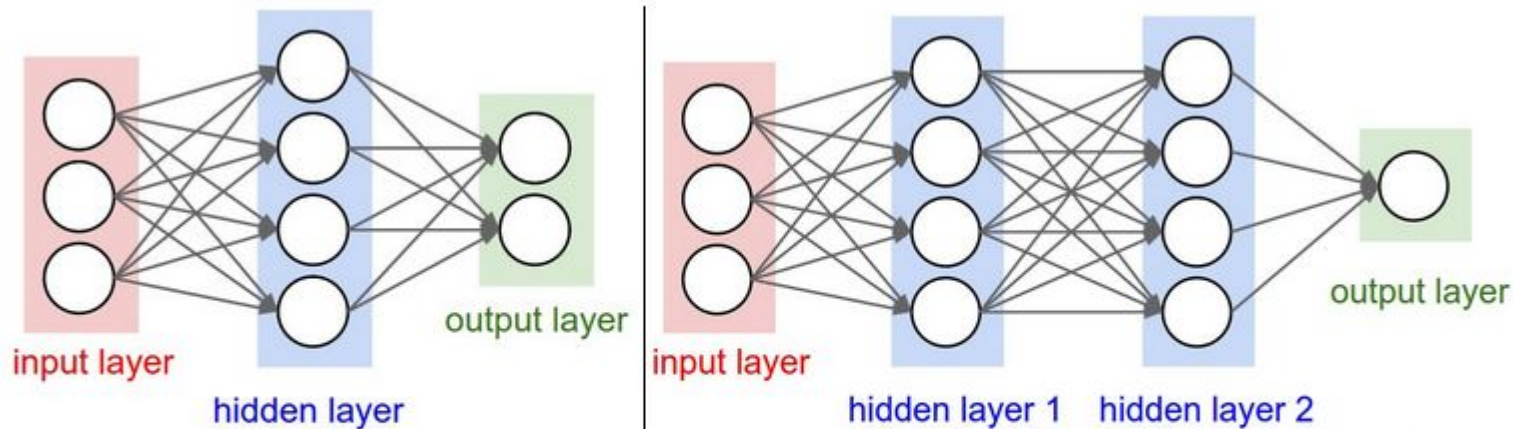


DNN_02

1. Neural Network Architecture
2. Forward Pass
3. Loss Function
4. Optimization
5. Gradient Descent
6. Learning Rate
7. Tips of Gradient Descent
8. Batch Gradient Descent
9. Stochastic Gradient Descent
10. Mini Batch Gradient Descent
11. Exponential Weighted Average
12. Gradient Descent with momentum
13. RMSprop
14. Adam
15. Learning Rate Decay
16. Local Optima

Neural Network Architectures



Left: A 2-layer Neural Network (one hidden layer of 4 neurons (or units) and one output layer with 2 neurons), and three inputs.

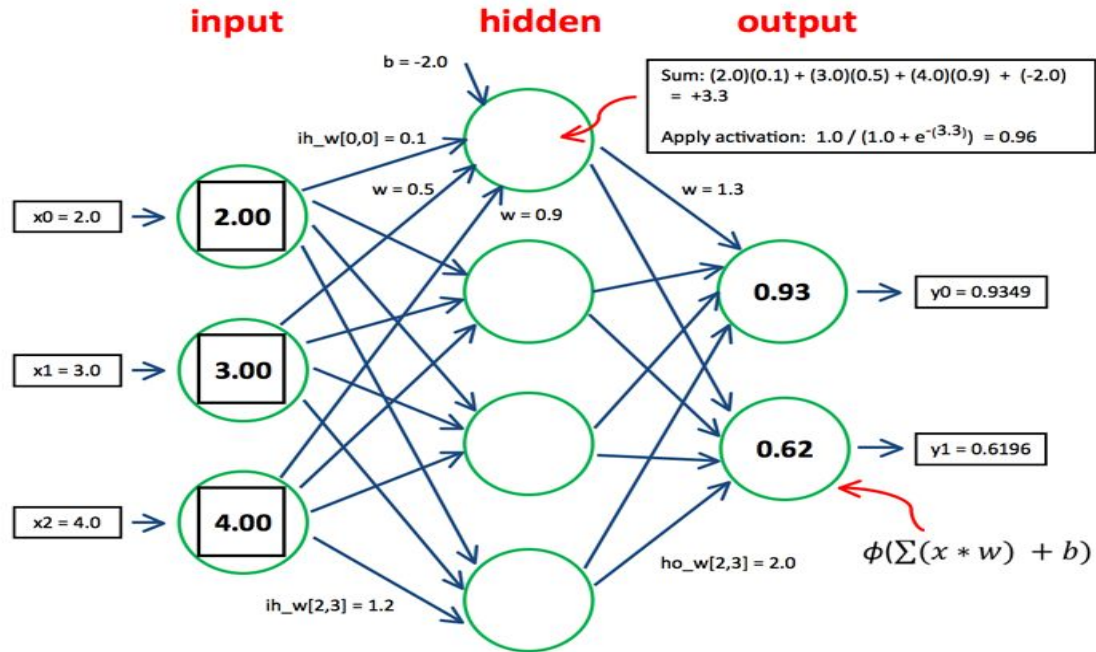
Right: A 3-layer neural network with three inputs, two hidden layers of 4 neurons each and one output layer. Notice that in both cases there are connections (synapses) between neurons across layers, but not within a layer.

Neural Network Architectures ..cont

- The first network (left) has $4 + 2 = 6$ neurons (not counting the inputs), $[3 \times 4] + [4 \times 2] = 20$ weights and $4 + 2 = 6$ biases, for a total of 26 learnable parameters.
- The second network (right) has $4 + 4 + 1 = 9$ neurons, $[3 \times 4] + [4 \times 4] + [4 \times 1] = 12 + 16 + 4 = 32$ weights and $4 + 4 + 1 = 9$ biases, for a total of 41 learnable parameters.

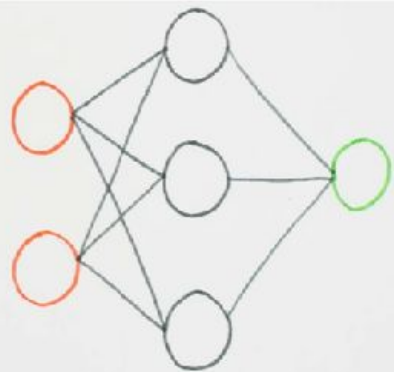
The forward pass of a fully-connected layer corresponds to one matrix multiplication followed by a bias offset and an activation function.

Forward Pass



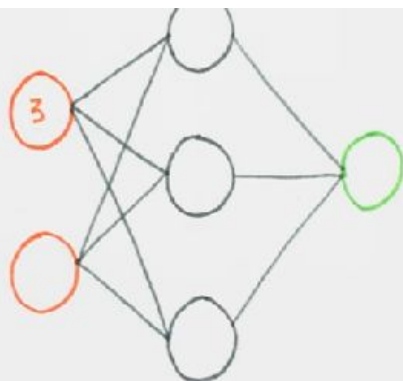
1

$$\begin{bmatrix} 3 & 5 \\ 5 & 1 \\ 10 & 2 \end{bmatrix}$$



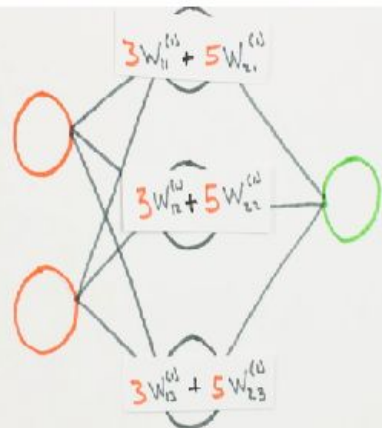
2

$$\begin{bmatrix} 3 & 5 \\ 5 & 1 \\ 10 & 2 \end{bmatrix}$$



3

$$\begin{bmatrix} 3 & 5 \\ 5 & 1 \\ 10 & 2 \end{bmatrix}$$



4

$$\begin{bmatrix} 3 & 5 \\ 5 & 1 \\ 10 & 2 \end{bmatrix} \begin{bmatrix} W_{11}^{(1)} & W_{12}^{(1)} & W_{13}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} & W_{23}^{(1)} \end{bmatrix} = \begin{bmatrix} 3W_{11}^{(1)} + 5W_{21}^{(1)} & 3W_{12}^{(1)} + 5W_{22}^{(1)} & 3W_{13}^{(1)} + 5W_{23}^{(1)} \\ 5W_{11}^{(1)} + 1W_{21}^{(1)} & 5W_{12}^{(1)} + 1W_{22}^{(1)} & 5W_{13}^{(1)} + 1W_{23}^{(1)} \\ 10W_{11}^{(1)} + 2W_{21}^{(1)} & 10W_{12}^{(1)} + 2W_{22}^{(1)} & 10W_{13}^{(1)} + 2W_{23}^{(1)} \end{bmatrix}$$

