch.randn(512)

for i in range(100):
    a = torch.randn(512,512) * math.sqrt(1./512)
    x = tanh(a @ x)
x.mean(), x.std()
```

```
Out[26]: (tensor(-0.0034), tensor(0.0613))
```



# Gradient Checking

Helps in finding bugs in implementations of back propagation.



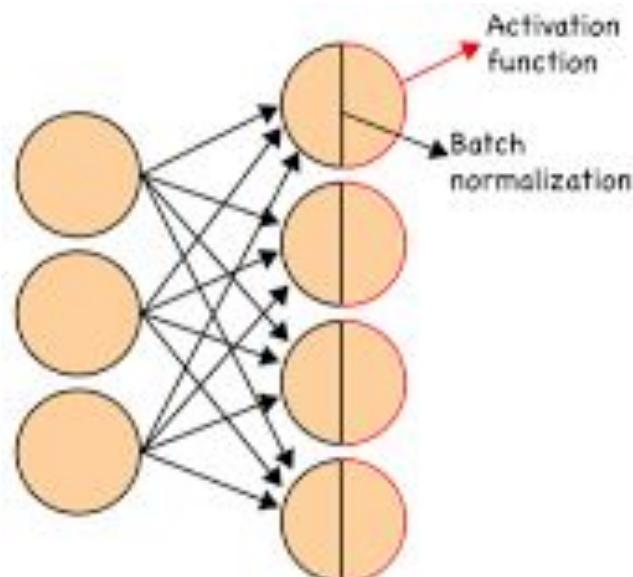
# Batch Normalization

1. Makes hyperparameter search problem much easier.
2. I.e. The choice of hyperparameters is a much bigger range of hyperparameters that work well, and will also enable you to much more easily train even very deep networks.



# Batch Normalization

1. Batch Normalizing after a layer  $l-1$  will help in efficient training of the weights and biases of the layer  $l$ .
2. Normalization need not be mean = 0 and variance = 1.
3. Batch Norm generally applied on the z's.



$$\begin{aligned}\hat{x}_i &\leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \\ y_i &\leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)\end{aligned}$$



1. These beta and gamma [of each layer] are trainable parameters. **Show over 2 layers\*\*\*\*\***
2. Any optimization algorithm can be used to train these parameters.
3. Programming Frameworks allows you just one line of code.
4. BN is applied generally on mini batches of data.
5. Show how bias parameter is irrelevant in applying batch norm. Beta (BN parameter) plays its role there then.



$$\mathbf{X} = \begin{bmatrix} X_{1,1} & X_{1,2} & \dots & X_{1,n} \\ X_{2,1} & X_{2,2} & \dots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{m,1} & X_{m,2} & \dots & X_{m,n} \end{bmatrix} \quad (1)$$

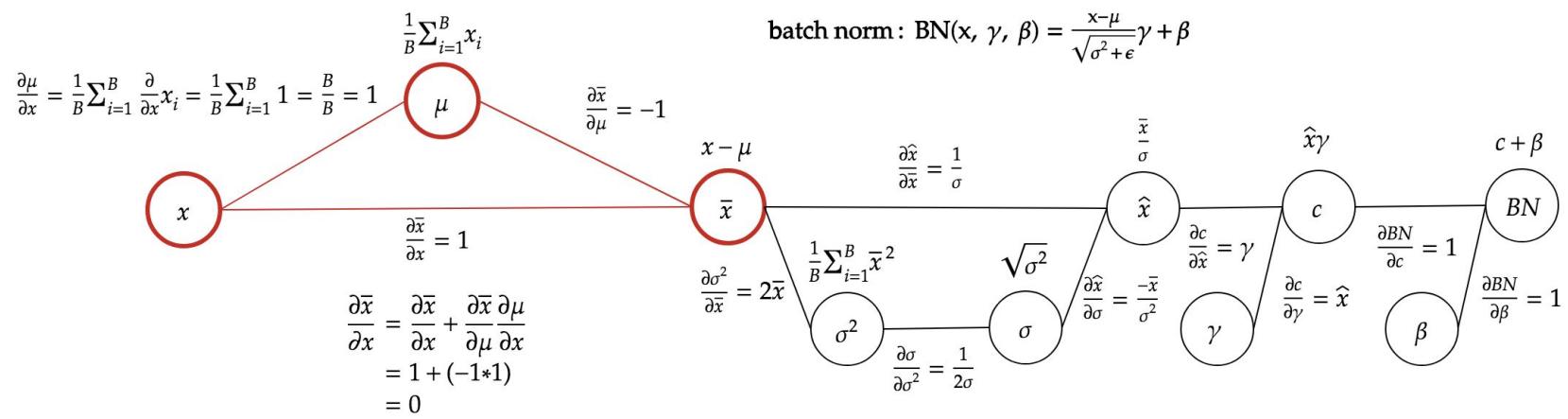
