

# Computational Intelligence

Samaneh Hosseini

Isfahan University of Technology

# Outline

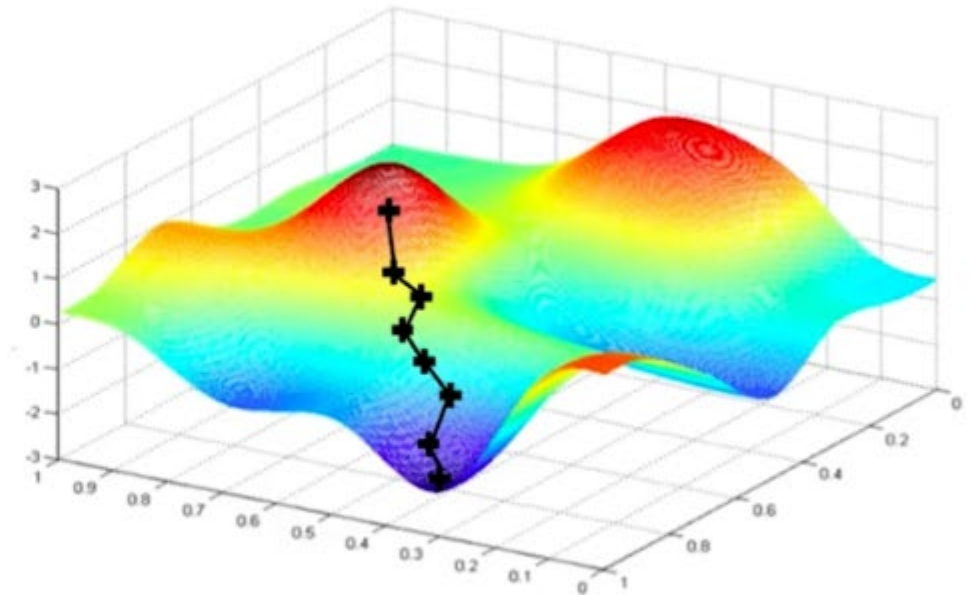
- Neural Networks in Practice:
  - Mini-batches
  - Mini-batch Gradient Descent
  - Understanding Mini-batch Gradient Descent
  - • Exponentially Weighted Averages