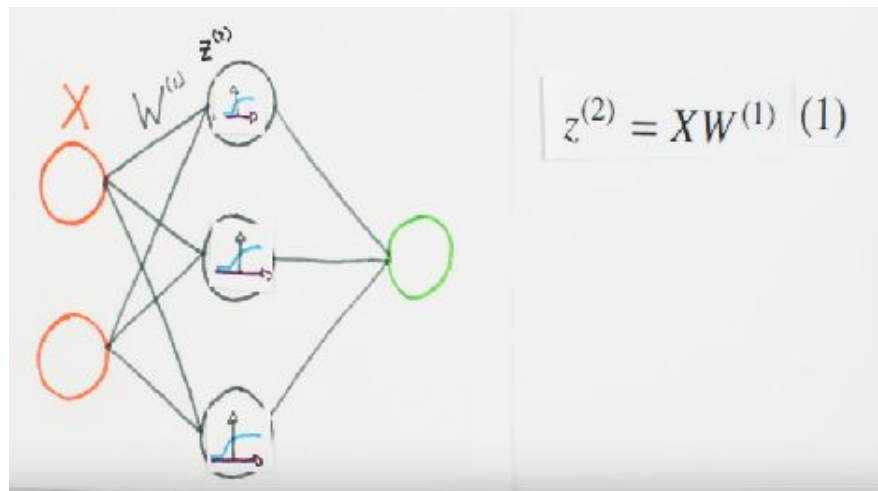
5

X

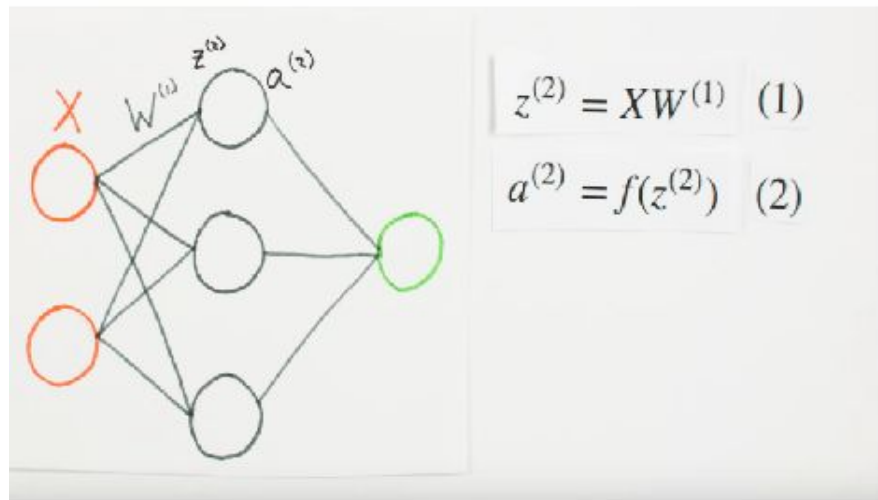
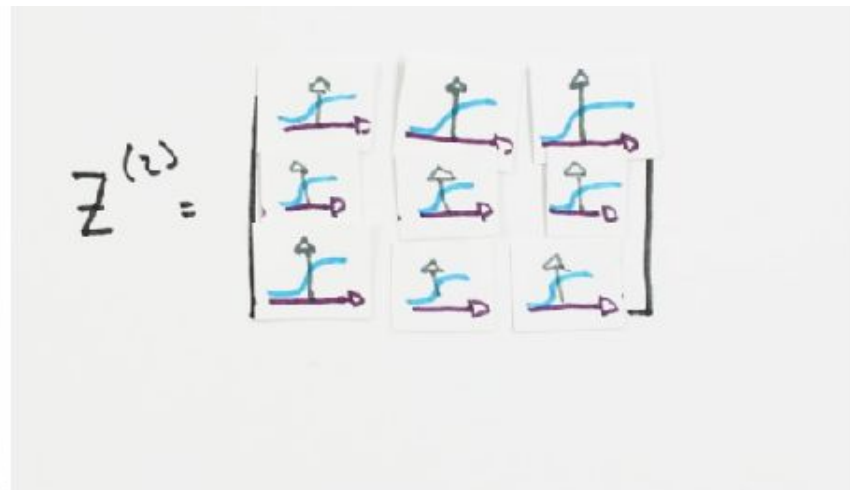
 $W^{(1)}$

=

 $Z^{(2)}$

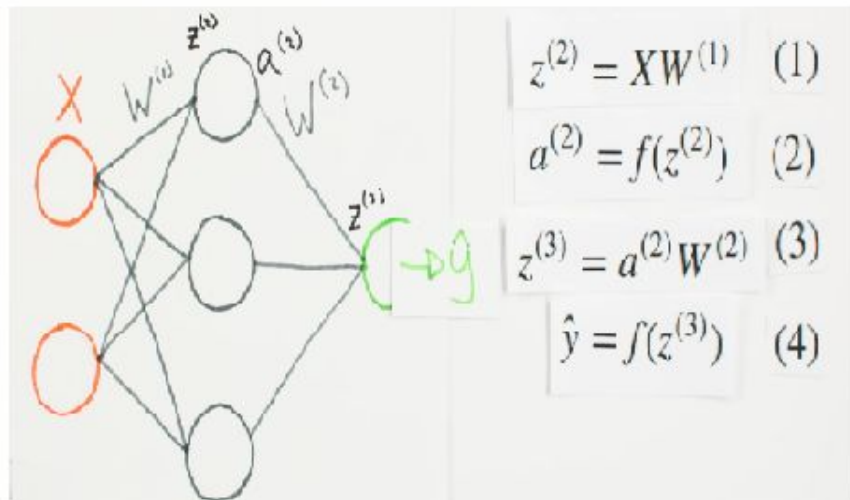


$$z^{(2)} = XW^{(1)} \quad (1)$$



$$z^{(2)} = XW^{(1)} \quad (1)$$

$$a^{(2)} = f(z^{(2)}) \quad (2)$$



$$z^{(2)} = XW^{(1)} \quad (1)$$

$$a^{(2)} = f(z^{(2)}) \quad (2)$$

$$z^{(3)} = a^{(2)}W^{(2)} \quad (3)$$

$$\hat{y} = f(z^{(3)}) \quad (4)$$


```
In [1]: class Neural_Network(object):
def __init__(self):
    #Define Hyperparameters
    self.inputLayerSize = 2
    self.outputLayerSize = 1
    self.hiddenLayerSize = 3

    #Weights (Parameters)
    self.W1 = np.random.randn(self.inputLayerSize, \
                                self.hiddenLayerSize)
    self.W2 = np.random.randn(self.hiddenLayerSize, \
                                self.outputLayerSize)

def forward(self, X):
    #Propagate inputs though network
    self.z2 = np.dot(X, self.W1)
    self.a2 = self.sigmoid(self.z2)
    self.z3 = np.dot(self.a2, self.W2)
    yHat = self.sigmoid(self.z3)
    return yHat

def sigmoid(self, z):
    #Apply sigmoid activation function to scalar, vector, or
    return 1/(1+np.exp(-z))
```

$$z^{(2)} = XW^{(1)} \quad (1)$$

$$a^{(2)} = f(z^{(2)}) \quad (2)$$

$$z^{(3)} = a^{(2)}W^{(2)} \quad (3)$$

$$\hat{y} = f(z^{(3)}) \quad (4)$$

Loss/Cost/Objective Function

<https://stats.stackexchange.com/questions/154879/a-list-of-cost-functions-used-in-neural-networks-alongside-applications>

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

Mean Square Error

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Diagram illustrating the Mean Square Error (MSE) formula:

- 1 : average over all results
- N : average over all results
- $\sum_{i=1}^N$: Summation over all data points
- y_i : true y
- \hat{y}_i : estimate of y
- 2 : makes result quadratic

Mostly used in linear regression.