$$\hat{\mathbf{X}} = (\mathbf{X} - \boldsymbol{\mu}^\top) \odot \hat{\mathbf{a}}^\top$$

$$\hat{\mathbf{X}} = \begin{bmatrix} \frac{X_{1,1}-\mu_1}{\sqrt{\sigma_1^2+\epsilon}} & \frac{X_{1,2}-\mu_2}{\sqrt{\sigma_2^2+\epsilon}} & \dots & \frac{X_{1,n}-\mu_n}{\sqrt{\sigma_n^2+\epsilon}} \\ \frac{X_{2,1}-\mu_1}{\sqrt{\sigma_1^2+\epsilon}} & \frac{X_{2,2}-\mu_2}{\sqrt{\sigma_2^2+\epsilon}} & \dots & \frac{X_{2,n}-\mu_n}{\sqrt{\sigma_n^2+\epsilon}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{X_{m,1}-\mu_1}{\sqrt{\sigma_1^2+\epsilon}} & \frac{X_{m,2}-\mu_2}{\sqrt{\sigma_2^2+\epsilon}} & \dots & \frac{X_{m,n}-\mu_n}{\sqrt{\sigma_n^2+\epsilon}} \end{bmatrix}$$

$$\mathbf{Y} = \boldsymbol{\gamma}^\top \odot \hat{\mathbf{X}} + \boldsymbol{\beta}^\top$$



# Back Prop with Batch Norm



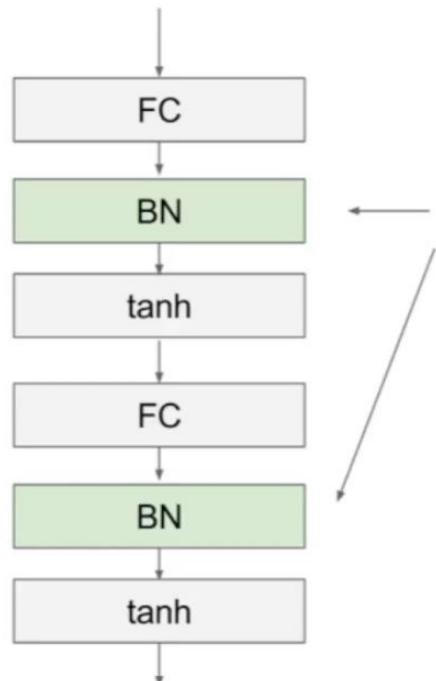
$$\frac{\partial BN}{\partial x} = \left( \frac{\partial BN}{\partial c} \frac{\partial c}{\partial \hat{x}} \frac{\partial \hat{x}}{\partial x} + \frac{\partial BN}{\partial \hat{x}} \frac{\partial c}{\partial \sigma} \frac{\partial \sigma}{\partial \sigma^2} \frac{\partial \sigma^2}{\partial x} \right) \frac{\partial \bar{x}}{\partial \mu} \frac{\partial \mu}{\partial x} + \frac{\partial \bar{x}}{\partial x}$$

<https://datascience.stackexchange.com/questions/30056/deriving-the-gradient-of-batch-normalization>



# Batch Normalization

[Ioffe and Szegedy, 2015]



Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Stanford



# Batch Norm over a mini batch

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

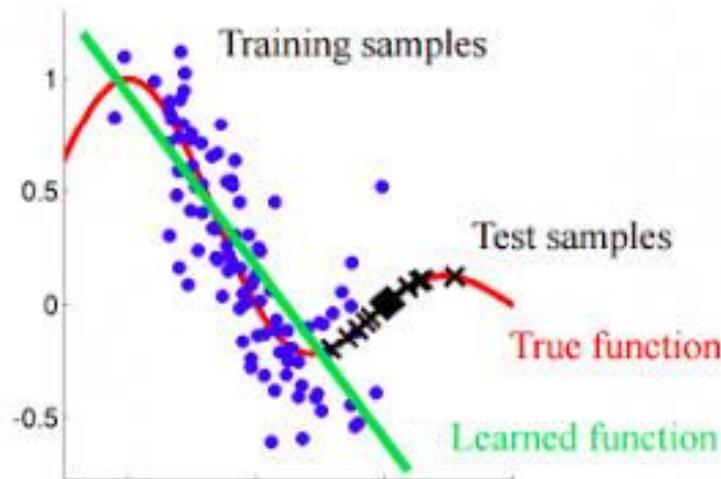
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