# Neural Networks in Practice: Mini-batches

# Gradient Descent

## Algorithm

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



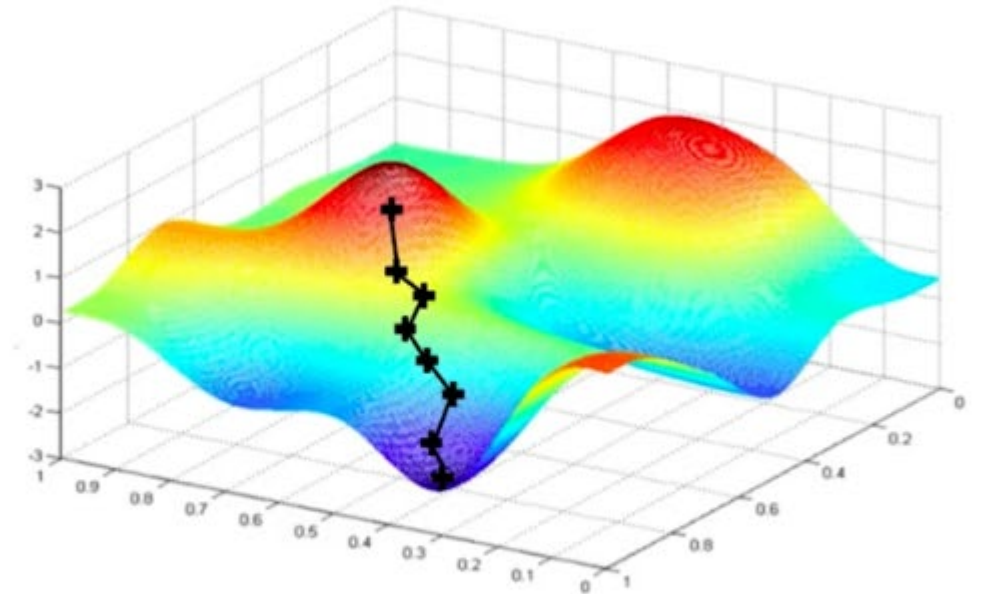
تکرار برای هر round

تعداد بار حلقه

# 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



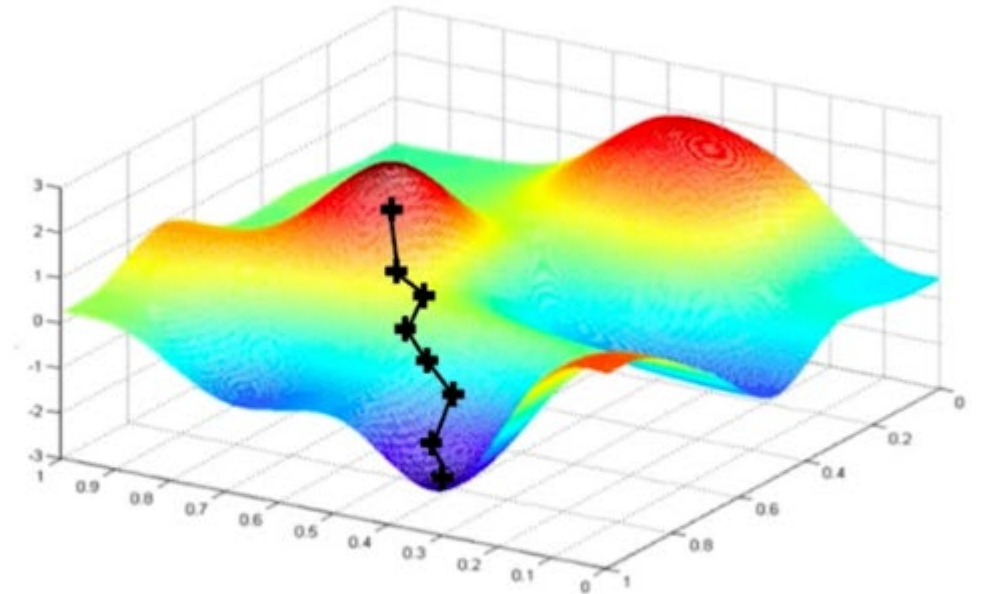
Can be very  
**computationally**  
intensive to compute!



# 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

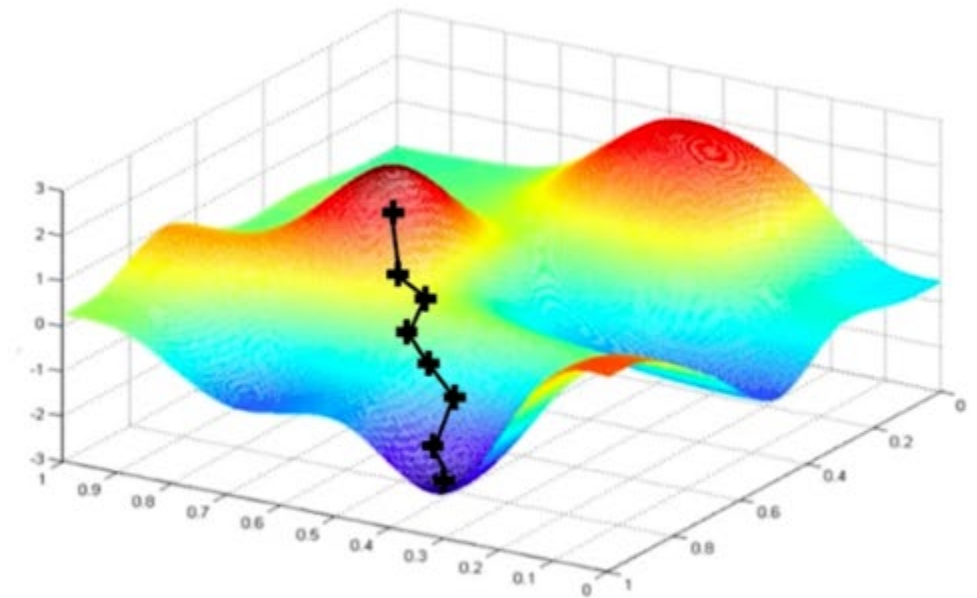


Can be very  
**computationally**  
intensive to compute!

