

A large, abstract neural network diagram is positioned on the left side of the slide. It consists of numerous small, colored dots (representing neurons) connected by thin, curved lines (representing synapses). The colors of the dots transition through a spectrum, including yellow, green, blue, purple, and red. The network is highly branched, with many dots having multiple connections to other dots.

Week 1:

Introduction to Deep Learning

PCLUB CS2020

Ankur Bhatia

January 25, 2020



Course Prerequisites - Linear Algebra Calculus

'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.



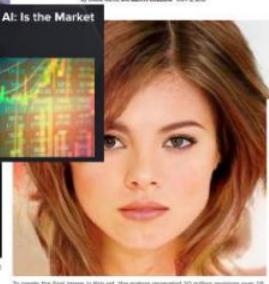
'Creative' AlphaZero leads way for chess computers and, maybe, science

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



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

By CAROL MITCHELL AND KEITH COLLINS | JAN. 2, 2018



To

create the final image in this set, the system generated 30 million images over 18 days.

Google's DeepMind aces protein folding

By Robert K. Serrano | Dec. 6, 2018, 12:05 PM

Complex of bacteria-infesting viral proteins modeled in CASP13. The complex consists of 10 proteins that were modeled individually. PRETERDATA BASE

protein folding

The Rise of Deep Learning

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



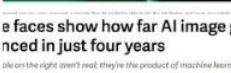
Technology outpacing security measures

| Facial Recognition | Features and Limitations



AI beats docs in cancer spotting

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



Automation And Algorithms: De-Risking Manufacturing With Artificial Intelligence

Sarah Geohrike Contributing Writer

Focus on the industrialization of additive manufacturing

TWEET THIS

The two key players in the market are

AI Can Help In Predicting Cryptocurrency Value



6.S191 Introduction to Deep Learning
introtodeeplearning.com



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/dlwikipclub



Why Deep Learning and Why Now?

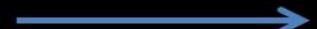
Things we want to do with data

Images



Label image

Audio



Speech recognition

Text

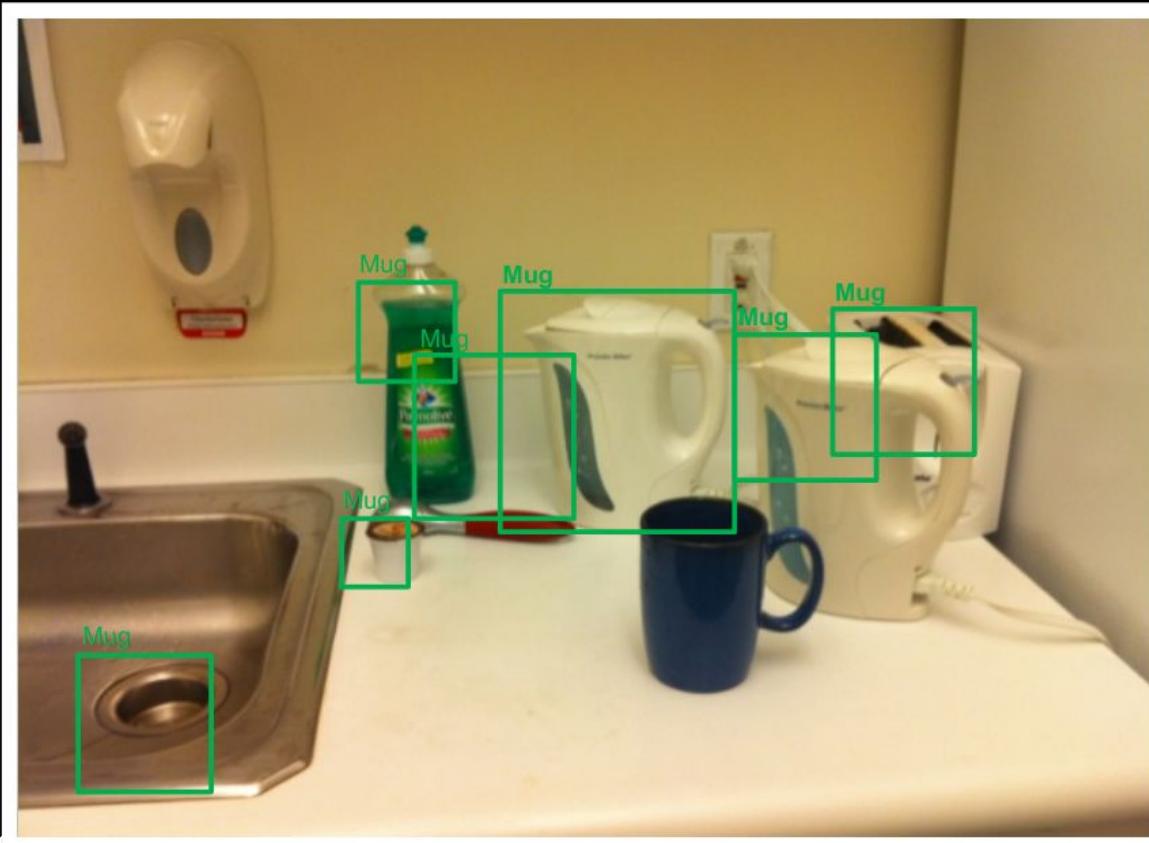


Web search

Computer vision: Identify coffee mug



Computer vision: Identify coffee mug



Why is computer vision hard?



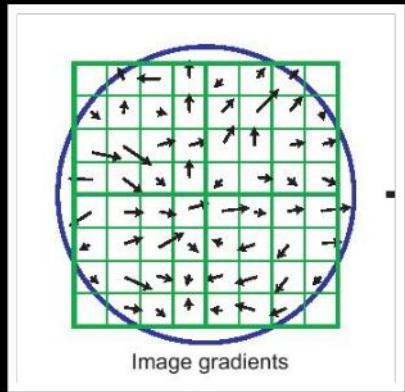
The camera sees :

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

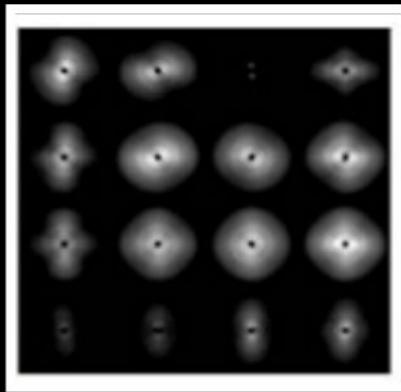
Computer vision



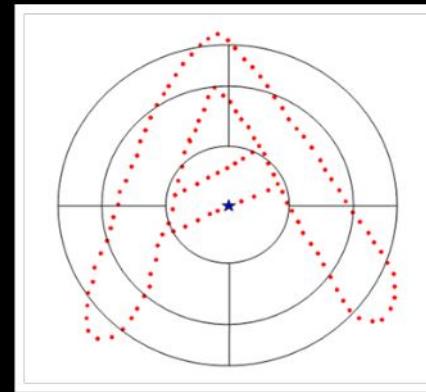
Features for vision



SIFT



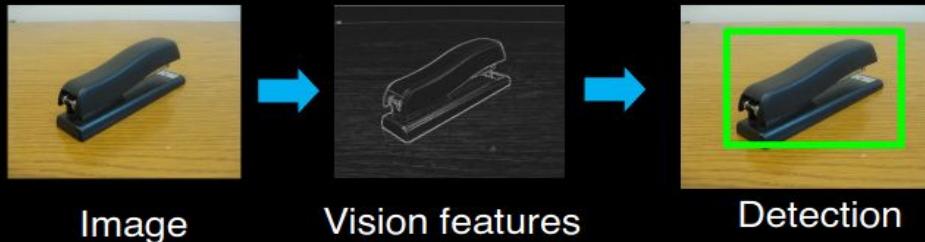
GIST



Shape context

Features for machine learning

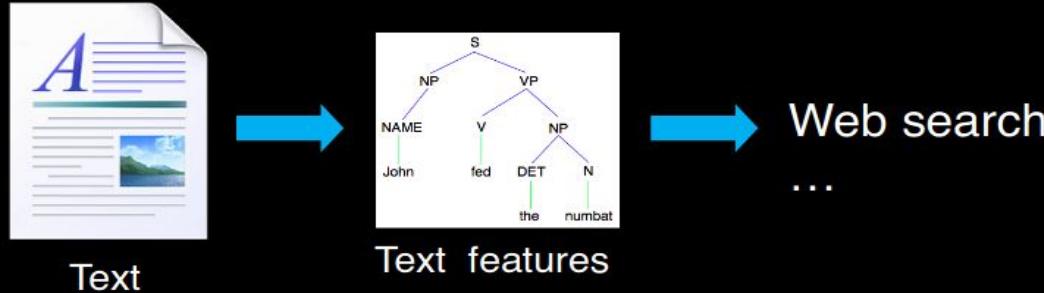
Images



Audio

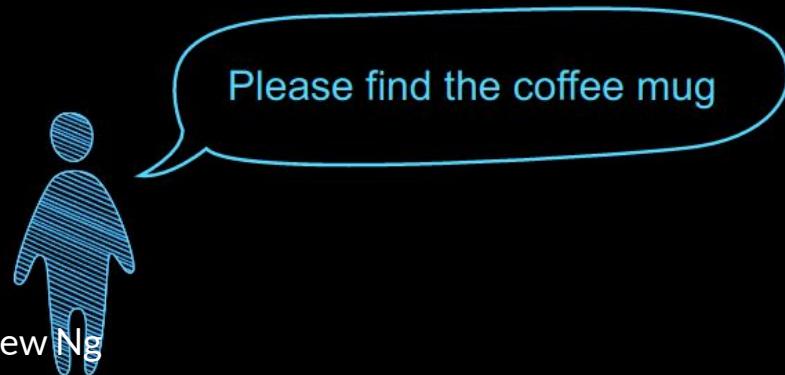
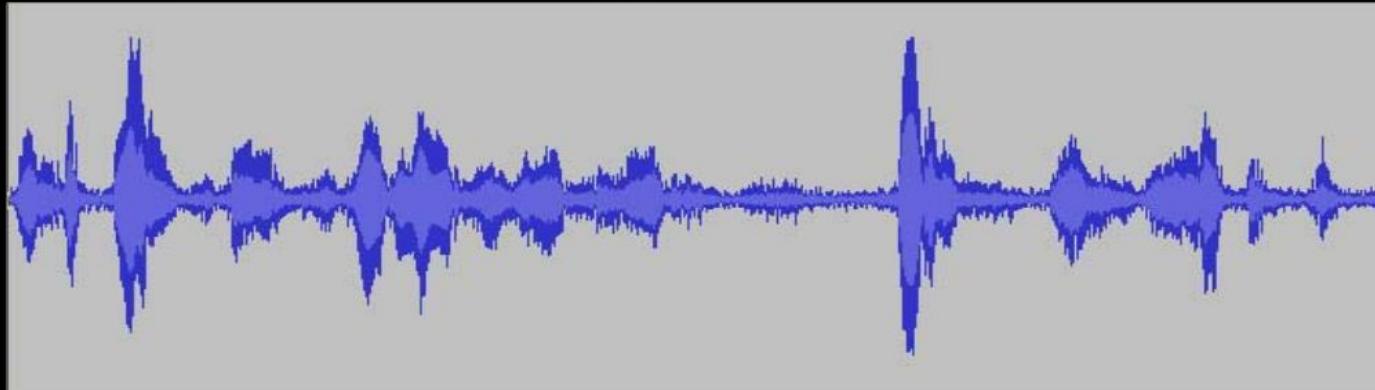


Text



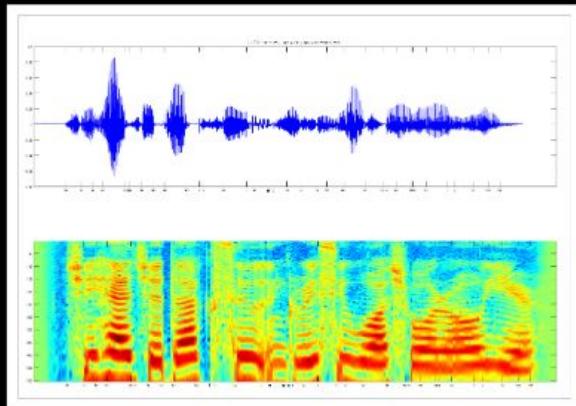
Why is speech recognition hard?

Microphone recording:

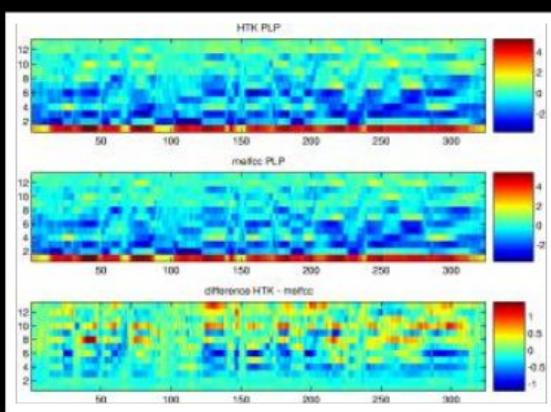


Andrew Ng

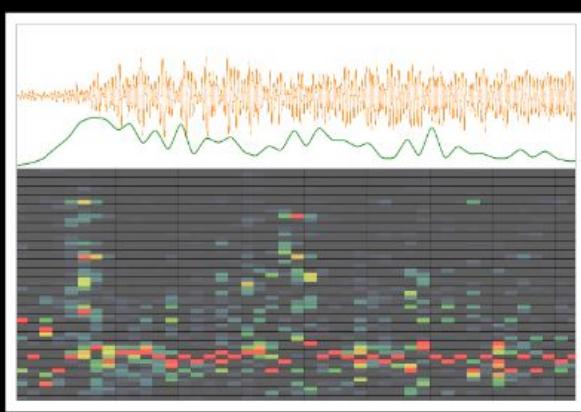
Features for audio



Spectrogram

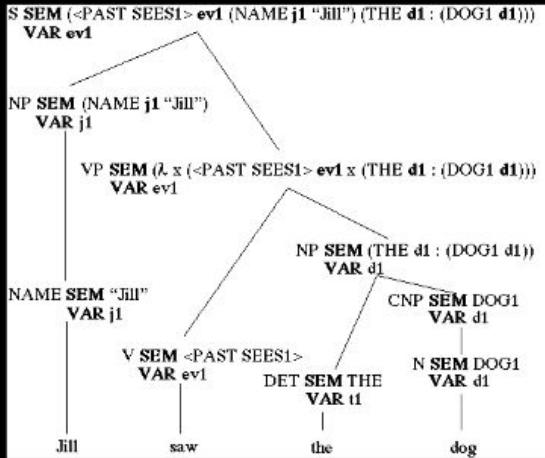


MFCC



Flux

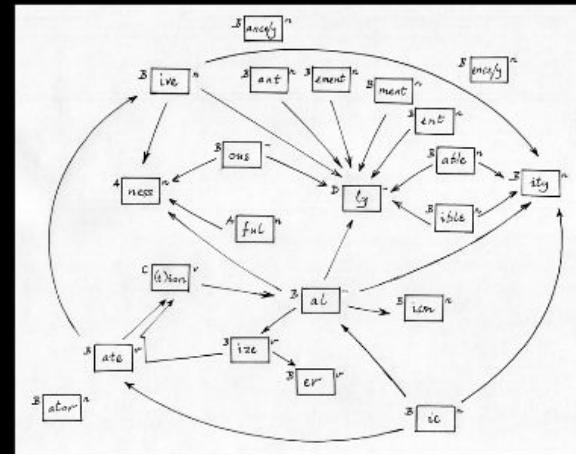
Features for text



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 o  
</p>  
</TXT>  
</DOC>
```

Named entity



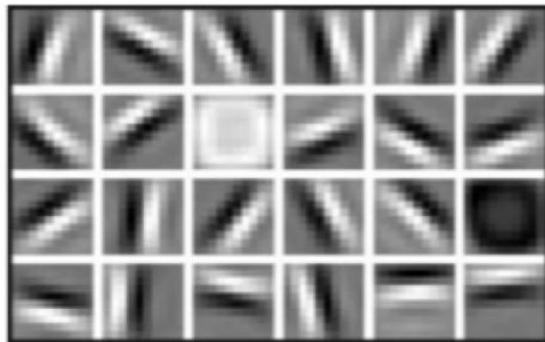
Stemming

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



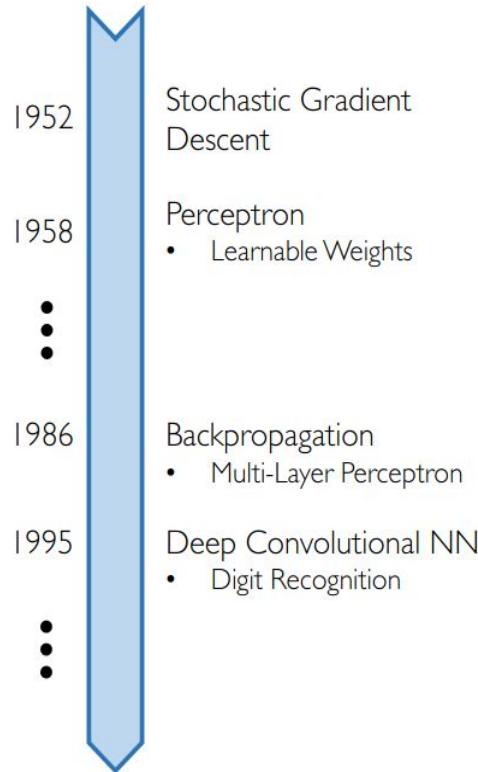
Eyes & Nose & Ears

High Level Features



Facial Structure

Why Now?



Neural Networks date back decades, so why the resurgence?

I. Big Data

- Larger Datasets
- Easier Collection & Storage



WIKIPEDIA
The Free Encyclopedia



2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



3. Software

- Improved Techniques
- New Models
- Toolboxes



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

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

Turing Award for Deep Learning



- 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

Untagged data (unsupervised learning)



Unknown



Unknown



Unknown



Unknown



Unknown

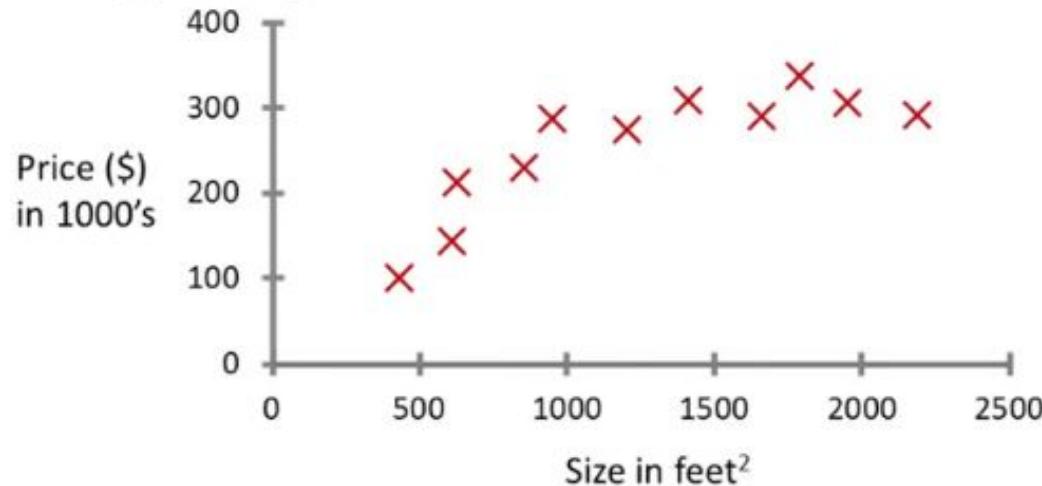


Unknown

Supervised Learning

EXAMPLE 1

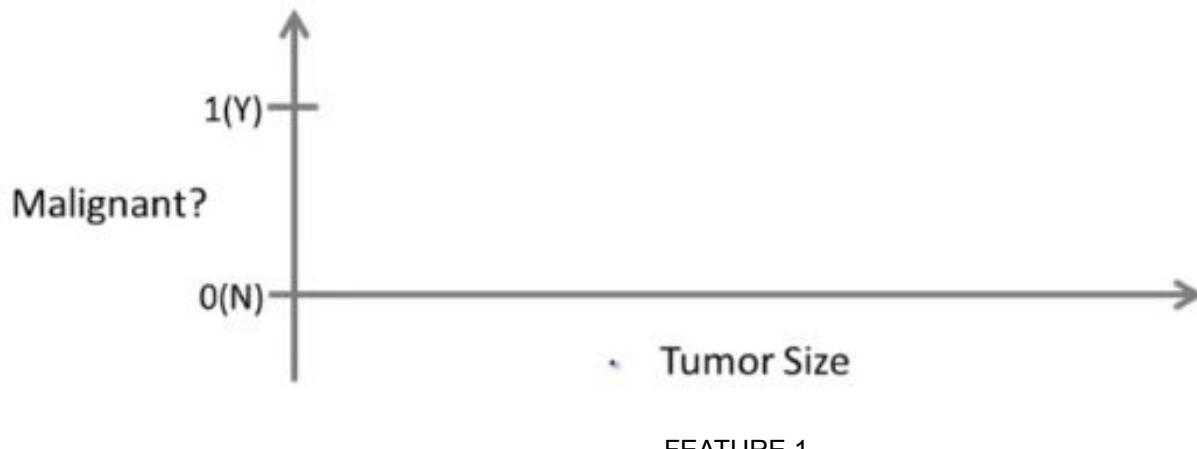
Housing price prediction.



Supervised Learning

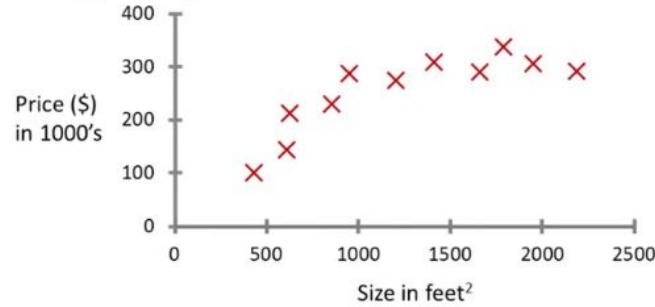
EXAMPLE 2

Breast cancer (malignant, benign)



Linear Regression:

Housing price prediction.



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

Parameters: θ_0, θ_1

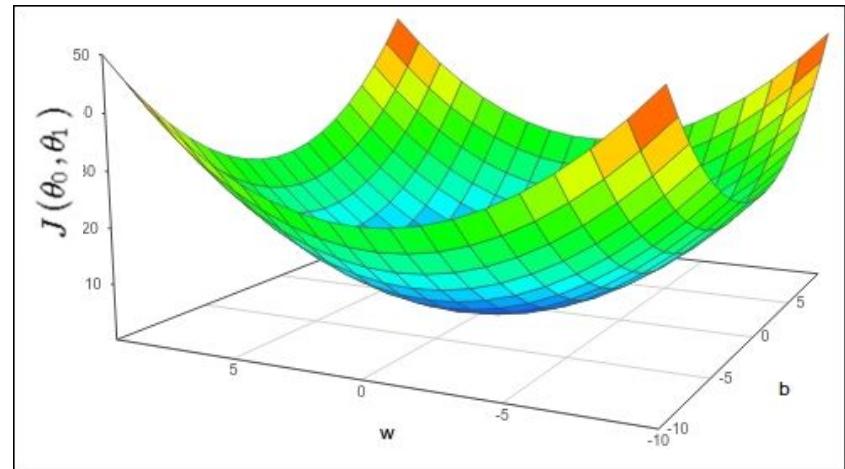
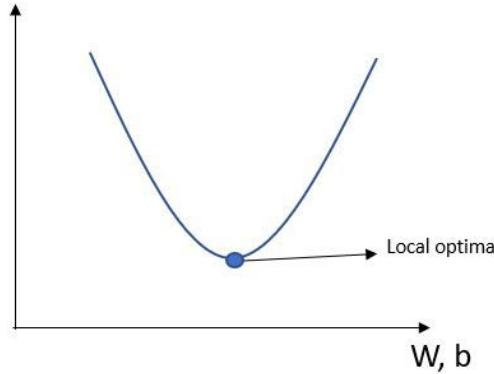
Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

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:

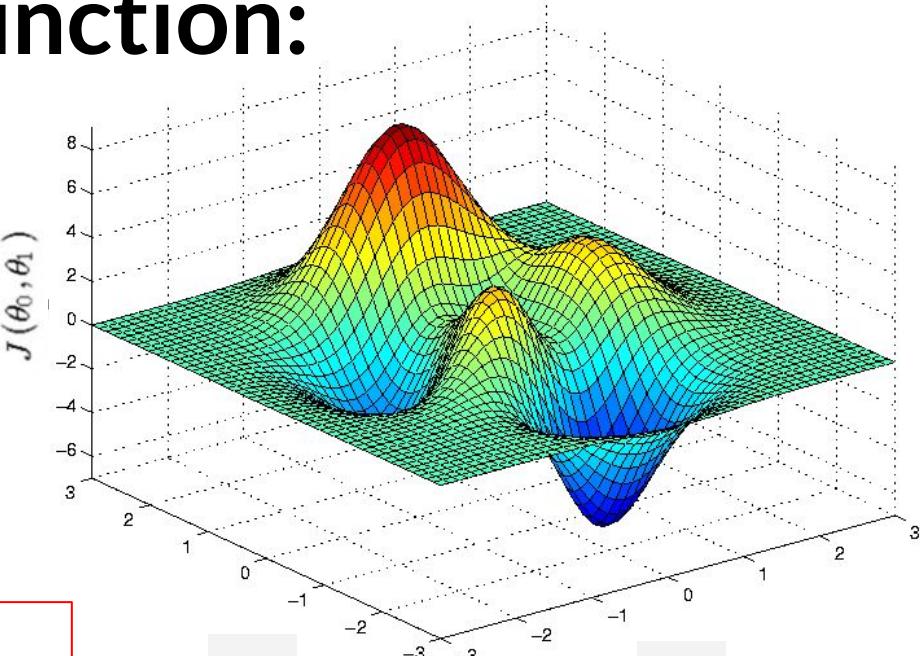
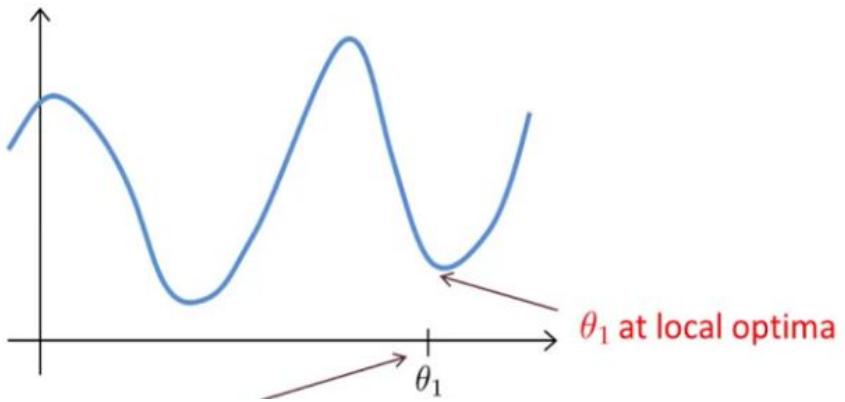
Cost Function (J)



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

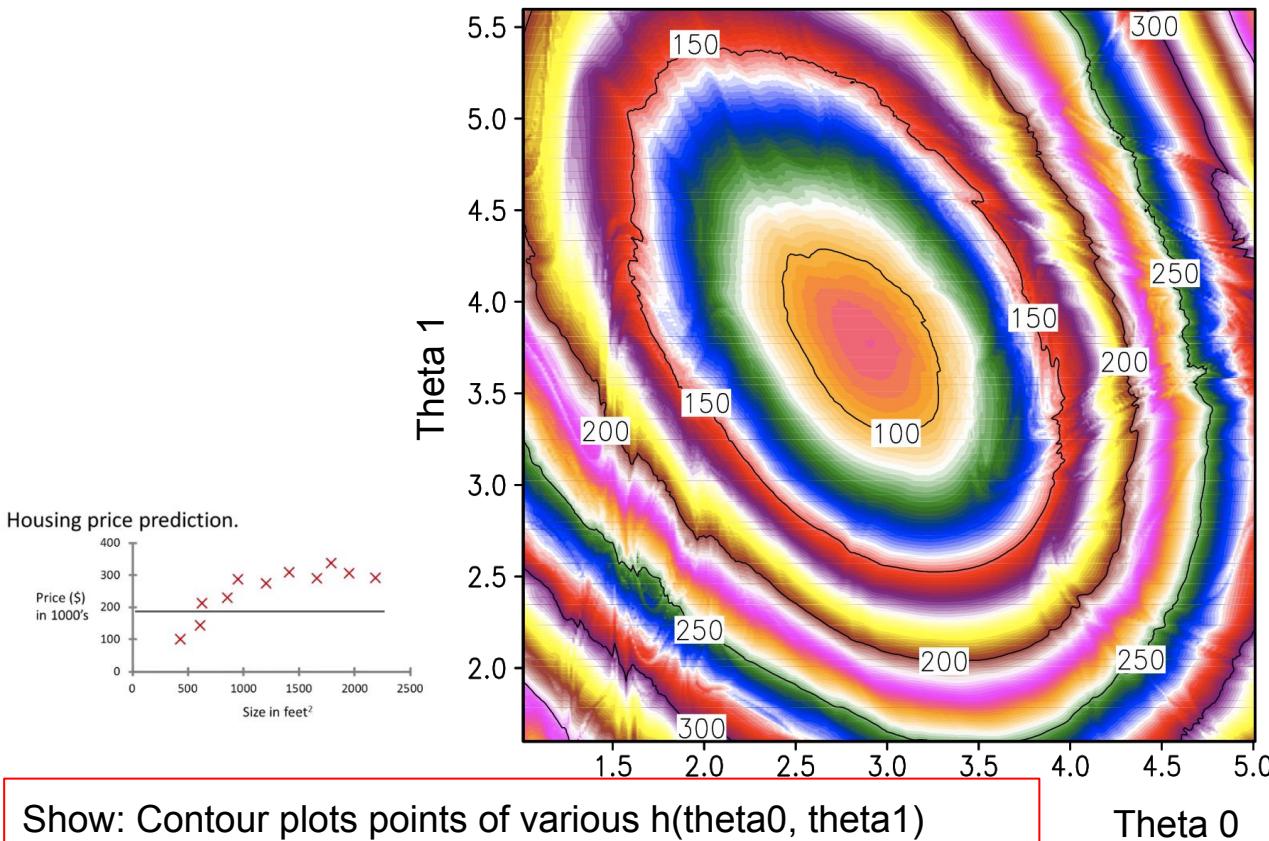
Convex Function

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:



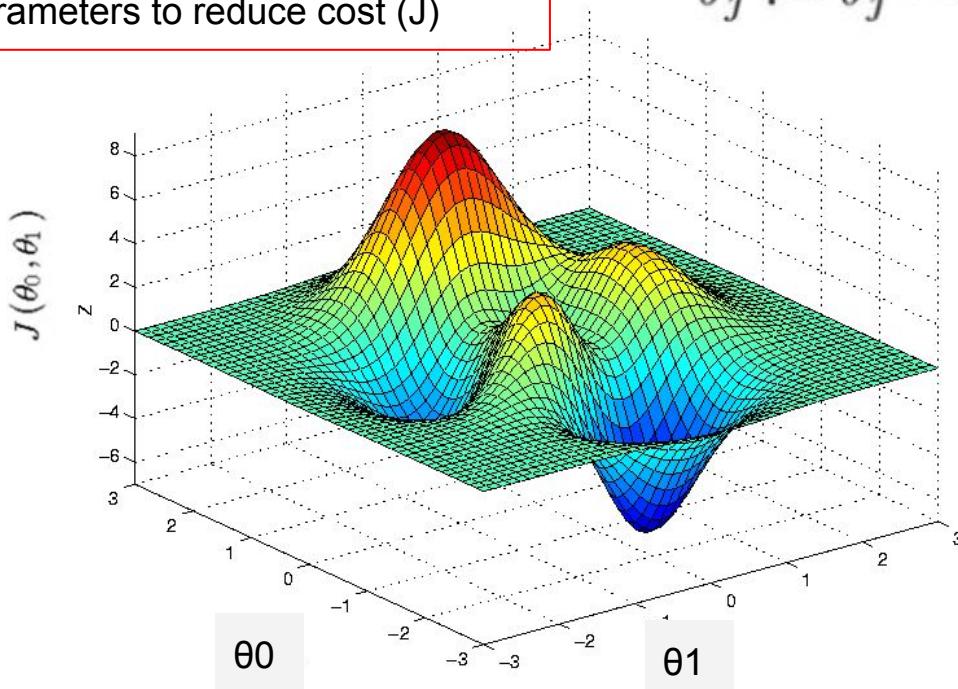
What does a point denote in this?

Gradient Descent “Batch”:

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