**Algorithm 1:** Batch Normalizing Transform, applied to activation  $x$  over a mini-batch.

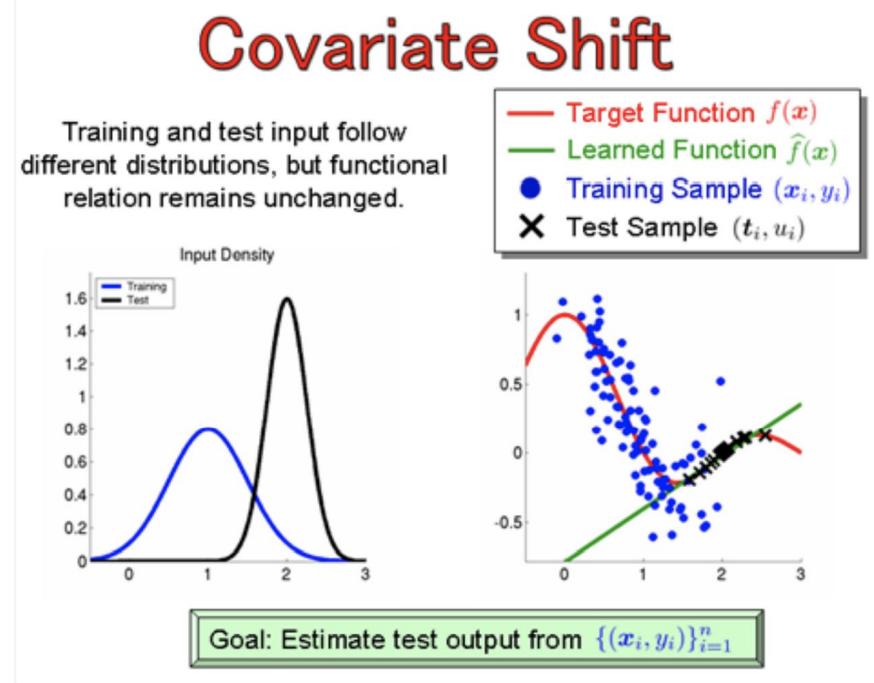
[From the original batch-norm paper](#)



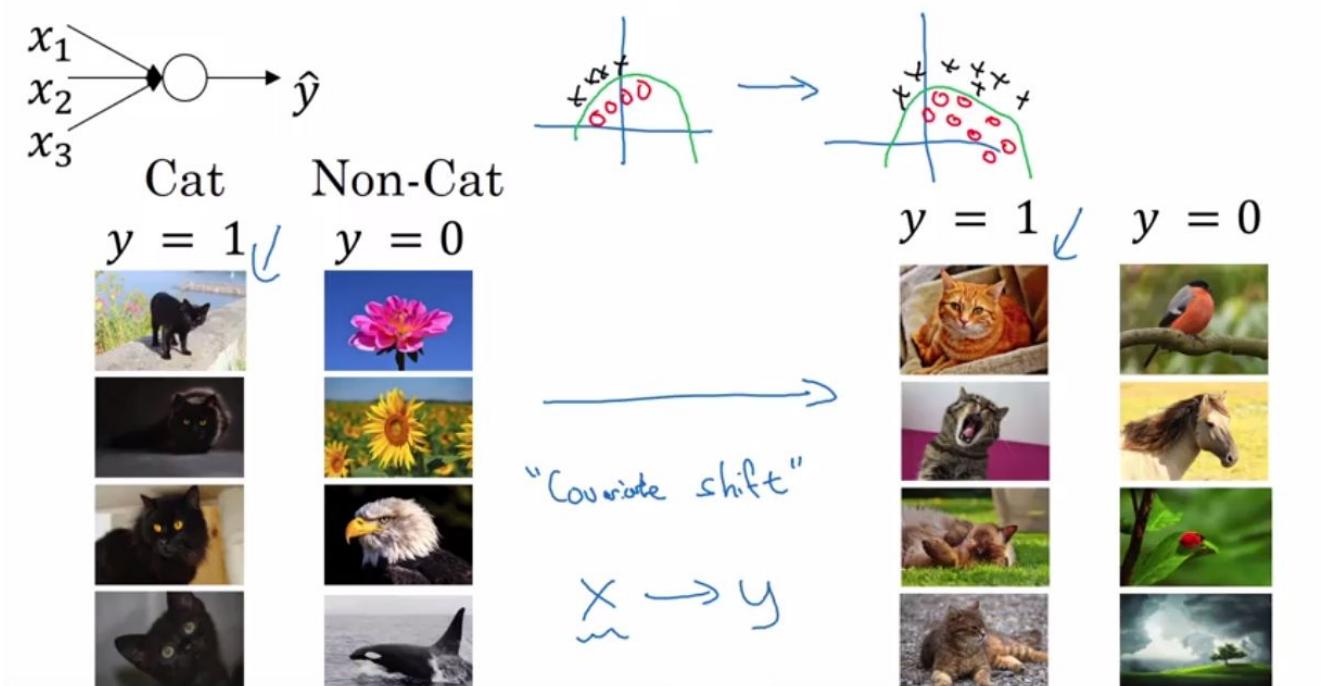
# Covariate Shift



# Covariate Shift



# Covariate Shift



Deeplearning.ai: Why Does Batch Norm Work? ([C2W3L06](#))

# Paper Presentations