# Stochastic Gradient Descent

## Algorithm

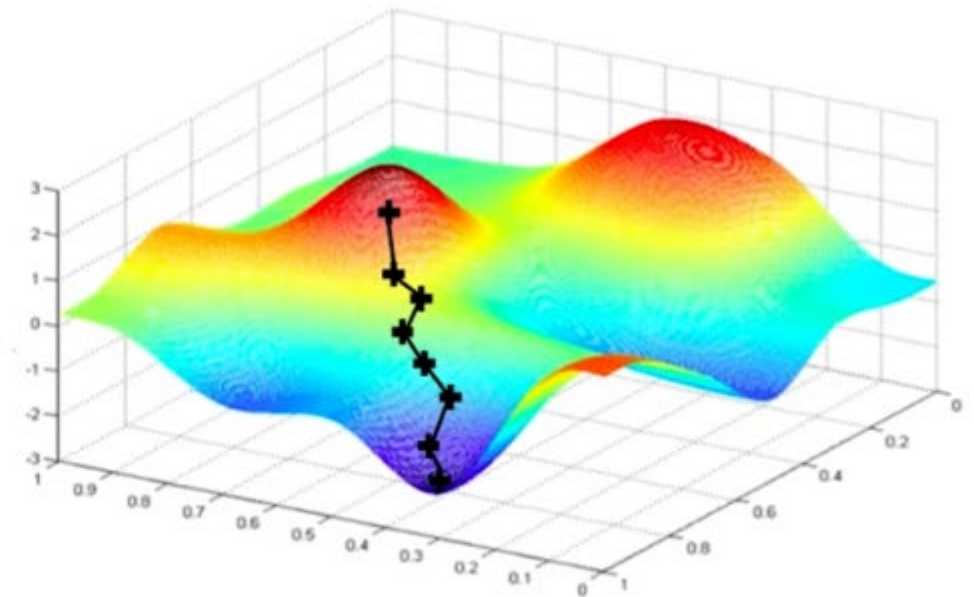
1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point  $i$
4. Compute gradient,  $\frac{\partial J_i(W)}{\partial W}$
5. Update weights,  $\underline{W} \leftarrow \underline{W} - \eta \frac{\partial J(W)}{\partial W}$
6. Return weights



# Stochastic Gradient Descent

## Algorithm

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



Easy to compute but  
very noisy (stochastic)!