Cross Entropy

$$H(x) = \sum_{i=1}^N \overset{\text{true label}}{p}(x) \log \underset{\text{estimate}}{q}(x)$$

Mostly used in classification problem / logistic regression.

Methods of Optimization

1. Random Search
2. Random Local Search
3. Following the gradient

Gradient Descent

Gradient descent is an optimization algorithm used to find the values of parameters of a function (f) that minimizes a cost function (cost).

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2.$$

Simplified

Hypothesis:

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

$$h_{\theta}(x) = \theta_1 x$$

Parameters:

$$\theta_0, \theta_1$$

$$\theta_1$$

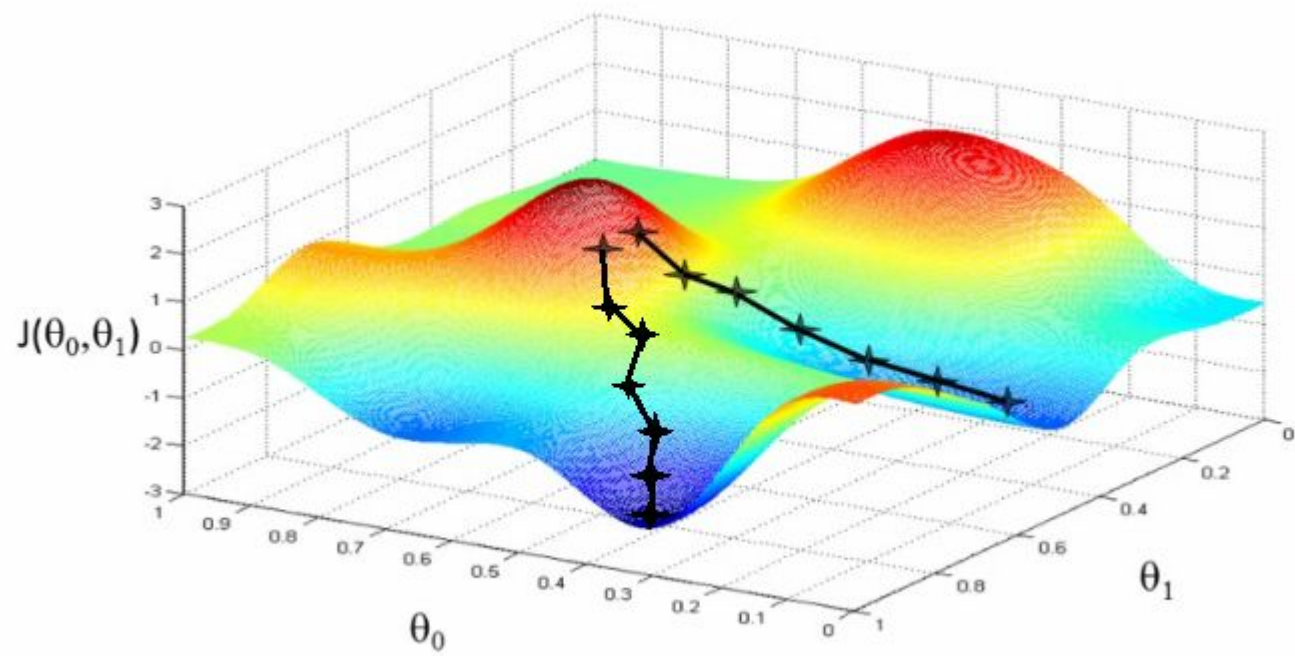
Cost Function:

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

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

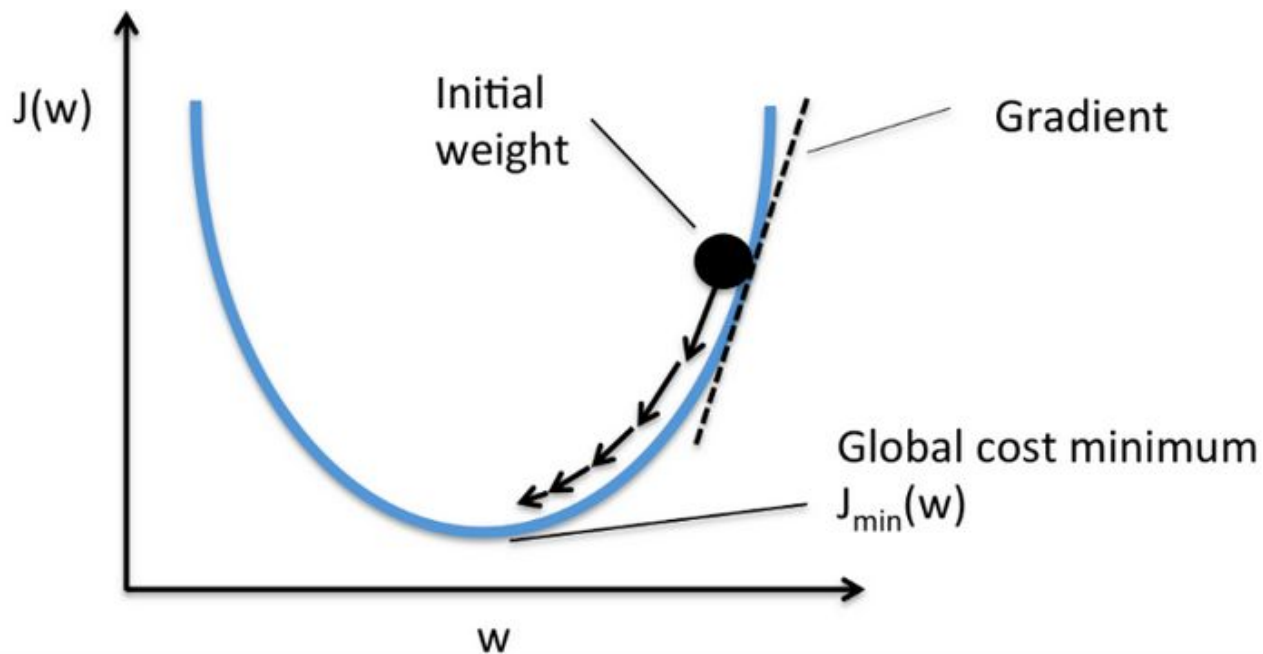
Goal: minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1