1. Start with some initial parameter value θ_1, θ_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.

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

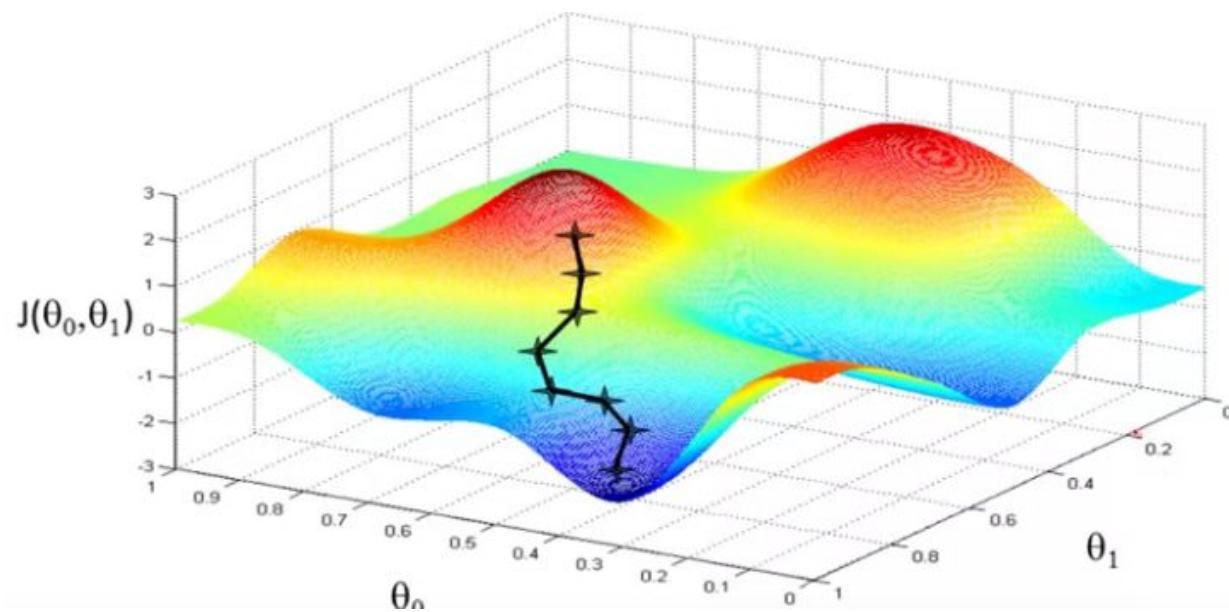
$$j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

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

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



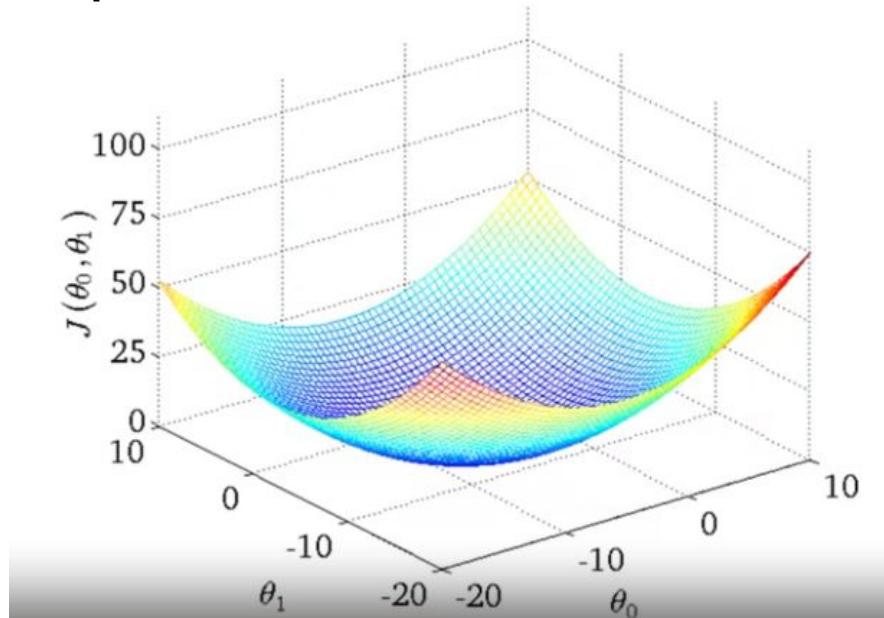
Gradient Descent for Linear Reg:

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

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

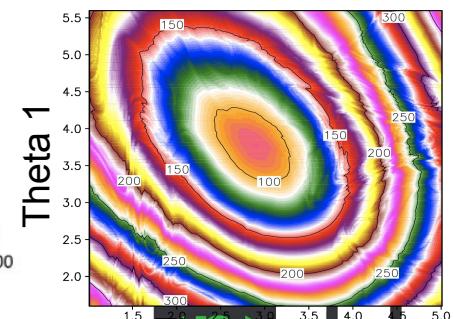
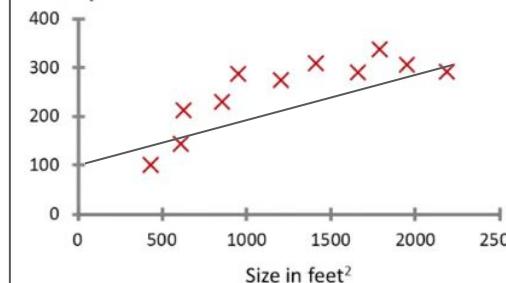
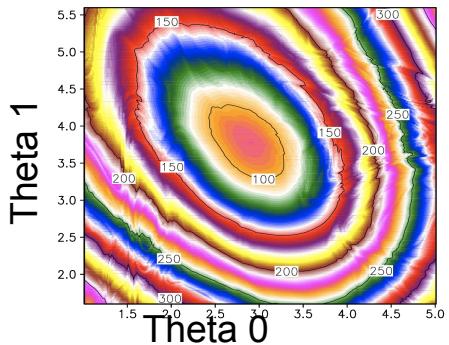
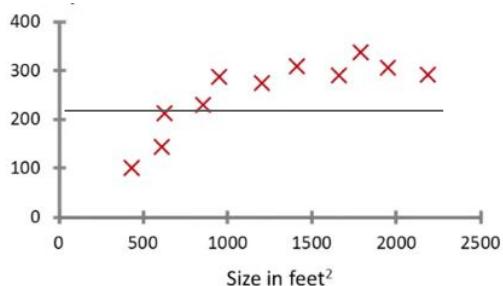
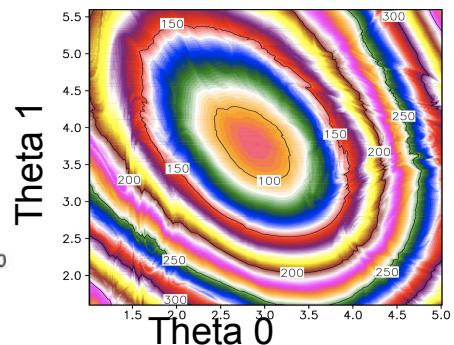
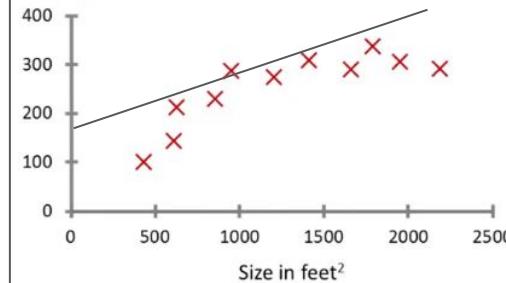
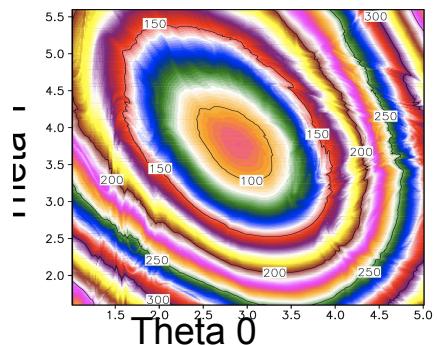
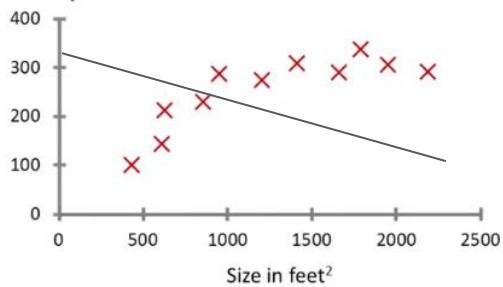
Ans. For Linear Reg., we have a bowl type cost function (Convex function).

Optimization problem is convex here.

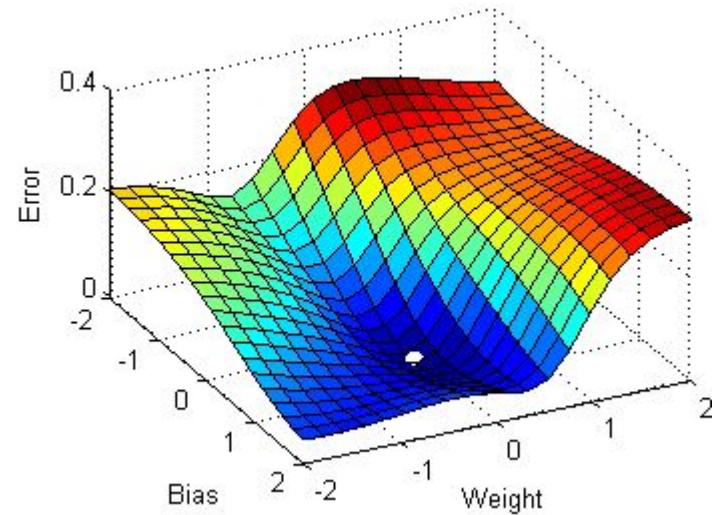
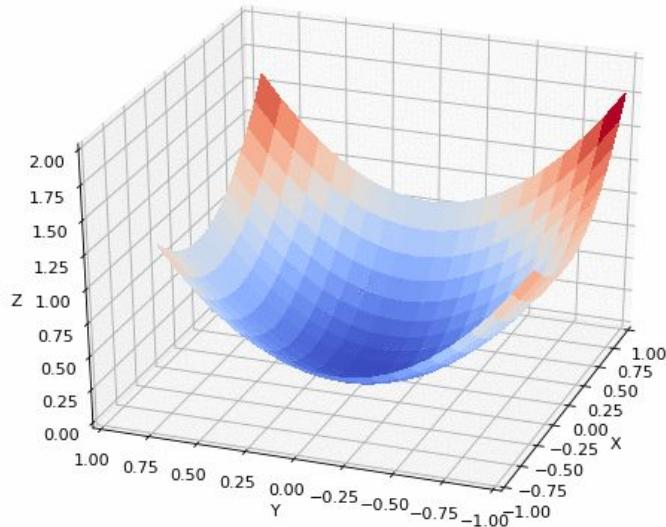


Gradient Descent for Linear Reg:

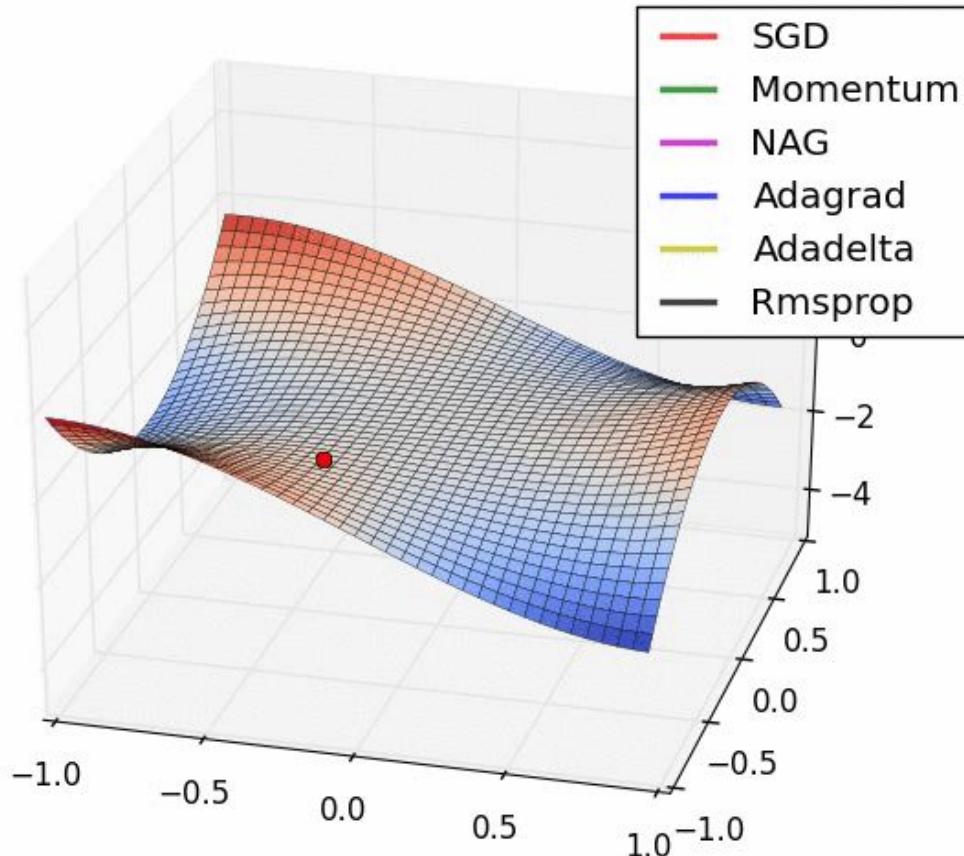
Task: Minimize $J(\theta_0, \theta_1) = \frac{1}{2m} (\sum(h(x) - y)^2)$



Gradient Descent for Linear Reg:



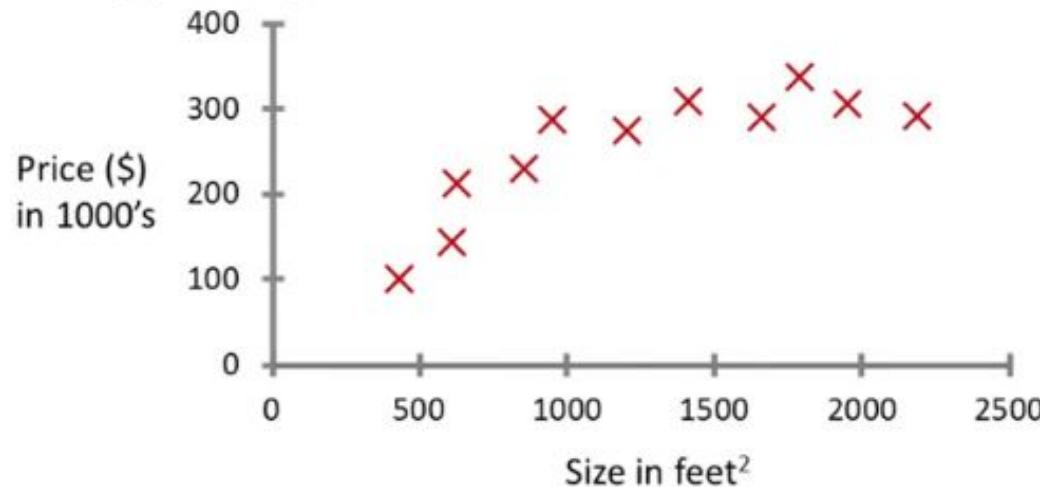
Various Descent's:



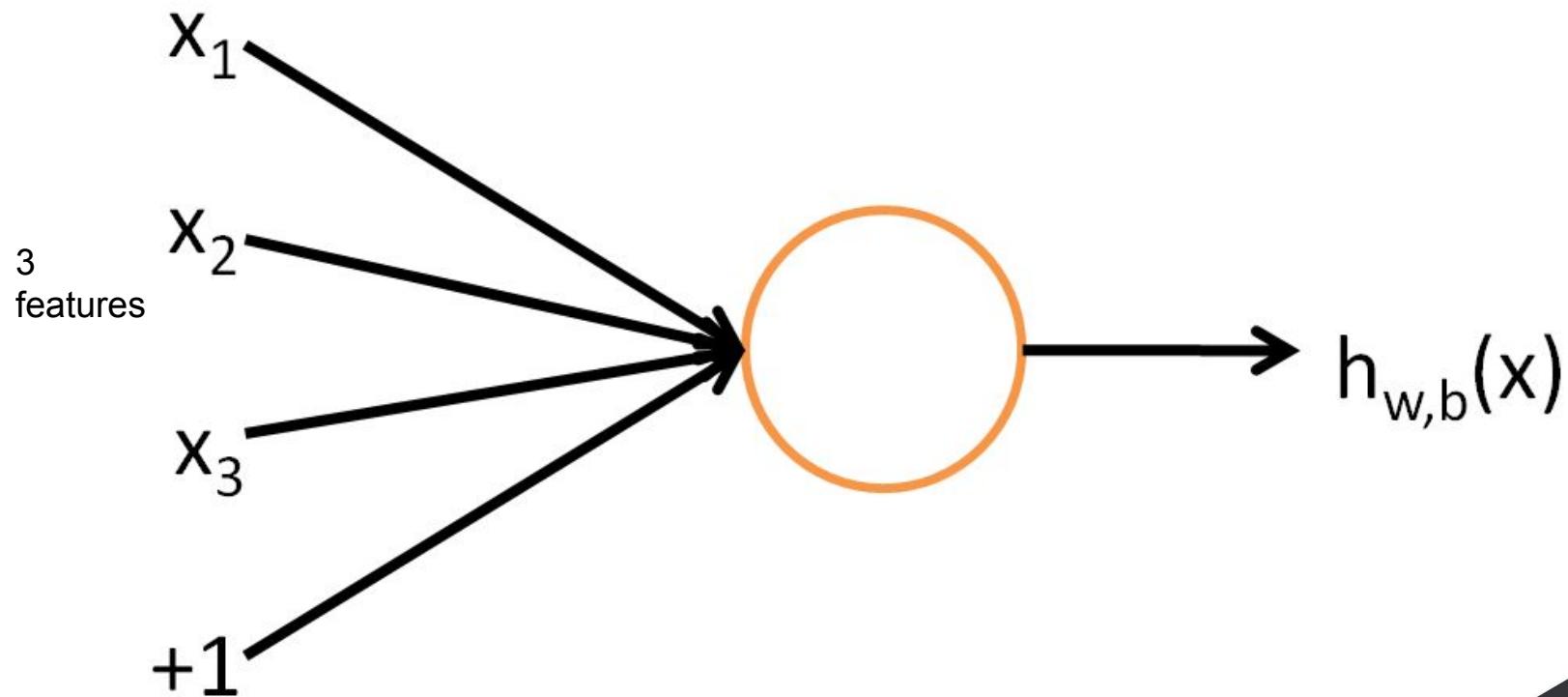
Now Let's Dive a bit more Deep.

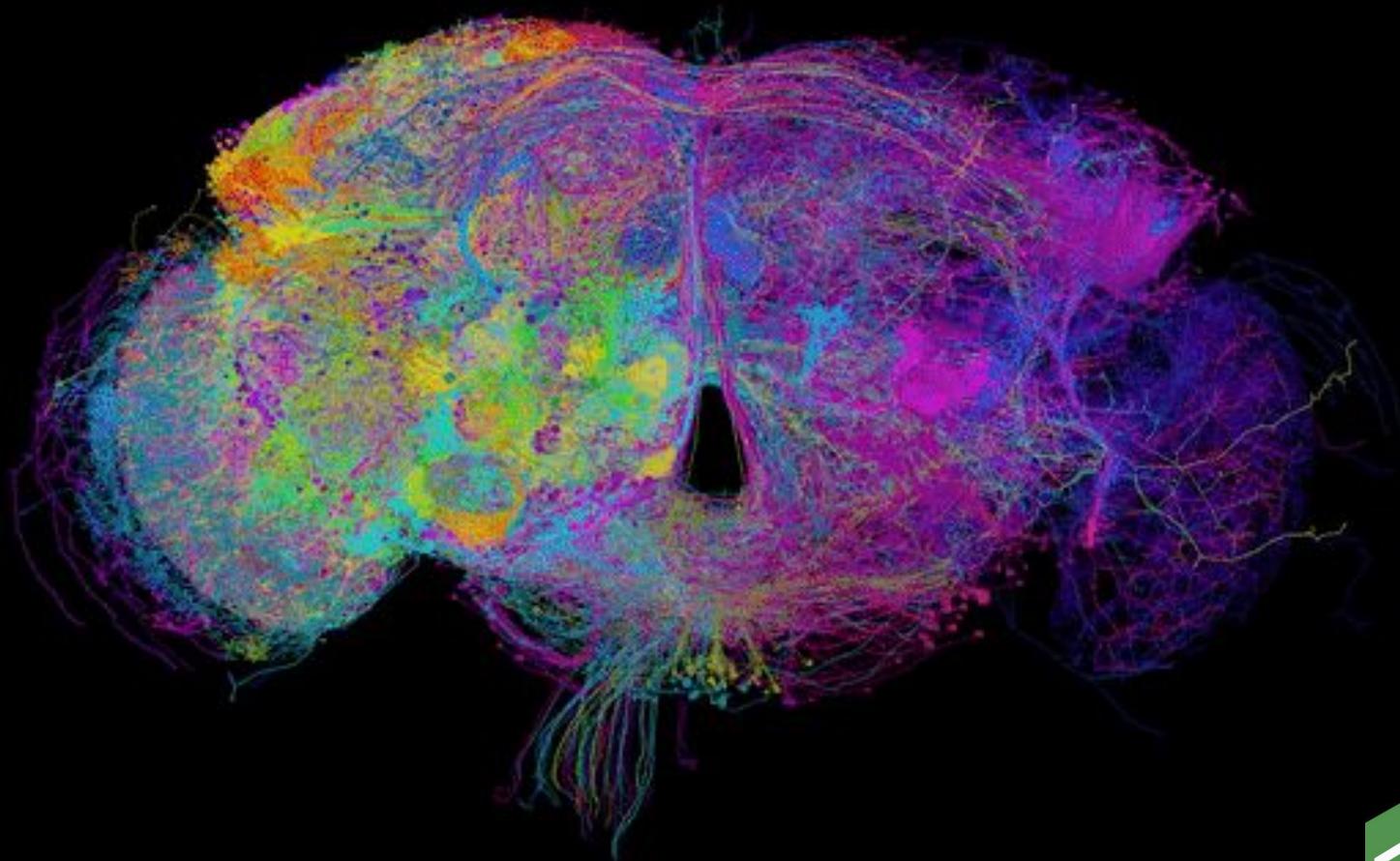
What is a Neural Network?

Housing price prediction.



What is a Neural Network?



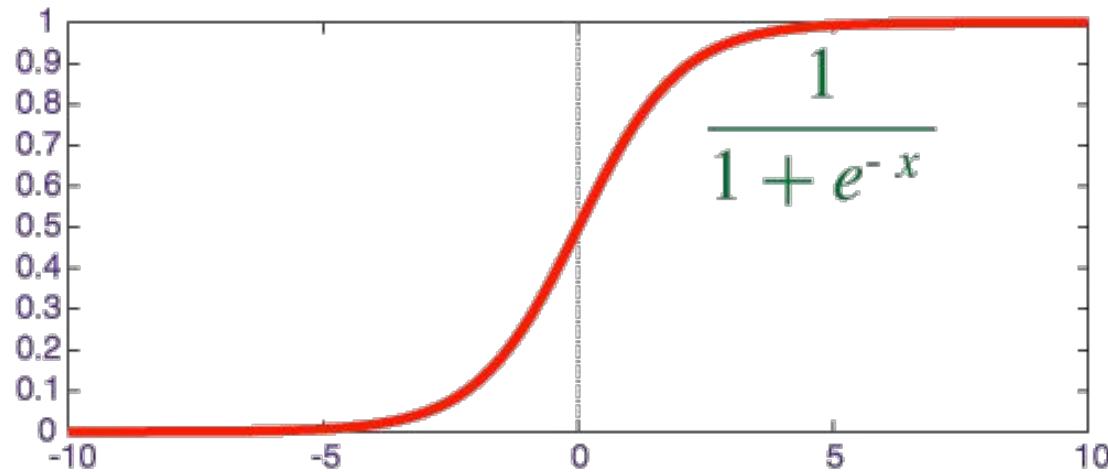


Binary Classification problem:



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

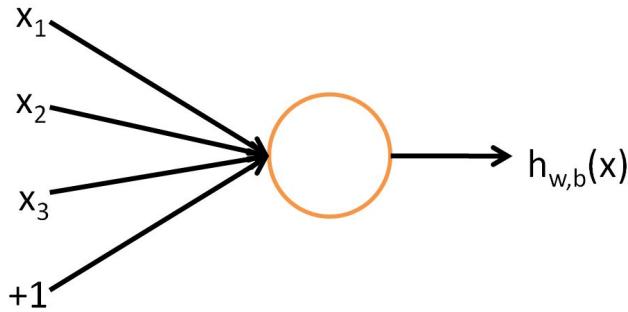
Understanding Logistic Regression problem:



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:

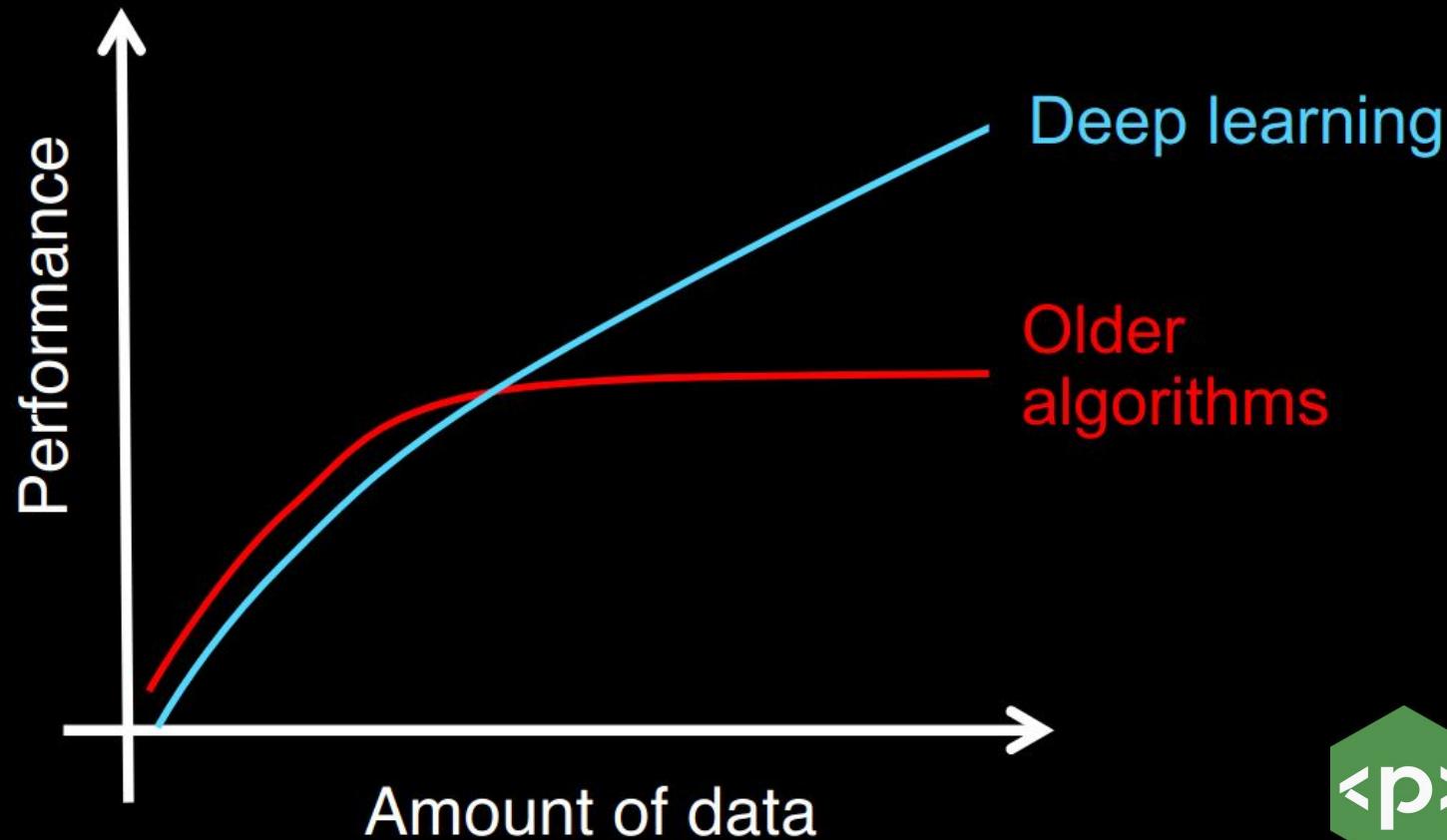
Backprop:

On m Training example

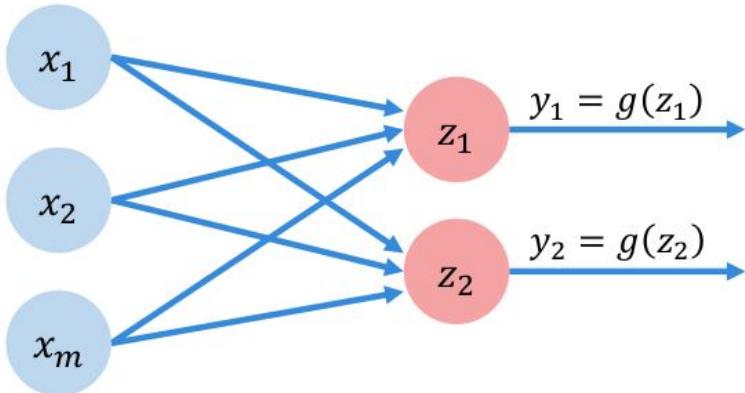
Now Let's Dive Deeeeeeeeep in?



Learning from tagged data

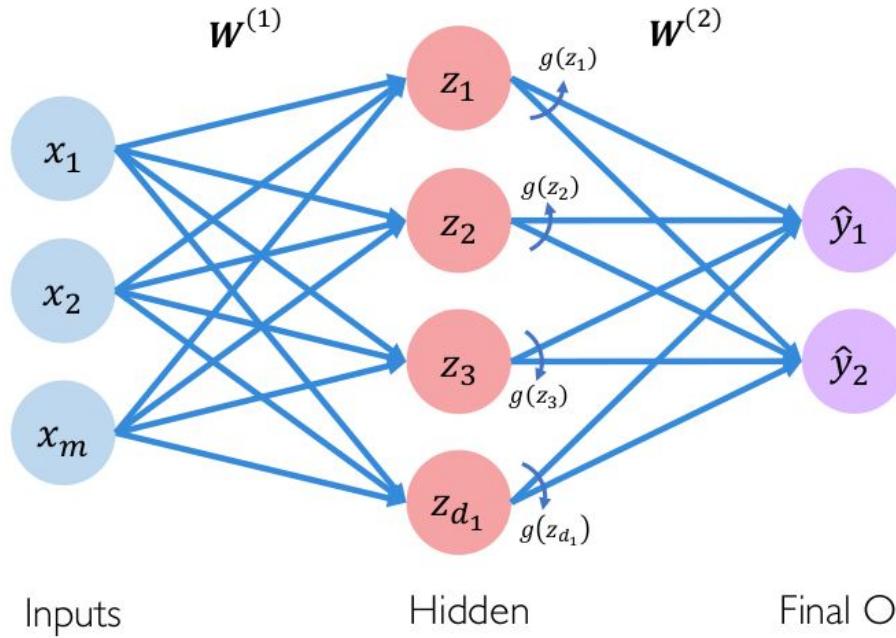


Multi Output Perceptron



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

Single Layer Neural Network

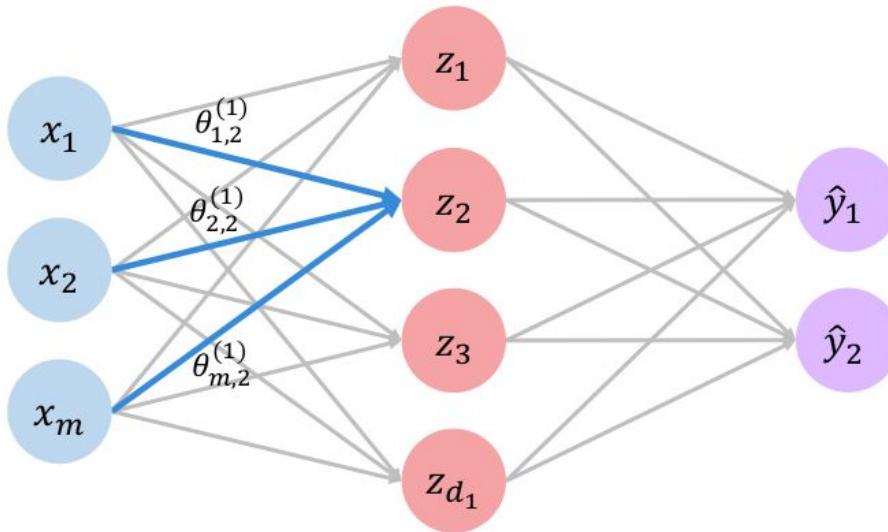