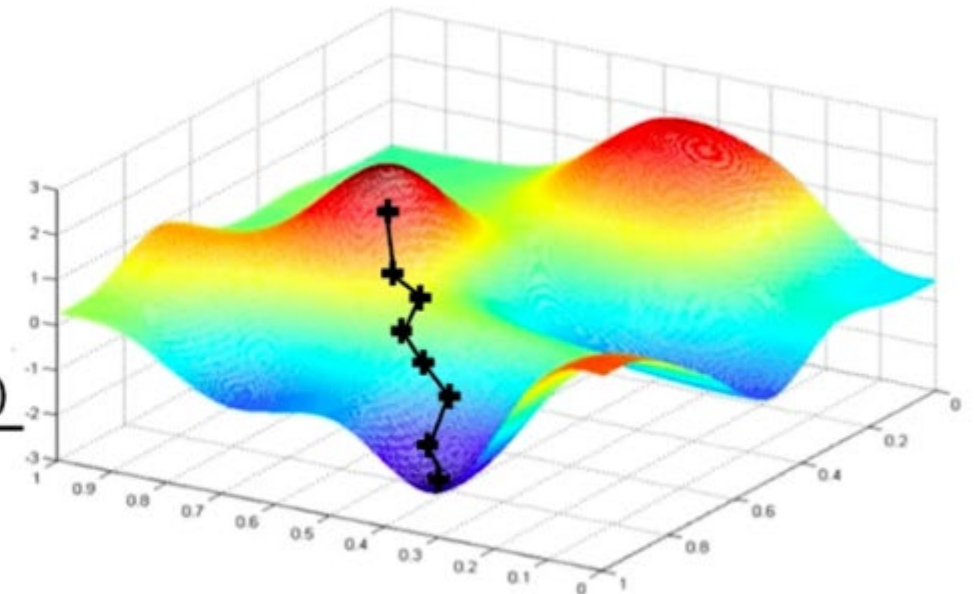


# Stochastic Gradient Descent

*mini-batch*

## Algorithm

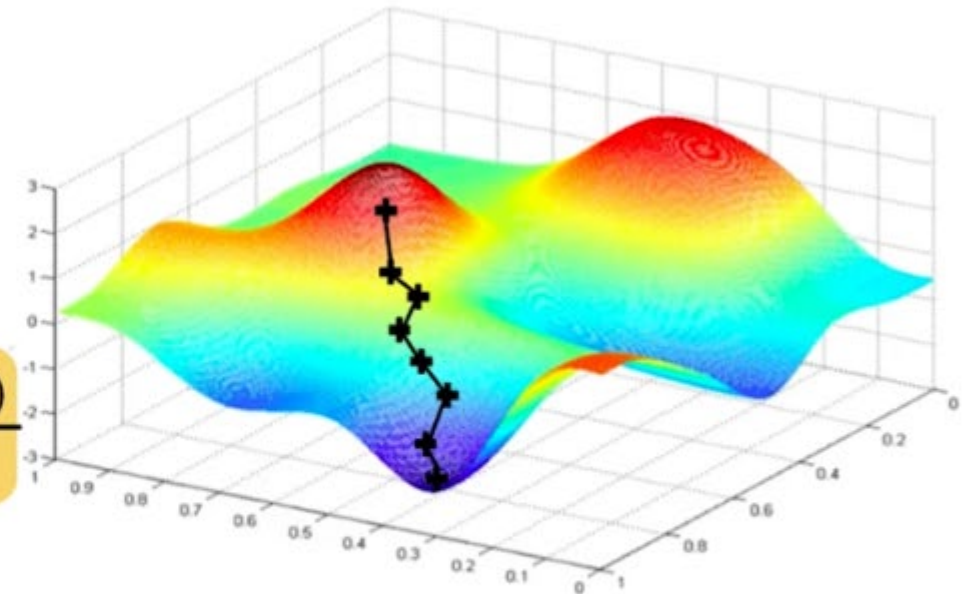
1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of  $B$  data points
4. Compute gradient,  $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights



# Stochastic Gradient Descent

## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3.     Pick batch of  $B$  data points
4.     Compute gradient,  $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5.     Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights



Fast to compute and a much better  
estimate of the true gradient!

# Batch vs. mini-batch gradient descent

$$X, Y$$

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

- Vectorization allows you to efficiently compute on m examples

$$X_{(n_x, m)} = \begin{bmatrix} X^{(1)} & X^{(2)} & X^{(3)} & \dots & X^{(1000)} \end{bmatrix} \quad \underbrace{\begin{bmatrix} X^{(1001)} & \dots & X^{(2000)} \end{bmatrix}}_{X^{\{2\}}_{(n_x, 1000)}} \quad X^{(2001)} \dots$$

$X^{\{1\}}_{(n_x, 1000)}$        $X^{\{2\}}_{(n_x, 1000)}$