minimize $J(\theta_1)$
 θ_1



$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

Gradient Descent

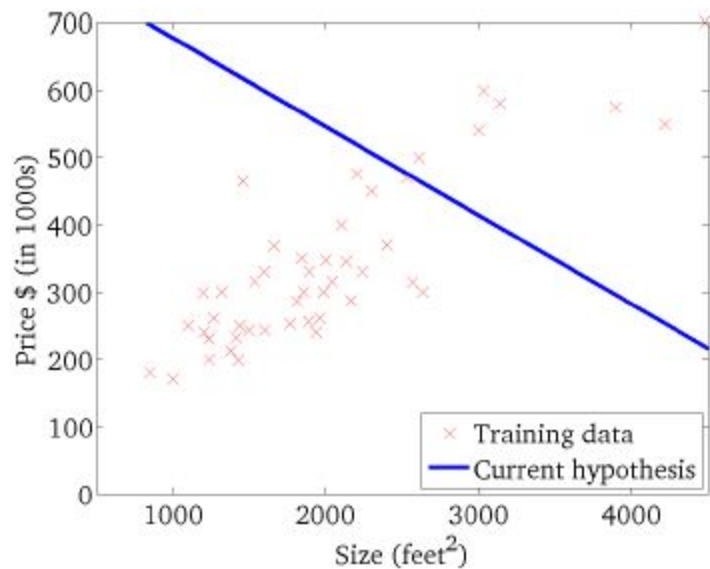


Gradient Descent Algorithm

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

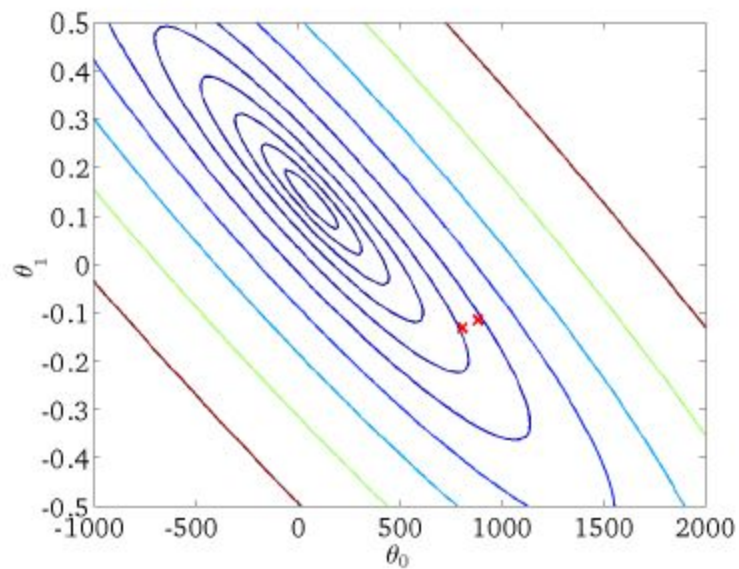
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



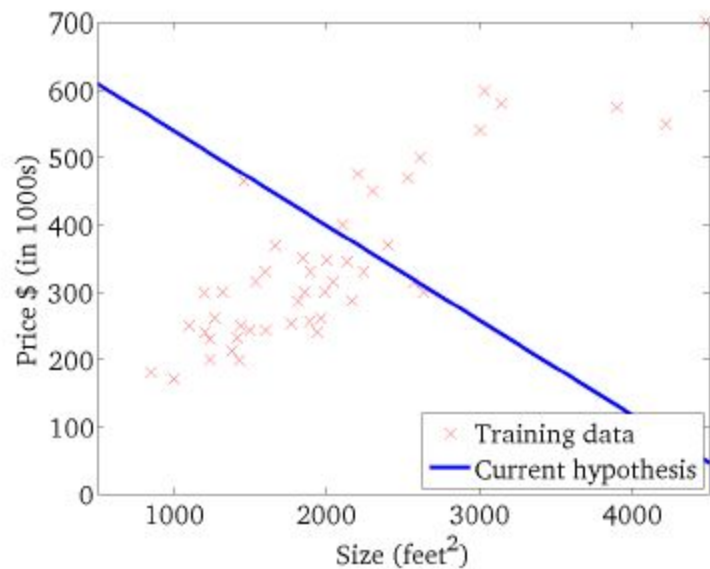
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



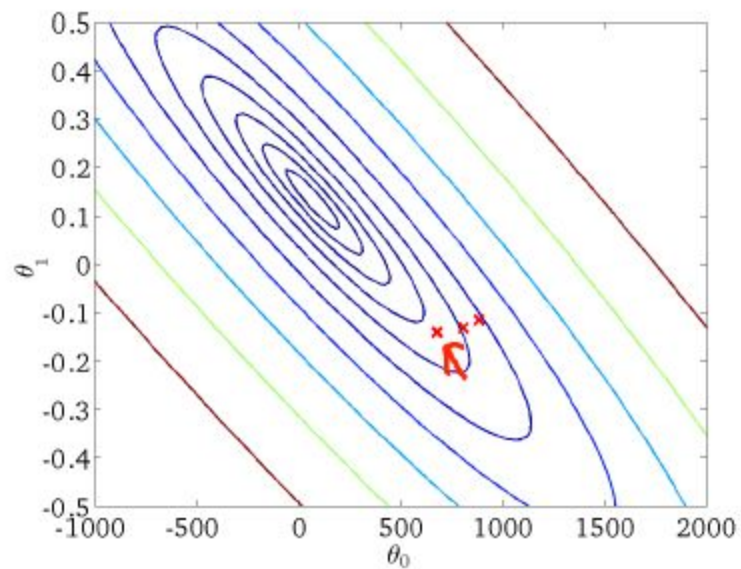
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



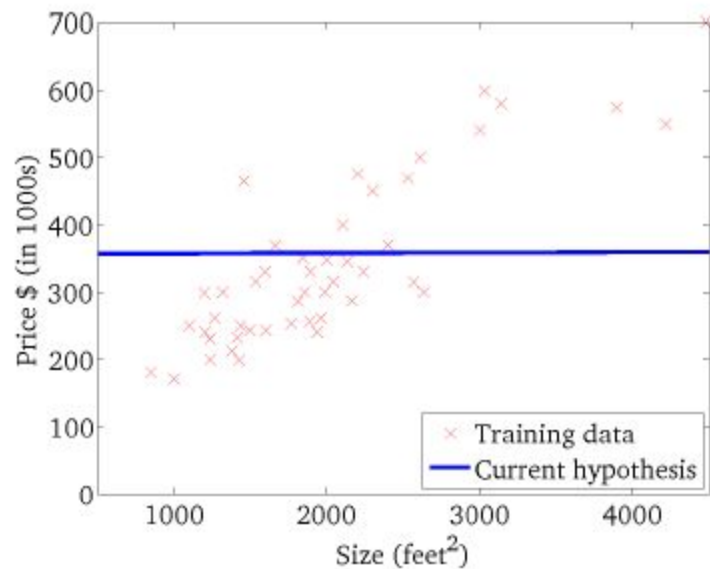
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



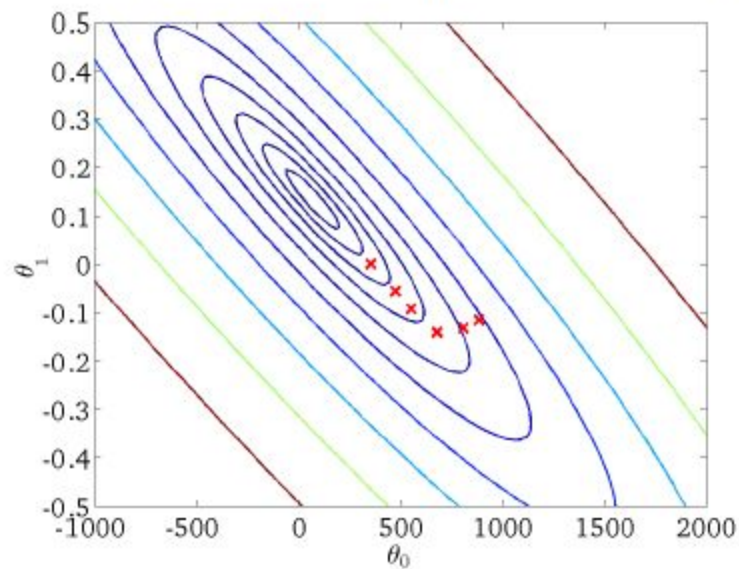
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