$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)} \quad \hat{y}_i = g\left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)}\right)$$

6.S191 Introduction to Deep Learning

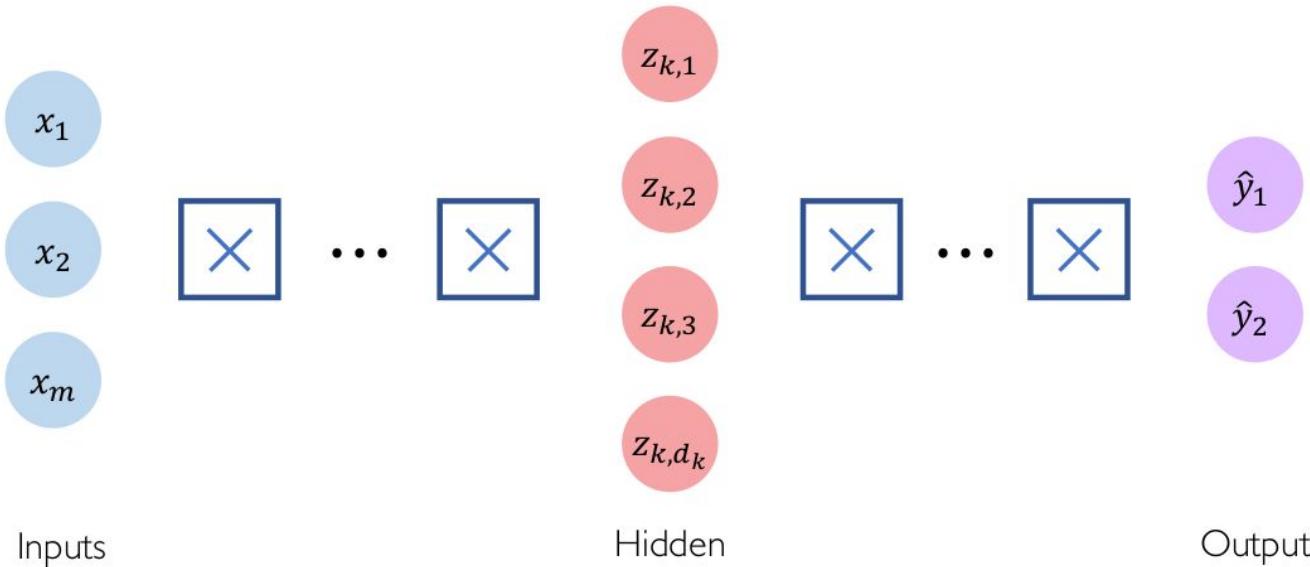
introtodeeplearning.com

Single Layer Neural Network



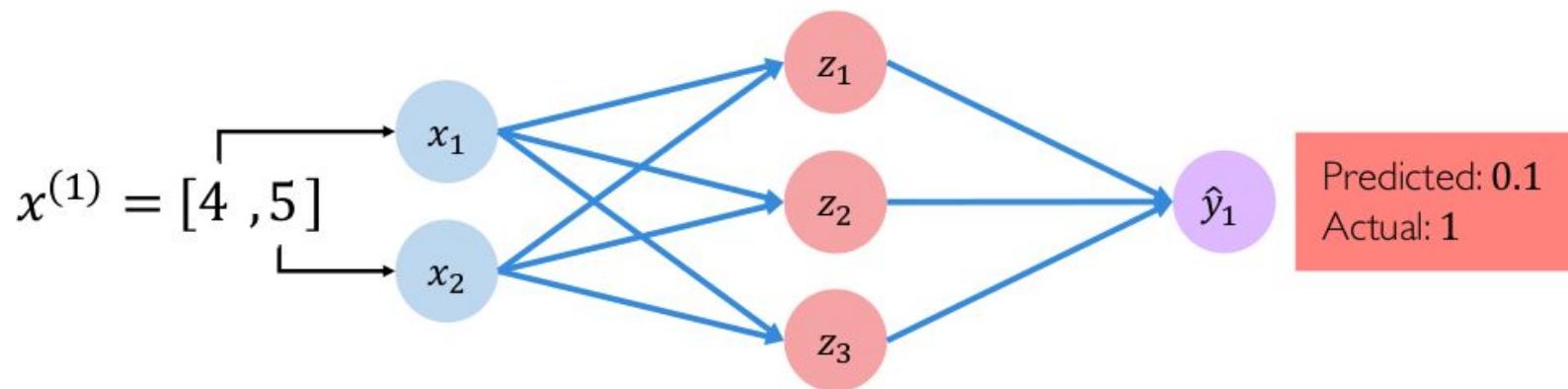
$$\begin{aligned} z_2 &= w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)} \\ &= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)} \end{aligned}$$

Deep Neural Network



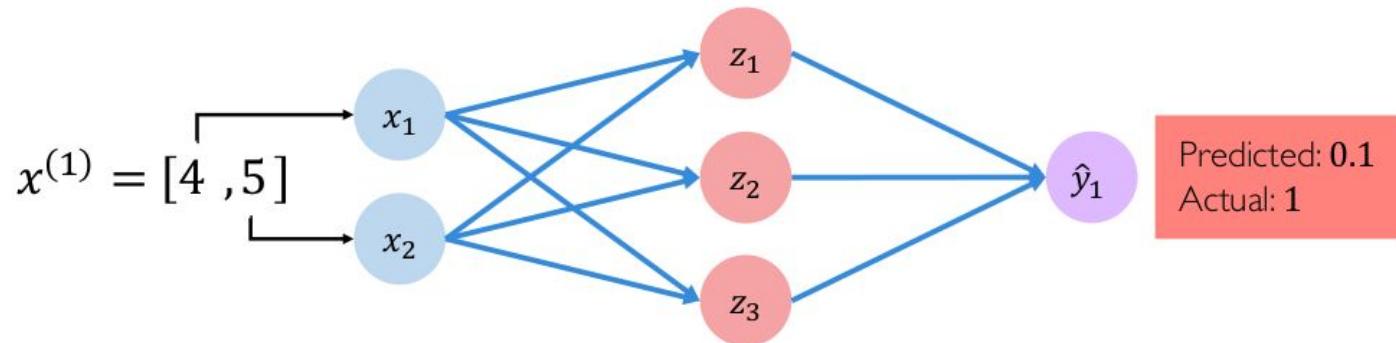
$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

Example Problem: Will I pass this class?



Quantifying Loss

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



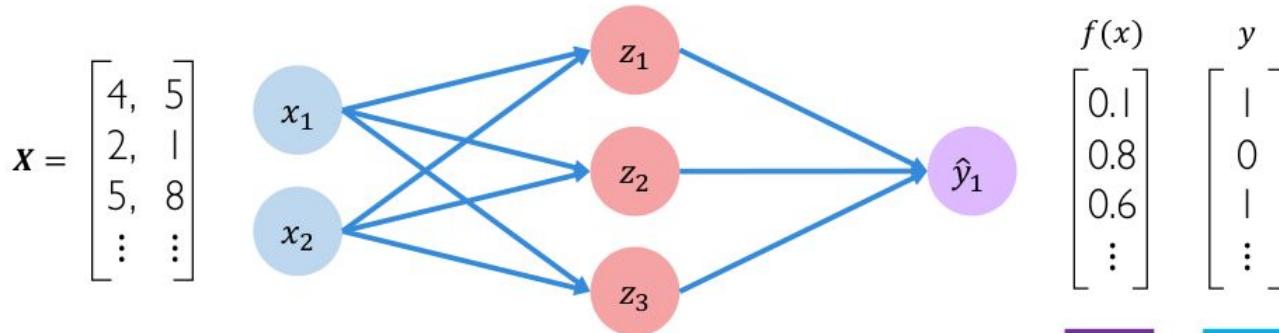
$$\mathcal{L} \left(\underline{f(x^{(i)}; \mathbf{W})}, \underline{y^{(i)}} \right)$$

Predicted Actual

Empirical Loss

(Cost Function)

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



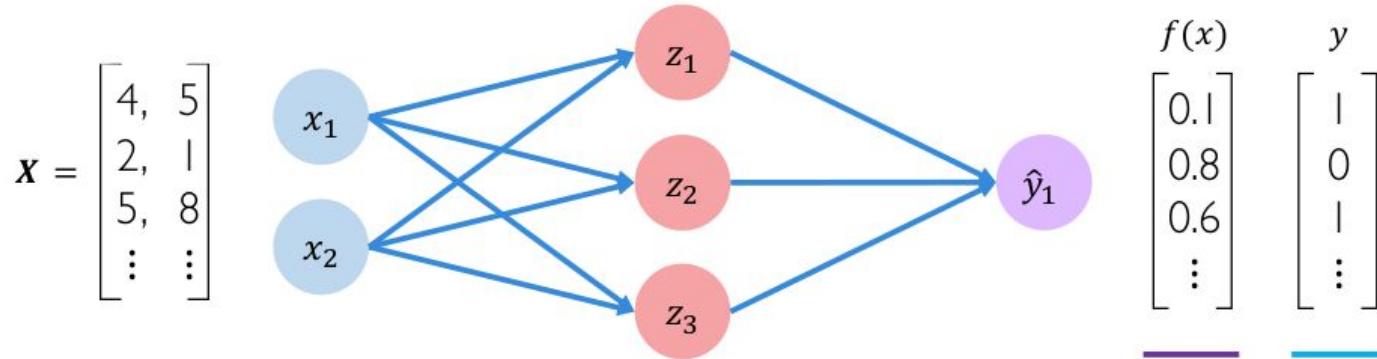
- Also known as:
- Objective function
 - Cost function
 - Empirical Risk

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

Predicted Actual

Binary Cross Entropy Loss

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

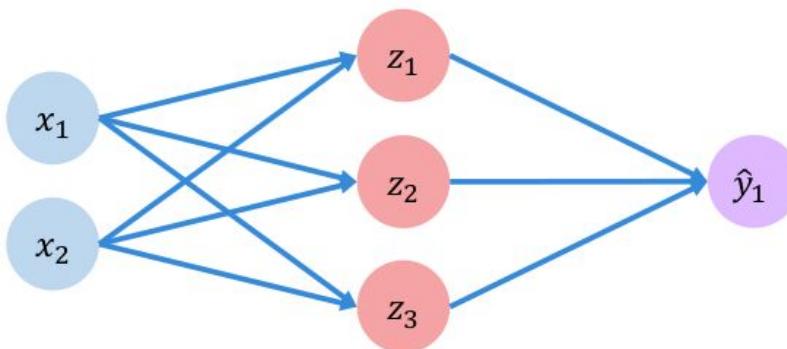


$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)} \log(f(x^{(i)}; \mathbf{W}))}_{\text{Actual}} + \underbrace{(1 - y^{(i)}) \log(1 - f(x^{(i)}; \mathbf{W}))}_{\text{Predicted}}$$

Mean Squared Error Loss

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

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



$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underbrace{(y^{(i)} - f(x^{(i)}; \mathbf{W}))^2}_{\text{Actual Predicted}}$$

$f(x)$	y
$\begin{bmatrix} 30 \\ 80 \\ 85 \\ \vdots \end{bmatrix}$	$\begin{bmatrix} 90 \\ 20 \\ 95 \\ \vdots \end{bmatrix}$

Final Grades
(percentage)

Loss Optimization

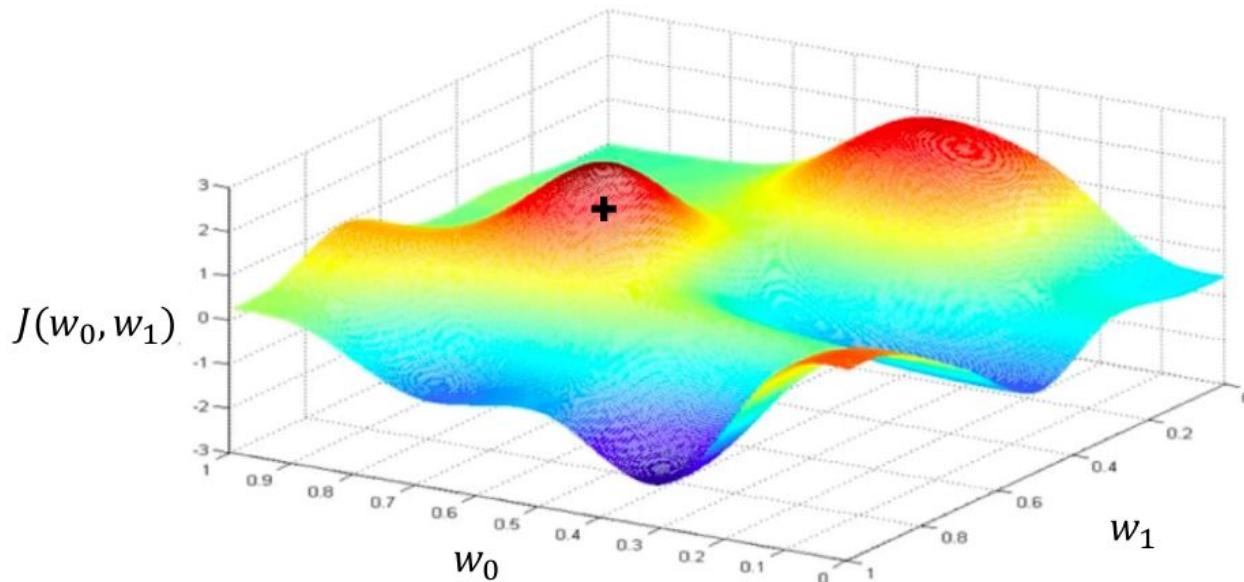
We want to find the network weights that **achieve the lowest loss**

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

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

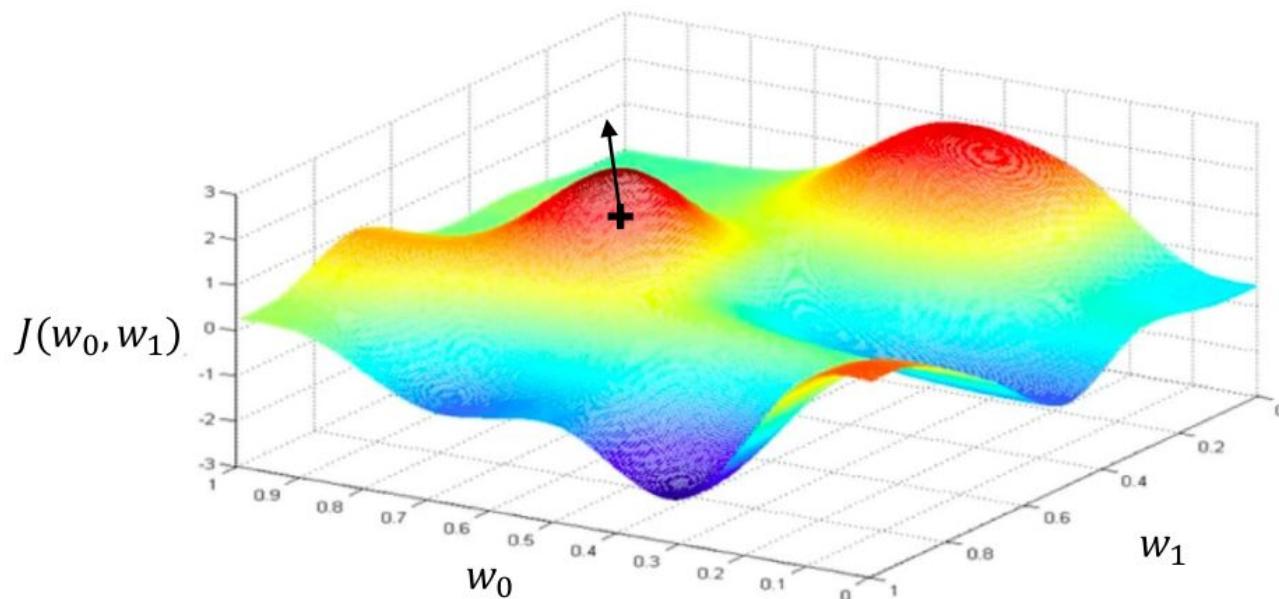
Loss Optimization

Randomly pick an initial (w_0, w_1)