$$Y_{(1, m)} = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} \end{bmatrix} \quad \underbrace{\begin{bmatrix} y^{(1001)} & \dots & y^{(2000)} \end{bmatrix}}_{Y^{\{2\}}_{(1, 1000)}} \quad \dots$$

$Y^{\{1\}}_{(1, 1000)}$        $Y^{\{2\}}_{(1, 1000)}$

$$\underbrace{\begin{bmatrix} \dots & X^{(m)} \end{bmatrix}}_{X^{\{5000\}}_{(n_x, 1000)}}$$

$$\underbrace{\begin{bmatrix} \dots & y^{(m)} \end{bmatrix}}_{Y^{\{5000\}}_{(1, 1000)}}$$

$$m = 5,000,000$$

5000 mini batches of 1000 each

mini\_batch t :  $X^{\{t\}}, Y^{\{t\}}$

$$\begin{matrix} X^{(i)} \\ Z^{[L]} \\ X^{\{t\}}, Y^{\{t\}} \end{matrix}$$

# Mini-batch gradient descent

1 step of gradient descent  
using  $X^{(t)}$ ,  $Y^{(t)}$   
1000 ↙

$X, Y$

repeat {  
1 epoch  
↓  
1 pass through  
all training  
set

for  $t = 1 \dots 5000$  {

Forward pass on  $X^{(t)}$

$$Z^{(1)} = \underline{w}^{(1)} X^{(t)} + b^{(1)}$$

$$A^{(1)} = g^{(1)}(Z^{(1)})$$

⋮

$$A^{(L)} = g^{(L)}(Z^{(L)})$$

vectorized implementation  
(1000 example)

Compute cost  $J^{(t)} = \frac{1}{1000} \sum_{i=1}^L \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|w^{(l)}\|_F^2$

Backprop to compute gradient wrt  $J^{(t)}$  (using  $X^{(t)}, Y^{(t)}$ )