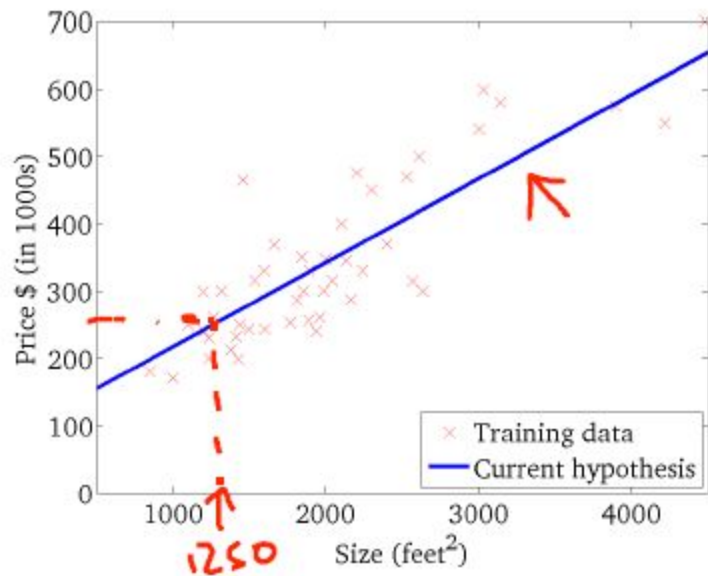
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



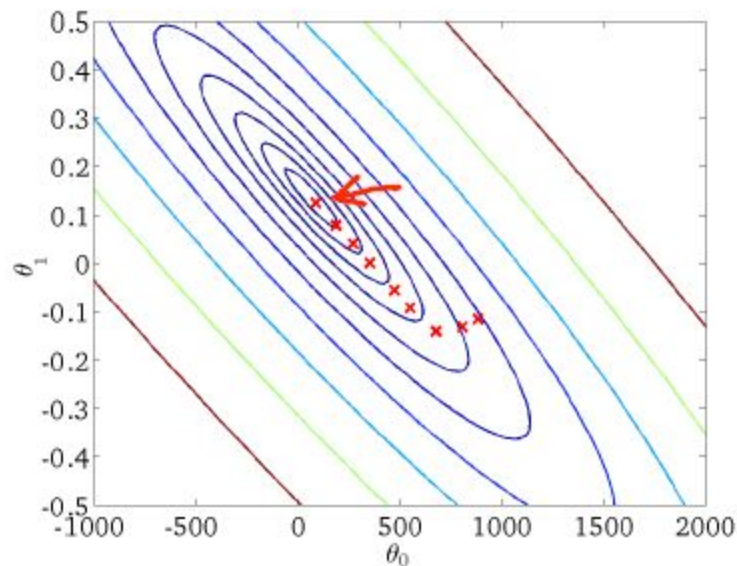
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



Implementation

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis} = 1, \text{keepdims} = \text{True})$$

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

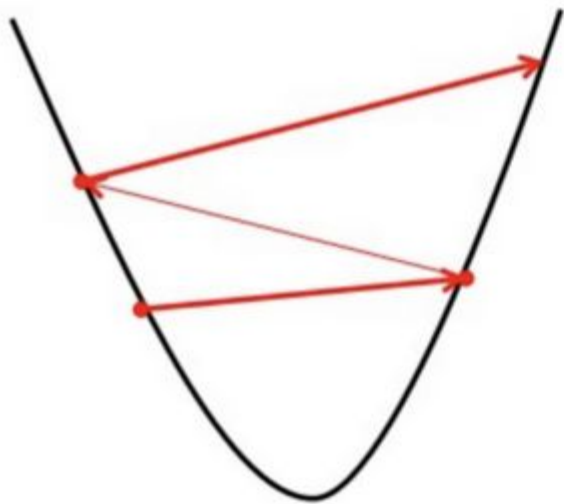
$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis} = 1, \text{keepdims} = \text{True})$$

Importance of Learning Rate

How big the steps are that Gradient Descent takes into the direction of the local minimum are determined by the so-called learning rate. It determines how fast or slow we will move towards the optimal weights.

Big learning rate

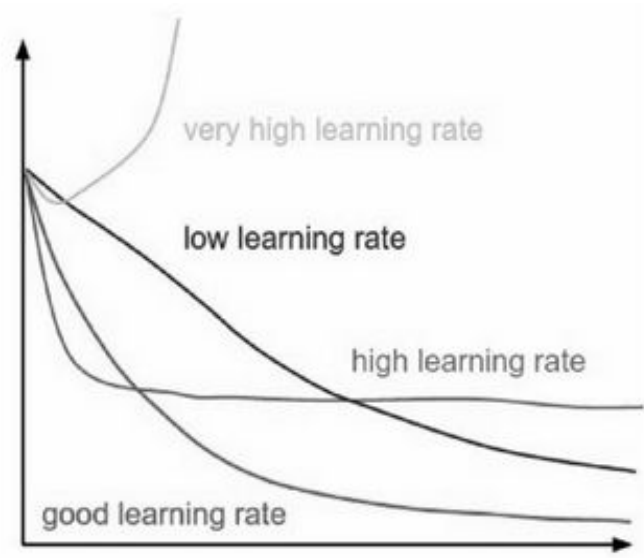
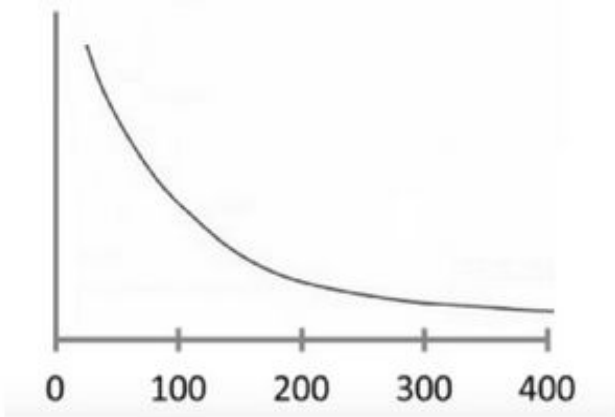


Small learning rate



Tips for Gradient Descent

- **Plot Cost versus Time:** Collect and plot the cost values calculated by the algorithm each iteration. The expectation for a well performing gradient descent run is a decrease in cost each iteration. If it does not decrease, try reducing your learning rate.
- **Learning Rate:** The learning rate value is a small real value such as 0.1, 0.001 or 0.0001. Try different values for your problem and see which works best.
- **Rescale Inputs:** The algorithm will reach the minimum cost faster if the shape of the cost function is not skewed and distorted. You can achieved this by rescaling all of the input variables (X) to the same range, such as $[0, 1]$ or $[-1, 1]$.



Number of Iteration Vs Cost

Batch Gradient Descent

- Also known as vanilla gradient descent
- Calculates the error for each example within the training dataset, but only after all training examples have been evaluated, the model gets updated
- computational efficient
- it produces a stable error gradient and a stable convergence
- Stable error gradient may lead to local minima
- It also requires that the entire training dataset is in memory and available to the algorithm
- Use this is training set have approx 2000-3000 input data

Stochastic Gradient Descent

- Calculates the error for each example within the training dataset, and update the model for each example. This means that it updates the parameters for each training example, one by one.
- This can make SGD faster than Batch Gradient Descent, depending on the problem.
- Frequent updates are more computationally expensive.
- The frequency of those updates can also result in noisy gradients, which may cause the error rate to jump around, instead of slowly decreasing. This helps in solving the problem of convergence at local minima.

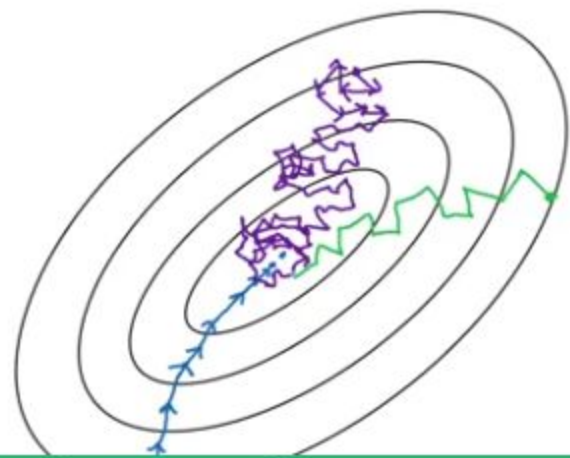
Mini Batch Gradient Descent

- Mini-batch Gradient Descent is the go-to method since it's a combination of the concepts of SGD and Batch Gradient Descent.
- It simply splits the training dataset into small batches and performs an update for each of these batches.
- Common mini-batch sizes range between 64, 256, 512

→ If mini-batch size = m : Batch gradient descent. $(X^{(1)}, Y^{(1)}) = (X, Y)$.