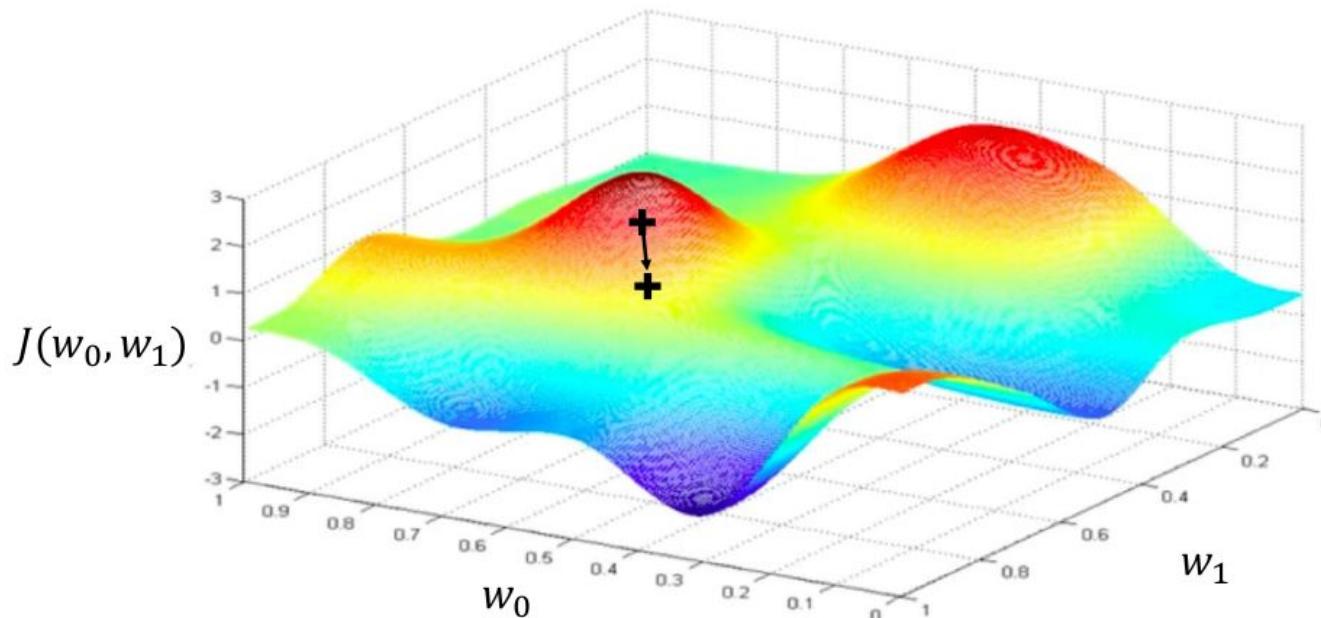
Loss Optimization

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



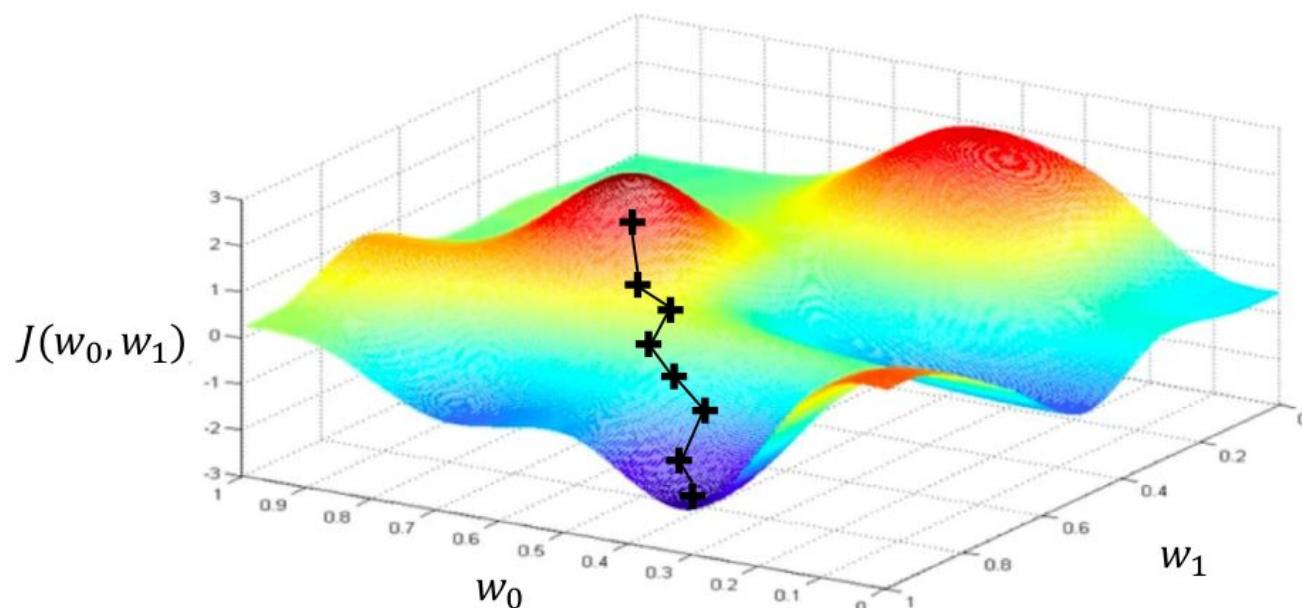
Loss Optimization

Take small step in opposite direction of gradient



Gradient Descent

Repeat until convergence



Gradient Descent

Algorithm

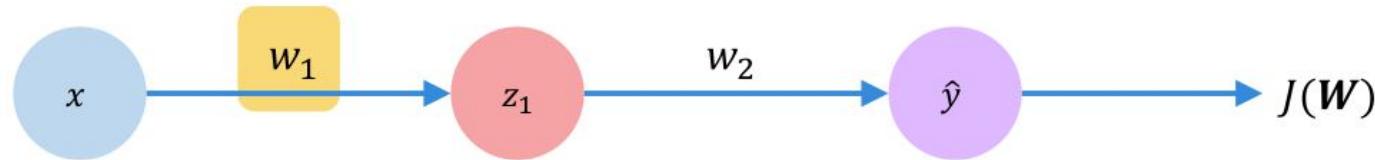
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
 3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
 4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
 5. Return weights

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



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underbrace{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{blue}}$$

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

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

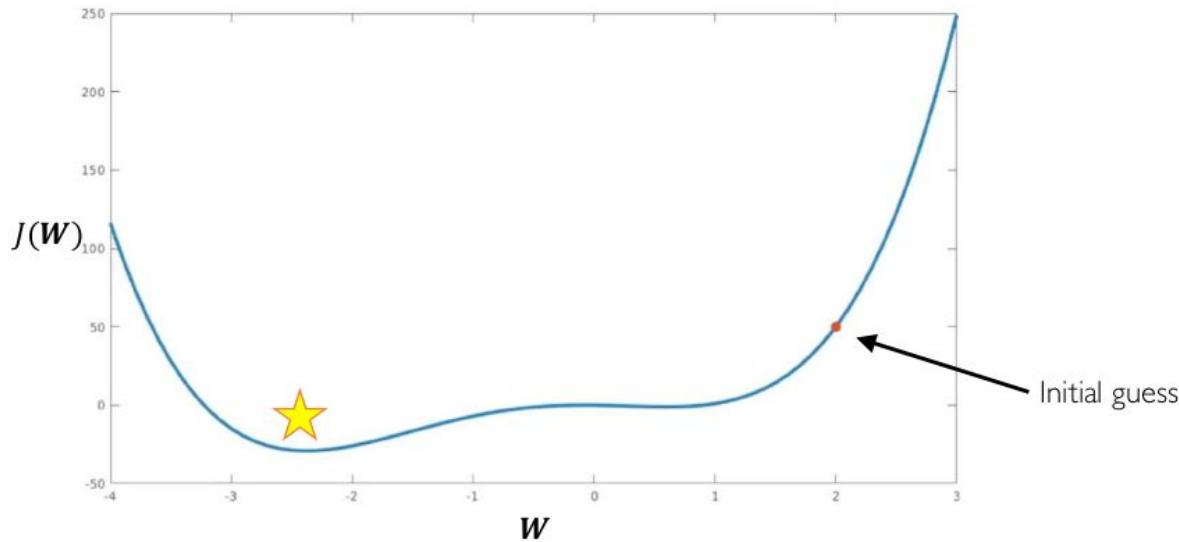


How can we set the
learning rate?

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

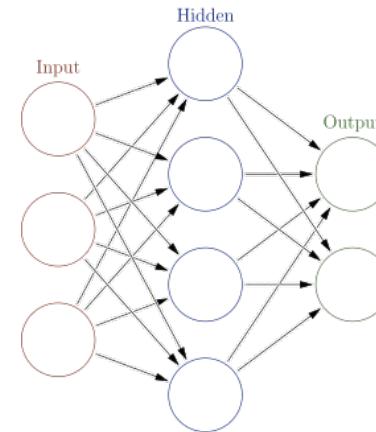
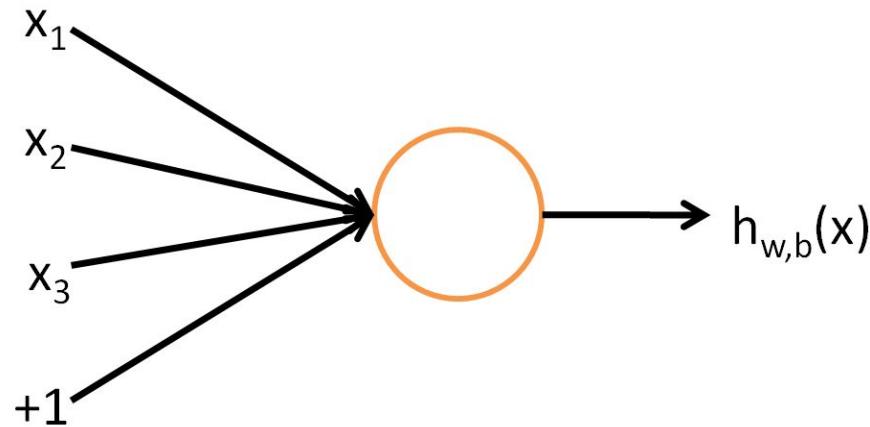
Adaptive Learning Rates

- 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?

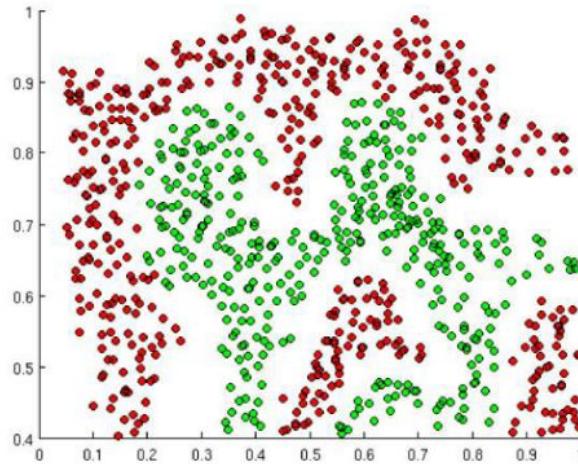


Linear hidden layer
is useless.

We actually need an AF to include non-linearities.

Importance of Activation Functions

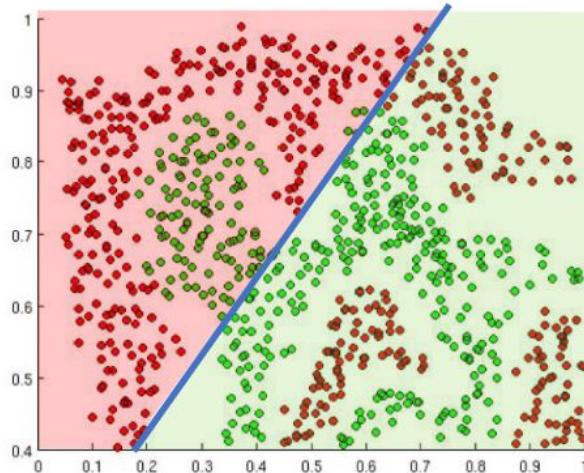
The purpose of activation functions is to *introduce non-linearities* into the network



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

Importance of Activation Functions

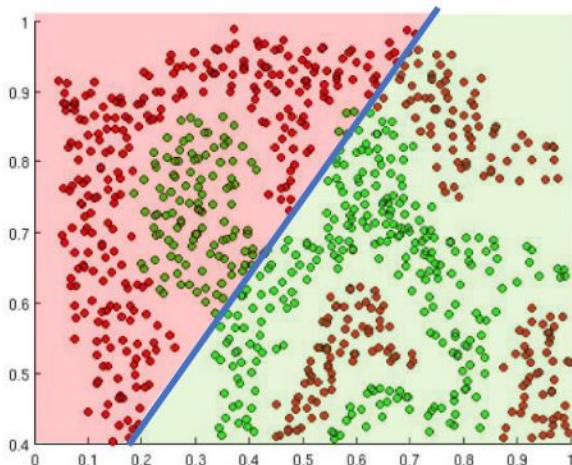
The purpose of activation functions is to **introduce non-linearities** into the network



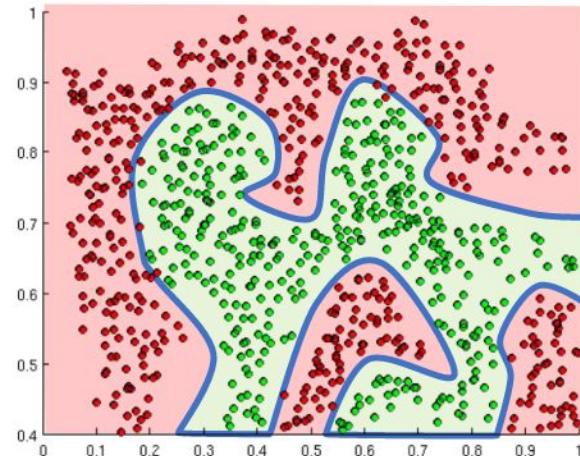
Linear Activation functions produce linear decisions no matter the network size

Importance of Activation Functions

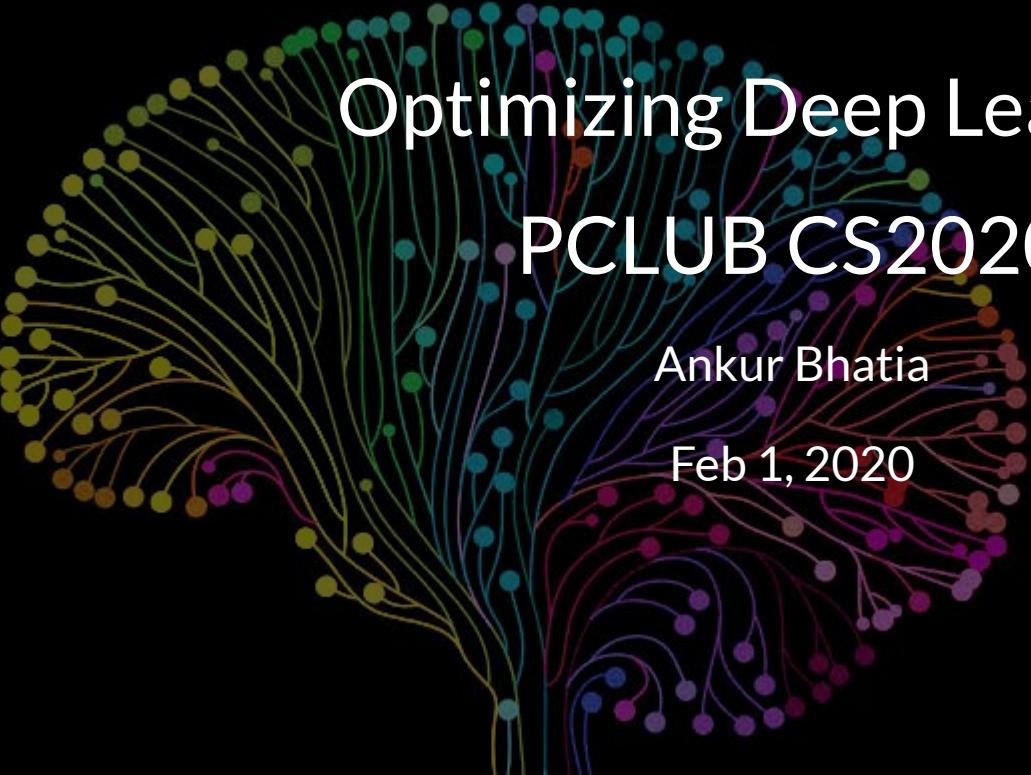
The purpose of activation functions is to *introduce non-linearities* into the network



Linear Activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions



A large, abstract neural network diagram is positioned on the left side of the slide. It consists of numerous small, colored circular nodes (representing neurons) connected by thin, curved lines (representing synapses). The nodes are primarily yellow, green, blue, purple, and red, and they are densely clustered in several distinct, radiating patterns across the dark background.