$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

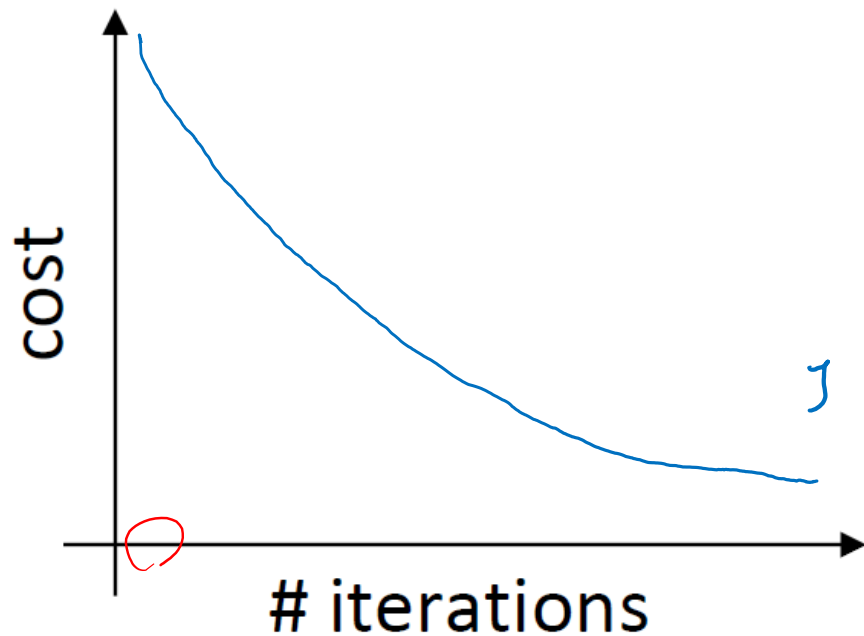
}  
= {

# Optimization Algorithms: Understanding Mini-batch Gradient Descent

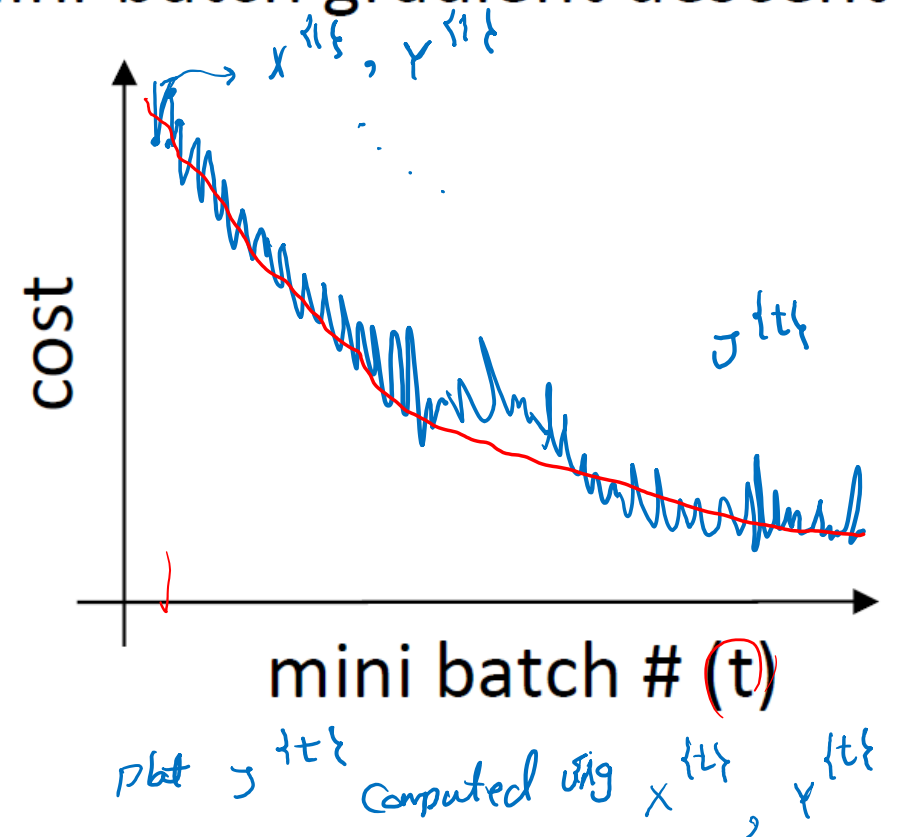


# Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent



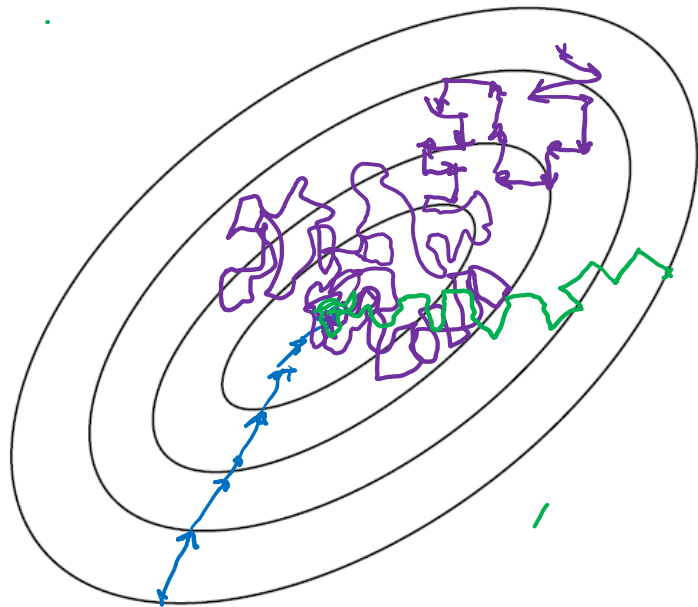
# Choosing your mini-batch size

→ if mini\_batch size =  $m$  : Batch gradient descent

→ if mini\_batch size = 1 : stochastic gradient descent  
 $(x^{(1)}, y^{(1)}) = (x^{(1)}, y^{(1)})$

In practice : somewhere between 1 and  $m$

$$(x^{(1)}, y^{(1)}) = (x, y)$$



Stochastic  
gradient  
descent



lose speedup  
from vectorization

In-between

(mini\_batch size  
not too big / small)



Fastest learning

- vectorization on per iteration (1000)

- Make progress without using entire trajectory

batch  
gradient descent  
(mini\_batch size =  $m$ )



too long

# Choosing your mini-batch size

if small training set : use batch gradient descent  
( $m \leq 2000$ )

Typical mini-batch size :