→ If mini-batch size = 1 : Stochastic gradient descent. Every example is its own mini-batch.
 $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$ mini-batch.

In practice: Somewhere in-between 1 and m



Stochastic
gradient
descent

}

Large spread
from vectorization

In-between
(mini-batch size
not too big/small)

}

Fastest learning.

- Vectorization.
(~ 1000)
- Make passes without

Batch
gradient descent
(mini-batch size = m)

}

Too long
per iteration

Exponential Weighted Average

Temperature in London

$$\theta_1 = 40^\circ\text{F} \quad 4^\circ\text{C} \leftarrow$$

$$\theta_2 = 49^\circ\text{F} \quad 9^\circ\text{C}$$

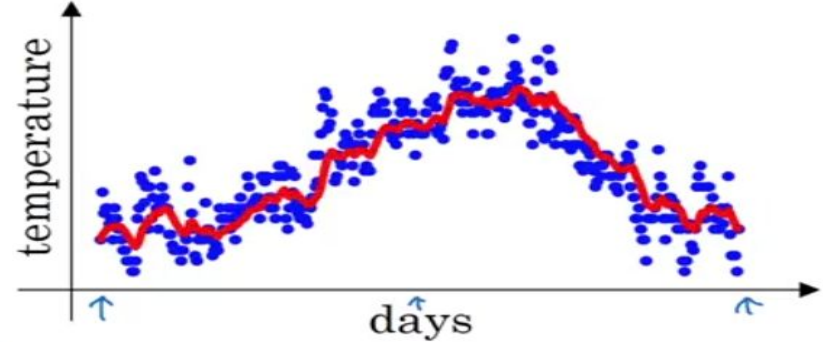
$$\theta_3 = 45^\circ\text{F} \quad \vdots$$

\vdots

$$\theta_{180} = 60^\circ\text{F} \quad 15^\circ\text{C}$$

$$\theta_{181} = 56^\circ\text{F} \quad \vdots$$

\vdots



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

$$V_3 = 0.9 V_2 + 0.1 \theta_3$$

\vdots

$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

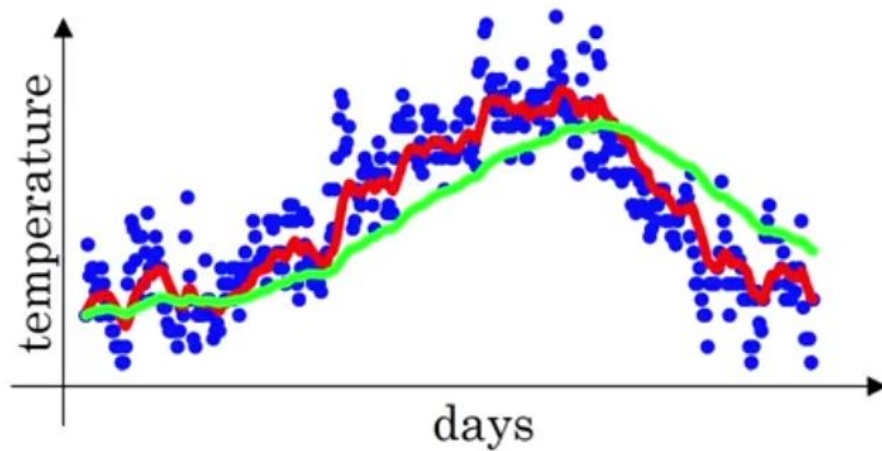
$$V_t = \beta V_{t-1} + (1-\beta) \Theta_t$$

$\beta = 0.9$: ≈ 10 days' temper.

$\beta = 0.98$: ≈ 50 days

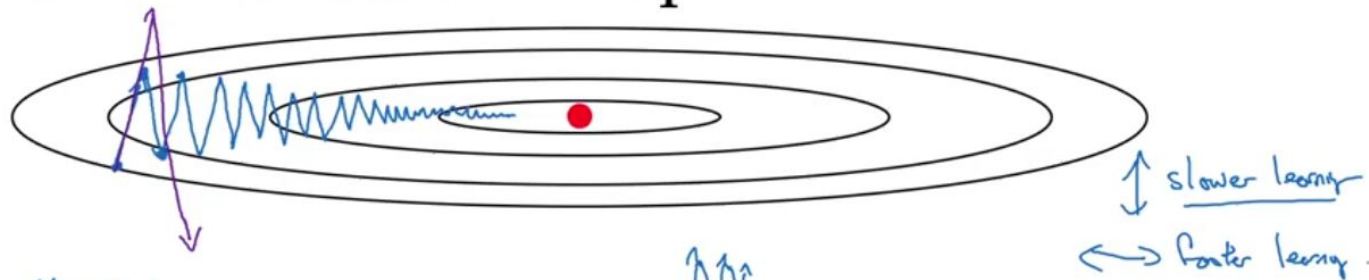
V_t is approximately
average over
 $\approx \frac{1}{1-\beta}$ days'
temperature.

$$\frac{1}{1-0.98} = 50$$



Gradient Descent with Momentum

Gradient descent example



Momentum:

On iteration t :

Compute $\Delta W, \Delta b$ on current mini-batch.

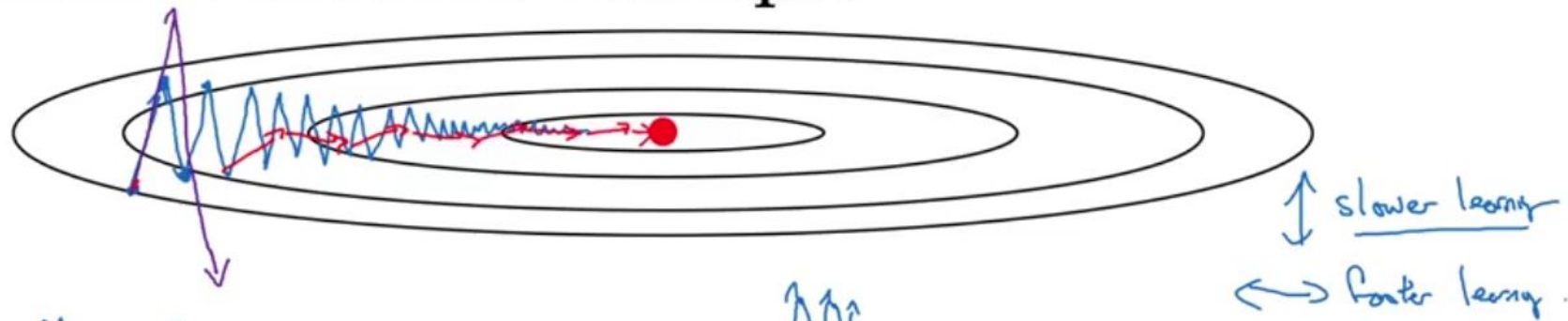
$$V_{\Delta W} = \beta V_{\Delta W} + (1-\beta) \Delta W$$

$$V_{\Delta b} = \beta V_{\Delta b} + (1-\beta) \Delta b$$

$$W := W - \alpha V_{\Delta W}, \quad b := b - \alpha V_{\Delta b}$$

$$V_{\theta} = \beta V_{\theta} + (1-\beta) \theta_c$$

Gradient descent example



Momentum:

On iteration t :

Compute $\underline{dW}, \underline{db}$ on current mini-batch.