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 {
 for $t = 1, \dots, 5000$ {
 } } {
 } } {
 forward prop on $X^{[t]}$
 $Z^{[t]} = W^{[t]} X^{[t]} + b^{[t]}$
 $A^{[t]} = g^{[t]}(Z^{[t]})$
 \vdots
 $A^{[t]} = g^{[t]}(Z^{[t]})$
} } } {
 vectorized implementation
 (1000 ex)
}

Compute cost $J^{[t]} = \frac{1}{1000} \sum_{n=1}^L \ell(\hat{y}^{[t]}, y^{[t]}) + \frac{\lambda}{2 \cdot 1000} \sum_{j=1}^L \|w_j\|^2$

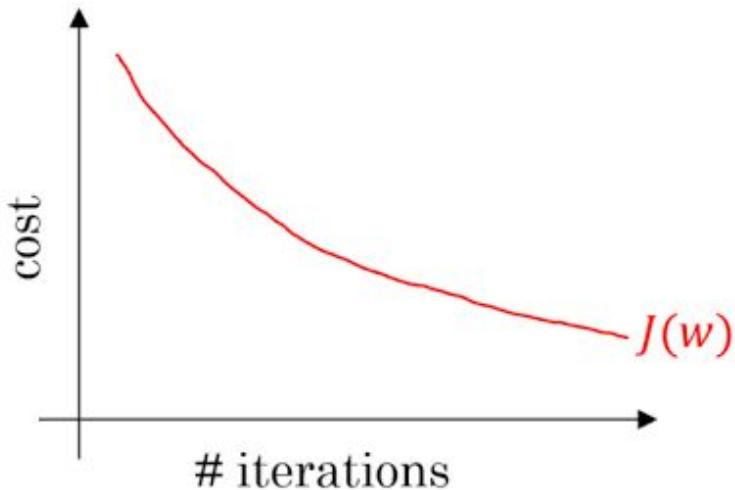
Back prop to compute gradients w.r.t $J^{[t]}$ (using $(X^{[t]}, y^{[t]})$)

$$w^{[t]} = w^{[t]} - \alpha \frac{dJ}{dw^{[t]}} , b^{[t]} = b^{[t]} - \alpha \frac{dJ}{db^{[t]}}$$

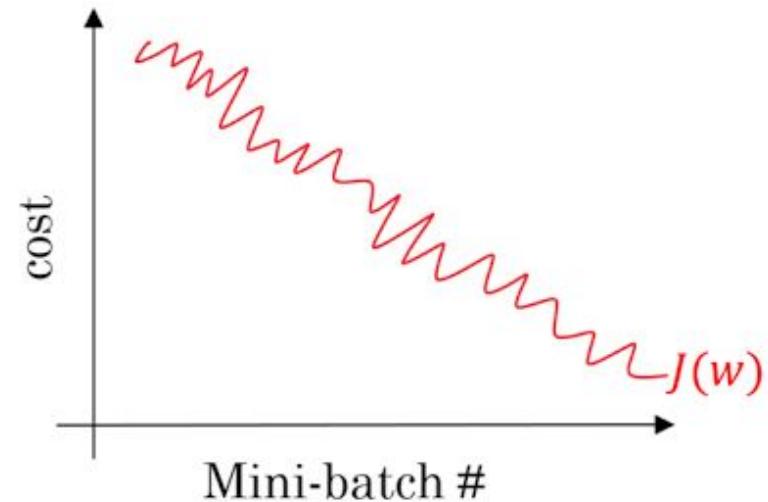
1 epoch (after full for loop completed).

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



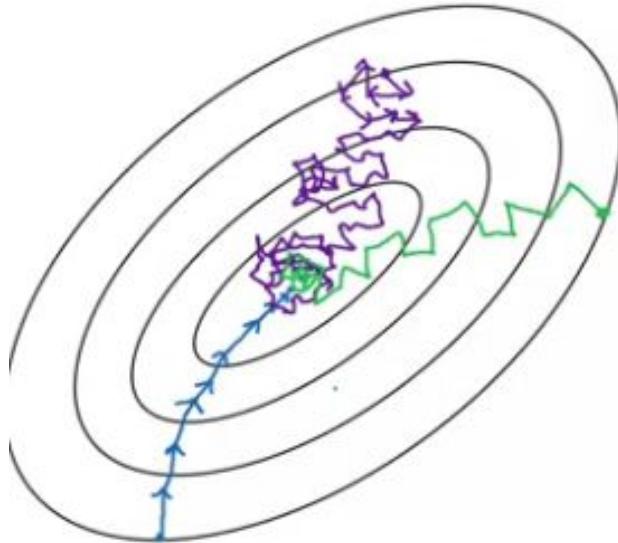
One Mini batch may be easy and other may be difficult to learn.

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

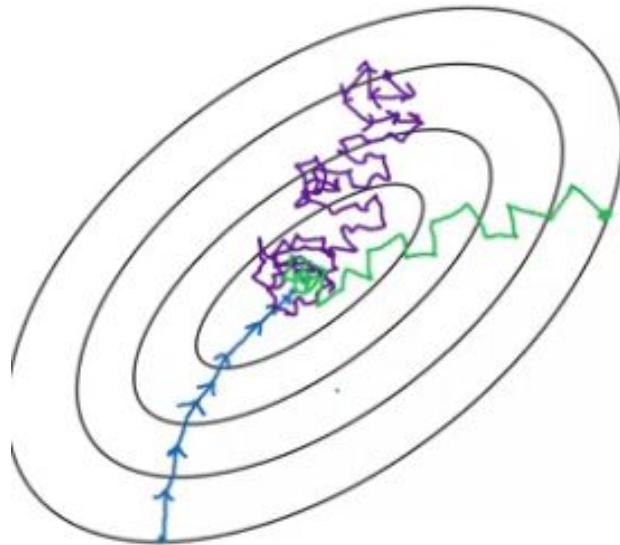


1. So What's the size to choose?.(bet 1 - m)
2. Stochastic GD loses speed of Vectorization.
3. In between mini-batch size gives:
 - a. Fast learning
 - b. Get good vectorization.
 - c. Makes progress without all data to process.
4. Batch-GD takes too long per iteration.

Mini Batch Gradient Descents

If $m = 1$; Stochastic Gradient Descent

If $m = m$; Batch Gradient Descent

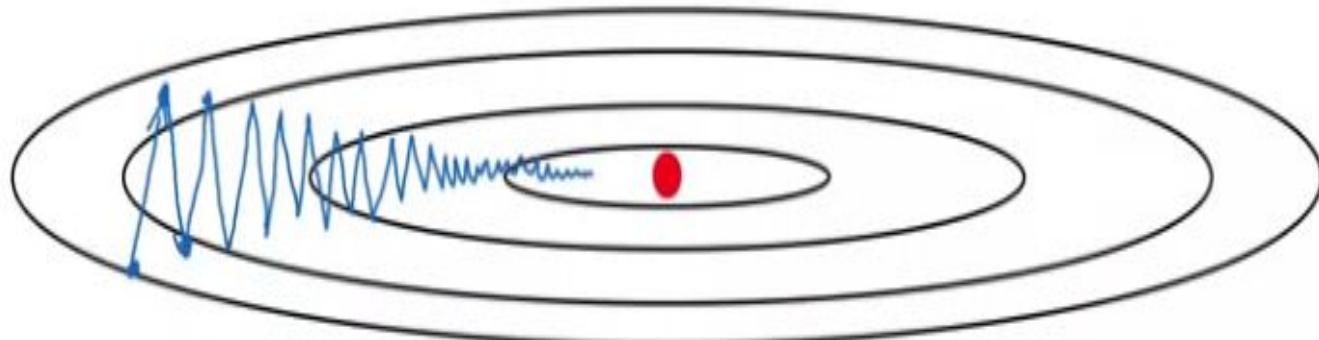


1. For small Trainx - $m = m$;
2. General **mini batch sizes** = 64, 128, 256...
3. mini - batch should fit in CPU - GPU m/m.

Optimization Algorithms faster than GD.

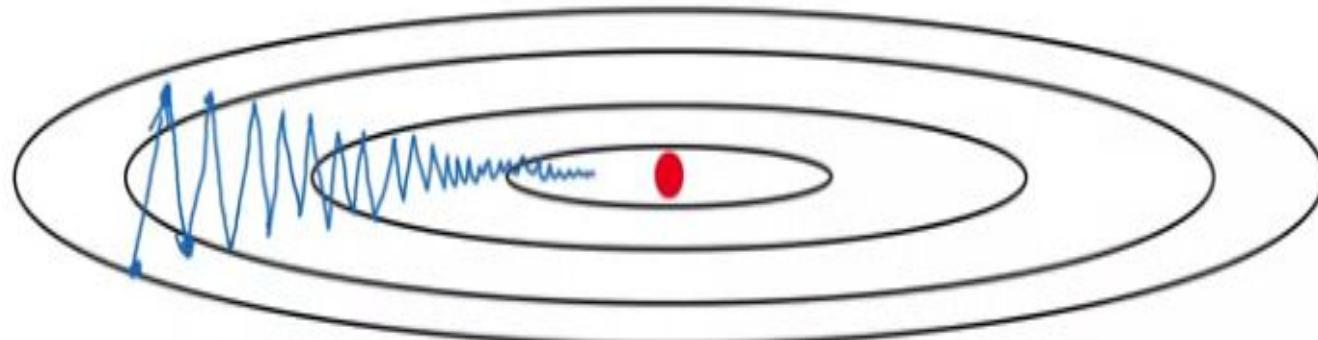
1. Gradient Descent with momentum

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:

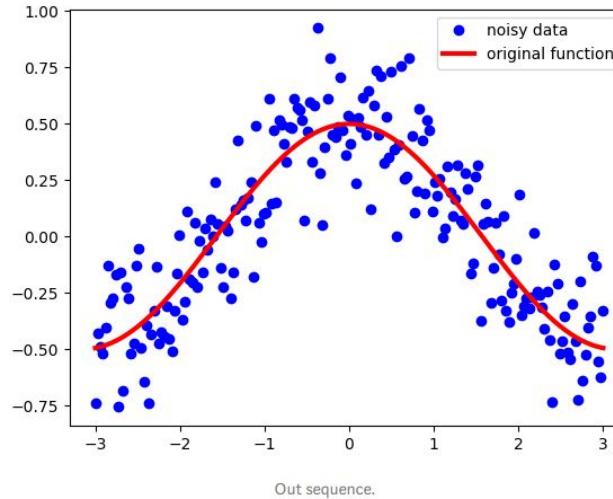
$$v_0 = 0$$

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

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

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

...



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:

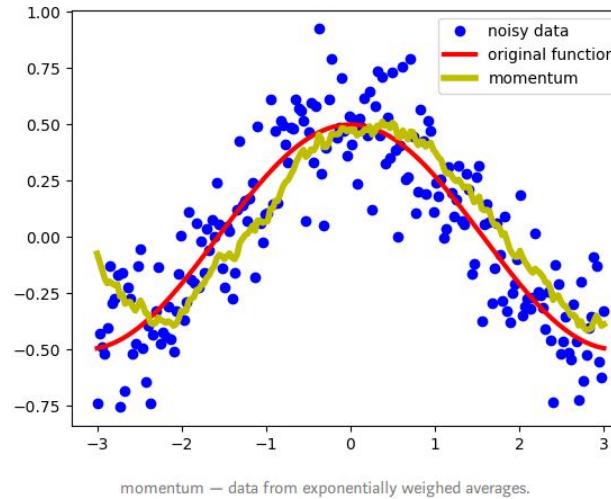
$$v_0 = 0$$

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

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

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

...



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

1. Gradient Descent with momentum : Exponentially weighted averages

Implementation \Rightarrow

$v_0 = 0$

Repeat { parameter
(can be weights/bias)

get current θ_t ;
 $v_t = \beta v_0 + (1 - \beta) \theta_t$

} ↓
moving averaged value according to β
parameters
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

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

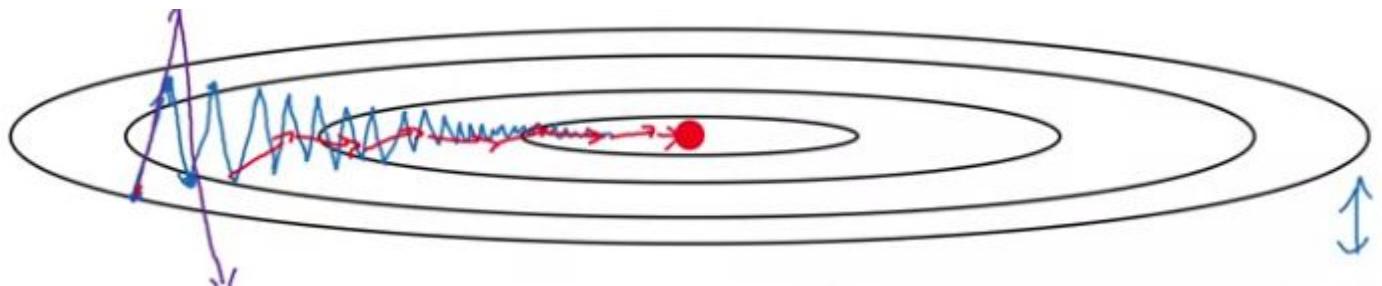
$$v_{db} = \beta v_{db} + (1 - \beta)db$$

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

Hyperparameters: α, β $\beta = 0.9$



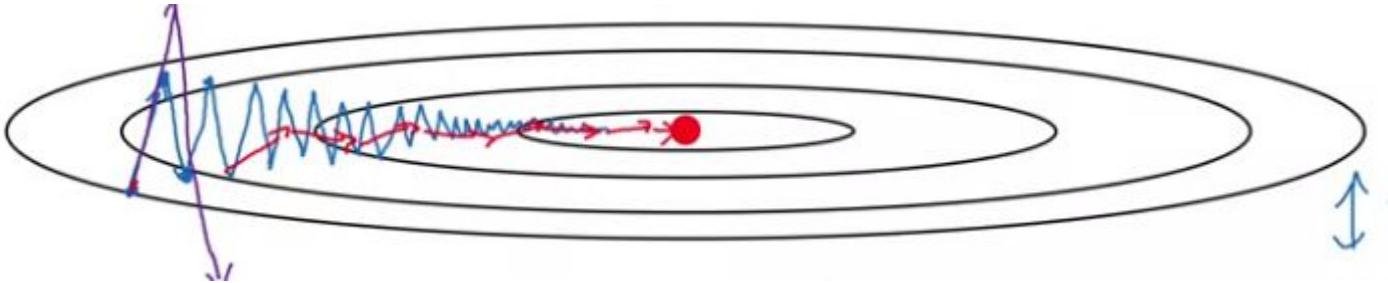
Gradient Descent with momentum



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

friction

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

momentum

velocity

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

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

Hyperparameters: α, β $\beta = 0.9$

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 V_{db} , V_{dw} ?

Now, Consider the intuition by: ball rolling

On iteration t :

Compute dW, db on the current mini-batch

friction

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

momentum

velocity

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

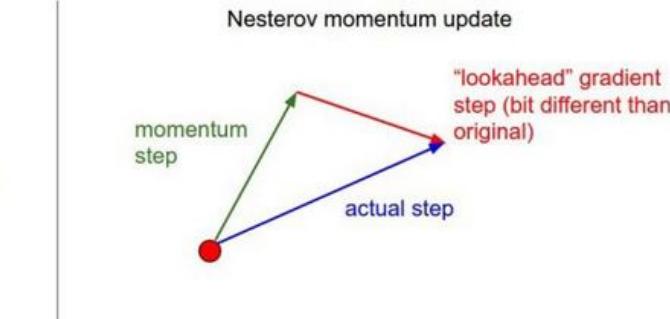
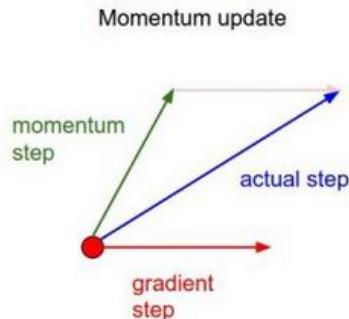
$$W = W - \alpha v_{dw}, \quad b = b - \alpha v_{db}$$

Hyperparameters: α, β $\beta = 0.9$

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

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.

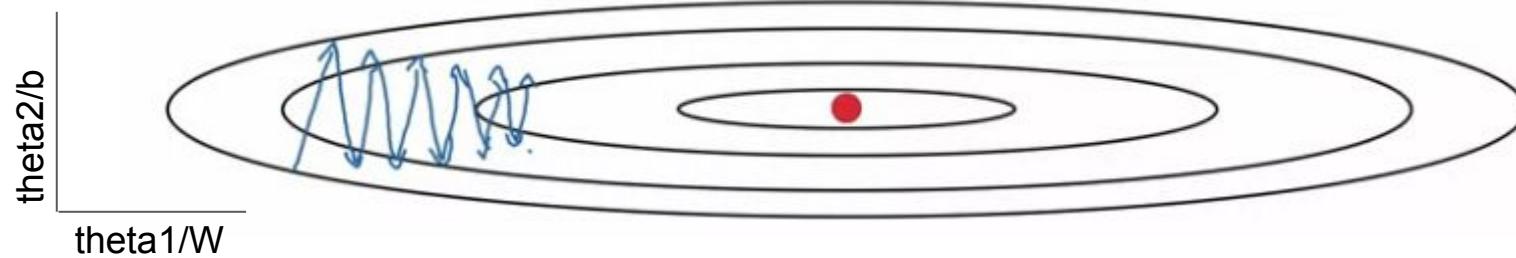


Source ([Stanford CS231n class](#))

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.

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

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

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

2. Gradient Descent with Root Mean Square Propagation:

Or mentioned Sdw
in place of Vdw