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

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



$$V_{\theta} = \beta V_{\theta} + (1-\beta) \theta_e$$

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

Implementation details

$$v_{dw} = 0, \quad v_{db} = 0$$

On iteration t :

Compute dW, db on the current mini-batch

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

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

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



Hyperparameters: α, β

$$\underline{\beta = 0.9}$$

Implementation details

$$v_{dw} = 0, \quad v_{db} = 0$$

On iteration t :

Compute dW, db on the current mini-batch

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

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

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

$$\underline{v_{dw}} = \beta v_{dw} + dW \leftarrow$$

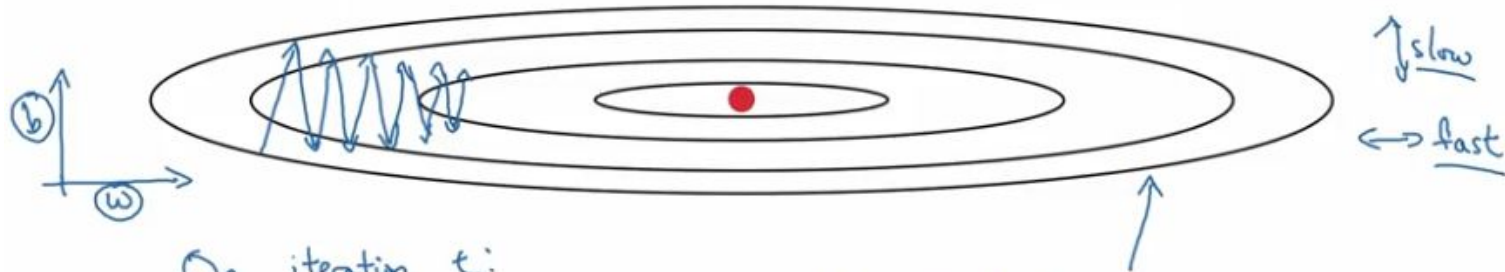
$$\frac{v_{dw}}{1 - \beta^t}$$

Hyperparameters: α, β

$$\underline{\beta = 0.9}$$

RMSprop

RMSprop



On iteration t :

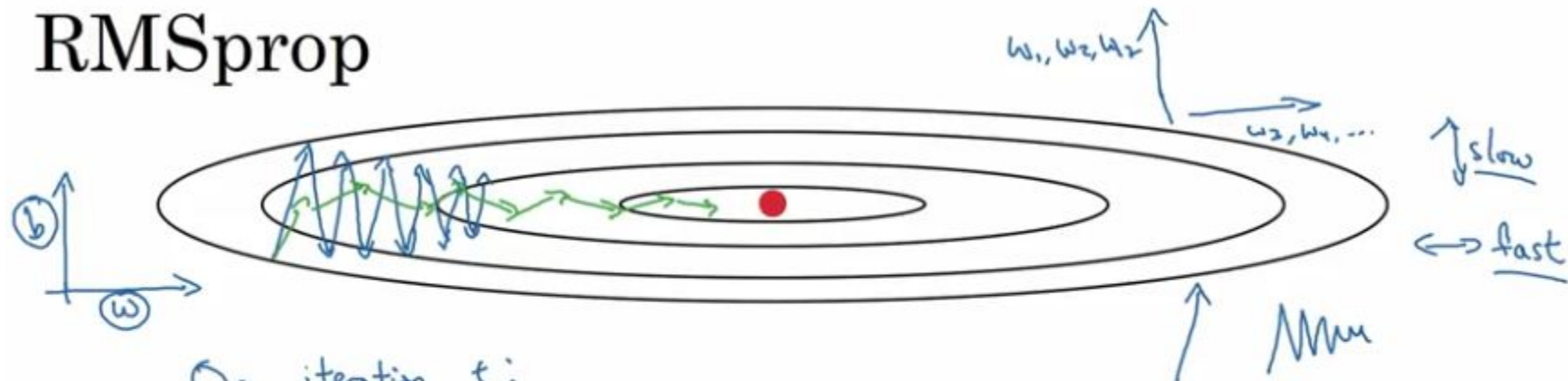
Compute $\Delta w, \Delta b$ on current mini-batch

$$\underline{S_{\Delta w}} = \beta \underline{S_{\Delta w}} + (1-\beta) \underline{\Delta w^2} \leftarrow \text{small}$$

$$\underline{S_{\Delta b}} = \beta \underline{S_{\Delta b}} + (1-\beta) \underline{\Delta b^2} \leftarrow \text{large}$$

$$w := w - \alpha \frac{\Delta w}{\sqrt{\underline{S_{\Delta w}}}} \leftarrow \quad b := b - \alpha \frac{\Delta b}{\sqrt{\underline{S_{\Delta b}}}} \leftarrow$$

RMSprop



On iteration t :

Compute dw, db on current mini-batch

$$\underline{S_{dw}} = \beta \underline{S_{dw}} + (1-\beta) \underline{dw^2} \leftarrow \text{small}$$

$$\rightarrow \underline{S_{db}} = \beta \underline{S_{db}} + (1-\beta) \underline{db^2} \leftarrow \text{large}$$

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw}}} \leftarrow$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db}}} \leftarrow$$

Adam optimization algorithm

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

On iteration t :

Compute dw, db using current mini-batch

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

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

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

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

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

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

Hyperparameters choice:

→ α : needs to be tune

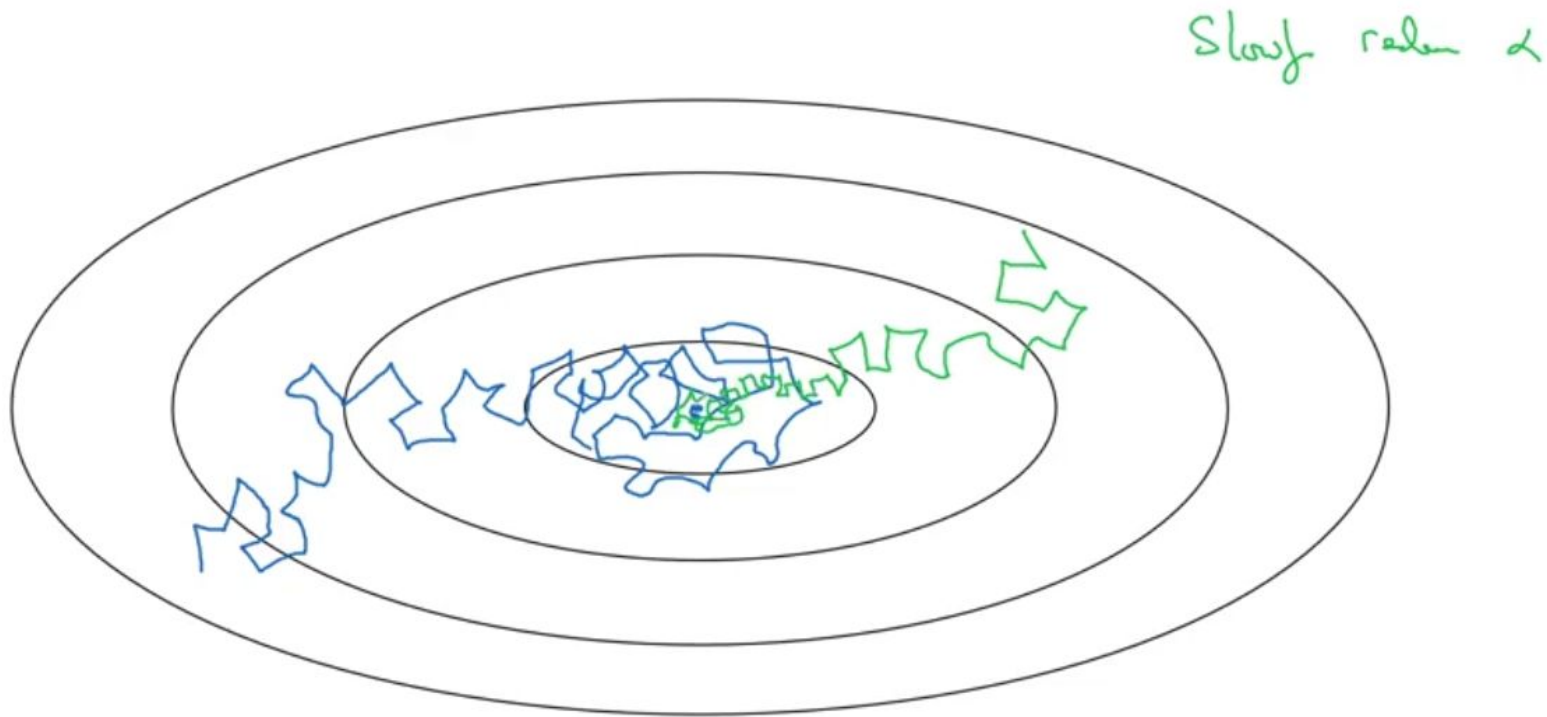
→ β_1 : 0.9 → (dw)