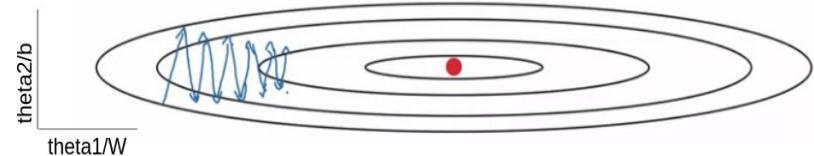
$$v_{dw} = \beta \cdot v_{dw} + (1 - \beta) \cdot dw^2 \quad \text{Elementwise}$$

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

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

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

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

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

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:

Or mentioned Sdw
in place of Vdw

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

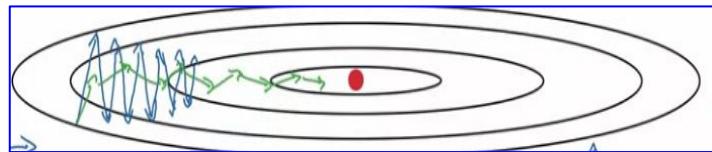
Vdw to be small

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

Vdb to be large

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

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



- 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

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

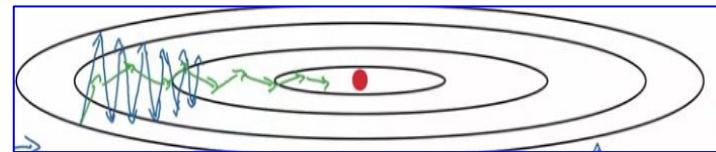
Vdw to be small

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

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 .

2. Gradient Descent with Root Mean Square Propagation:

Dimensions will be large.

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

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

classmate

Date _____

Page _____

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

On iteration t:

Compute d_w, d_b using minibatch $c_{1,2}$. β_1 momentum

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

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

~~Bias correction~~

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

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

$$w = w - \frac{d}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon}, b = b - \frac{d}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

Hyperparameters used

α = needs to be tuned - (Try own)

$$\beta_1 = 0.9 \quad (\text{dw})$$

$$\beta_2 = 0.999 \quad (\text{dw}^2)$$

$$\epsilon = 10^{-8}$$

(Adam paper)
(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>

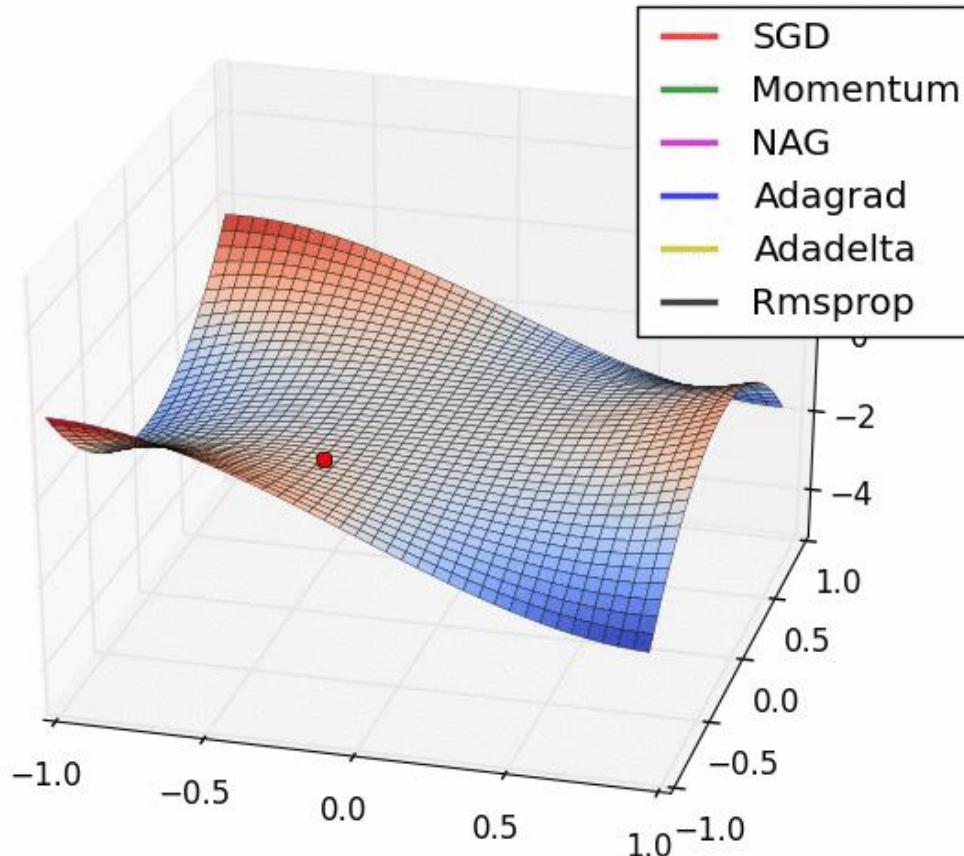
Algorithm 1: Adam, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_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