→ β_2 : 0.999 → (dw^2)

→ ϵ : 10^{-8}

Adam: Adaptive moment estimation

Learning Rate Decay

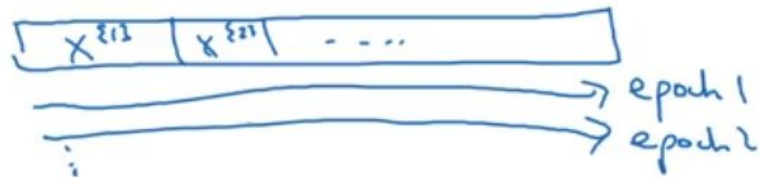


Learning rate decay

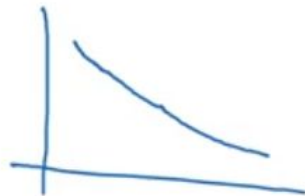
1 epoch = 1 pass through data.

$$\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch-num}} \alpha_0$$

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
\vdots	\vdots




$$\alpha_0 = 0.2$$
$$\text{decay-rate} = 1$$



Other learning rate decay methods

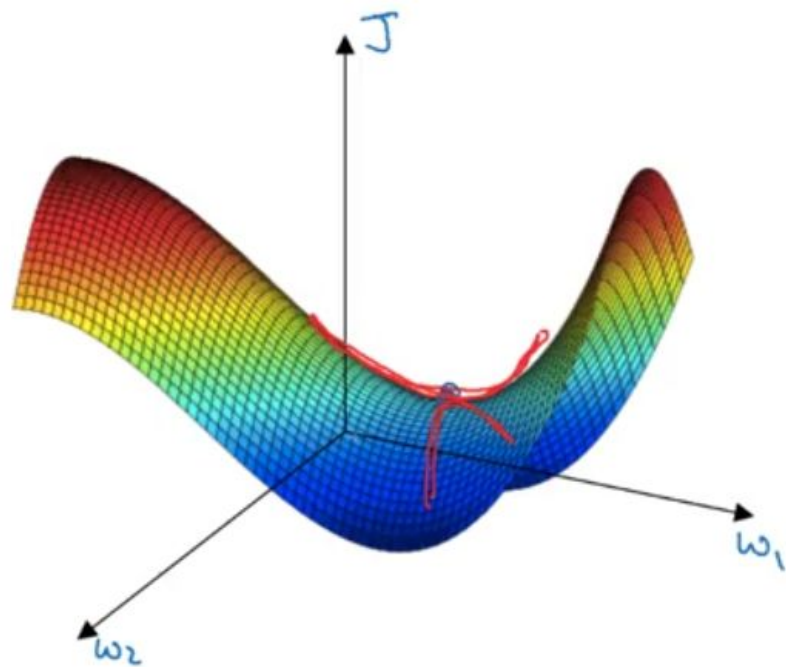
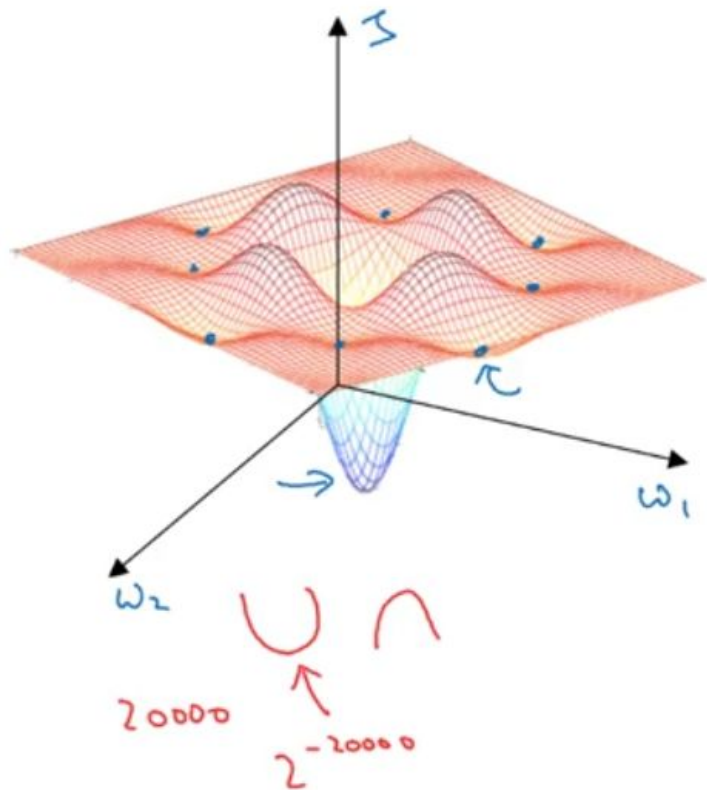
formula {

$$\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \quad - \text{exponentially decay.}$$
$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0$$


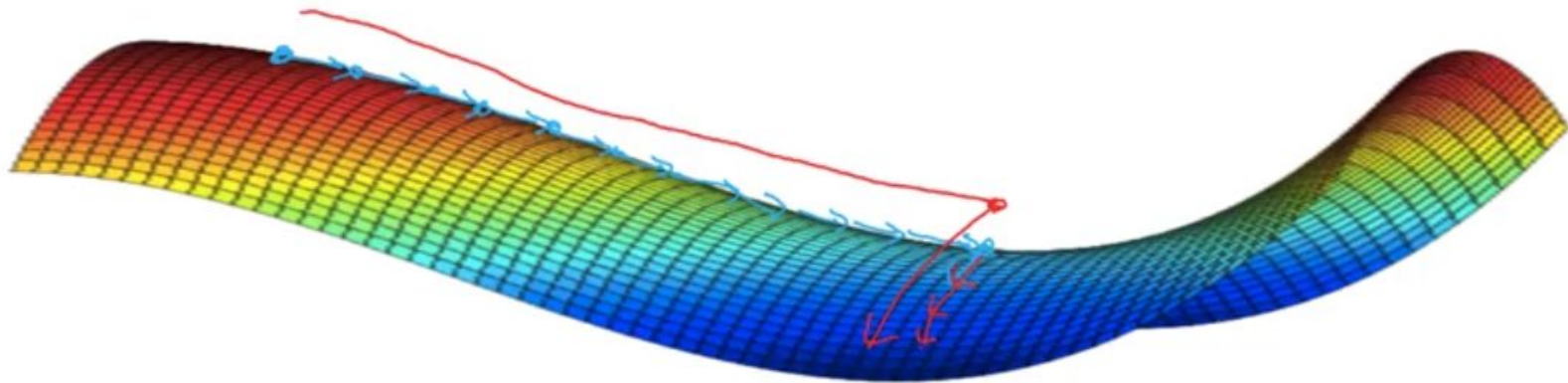
discrete staircase

Manual decay.

Local optima in neural networks



Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow

Reference

<https://machinelearningmastery.com/gradient-descent-for-machine-learning/>

<https://www.analyticsvidhya.com/blog/2017/03/introduction-to-gradient-descent-algorithm-along-its-variants/>

<http://neuralnetworksanddeeplearning.com/chap1.html>

<https://towardsdatascience.com/gradient-descent-in-a-nutshell-eaf8c18212f0>

<https://www.coursera.org/learn/deep-neural-network/home/week/2>