→ 64, 128, 256, 512  
 $2^6$      $2^7$      $2^8$      $2^9$

make sure mini-batch fit in CPU/GPU memory

$X^{\{t\}}$ ,  $y^{\{t\}}$

# Optimization Algorithms: Exponentially Weighted Average

# Temperature in London

$$\theta_1 = 40^\circ\text{F}$$

$$\theta_2 = 49^\circ\text{F}$$

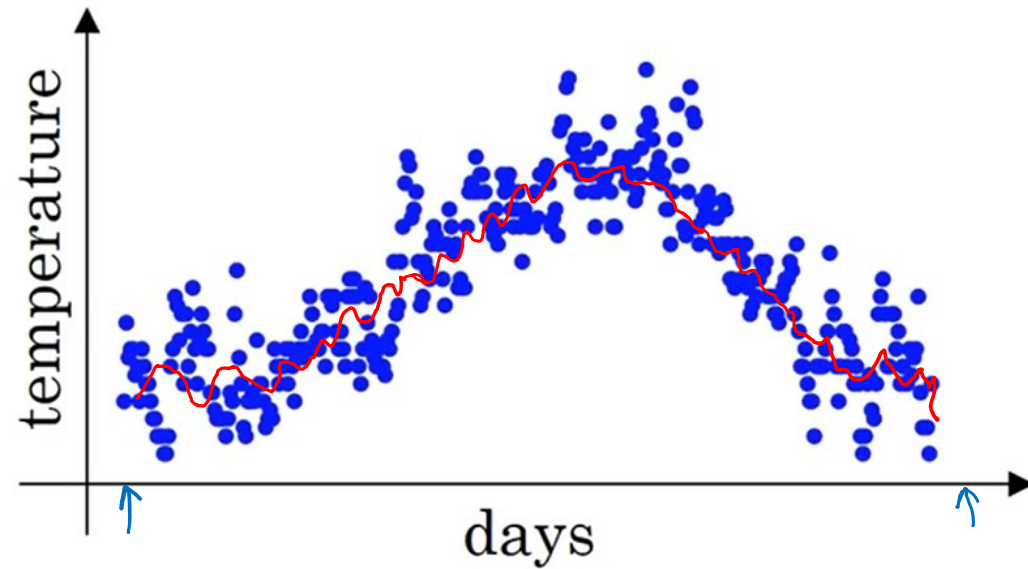
$$\theta_3 = 45^\circ\text{F}$$

$\vdots$

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

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

$\vdots$



$$v_0 = 0$$

$$v_1 = 0.9v_0 + 0.1\theta_1$$

$$v_2 = 0.9v_1 + 0.1\theta_2$$

$$v_3 = 0.9v_2 + 0.1\theta_3$$

$\vdots$

$$v_t = 0.9v_{t-1} + 0.1\theta_t$$



# Exponentially weighted averages

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

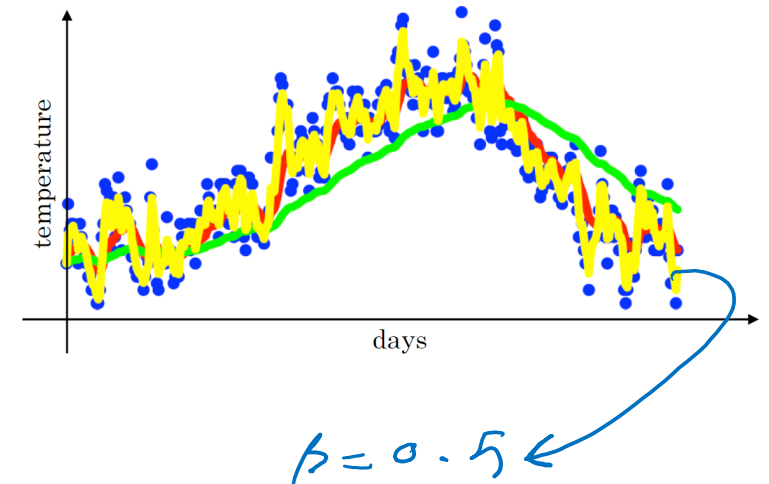
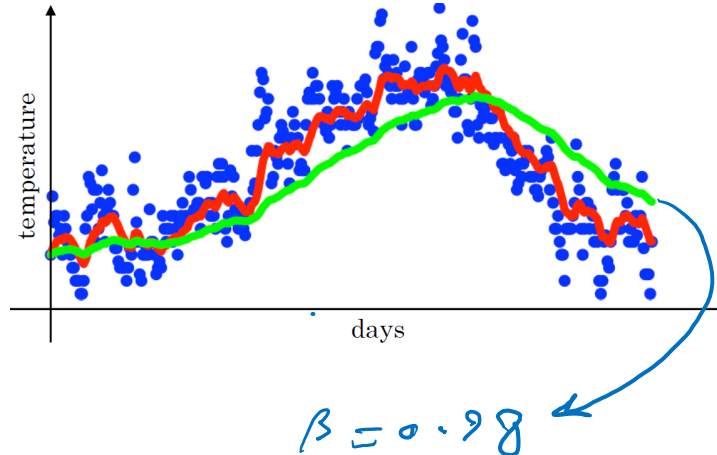
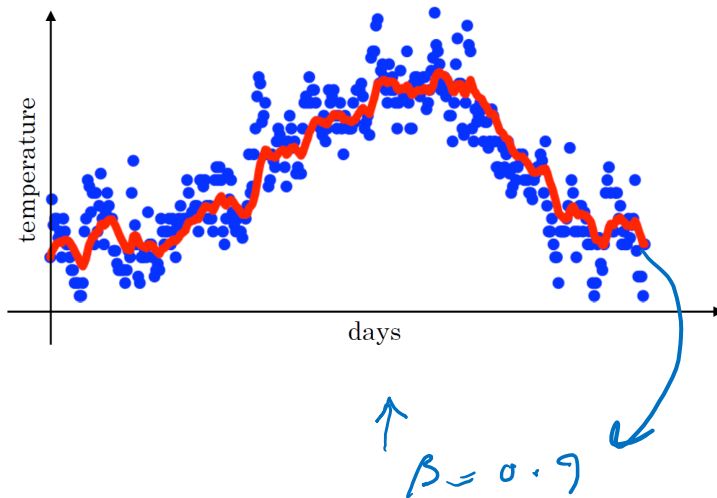
$\beta = 0.9$  :  $\approx 10$  datapoint temperature

$\beta = 0.98$  :  $\approx 50$  "

$\beta = 0.5$  :  $\approx 2$  "

$V_t$  average over

$\approx \frac{1}{1-\beta}$  datapoint



# Core Foundation Review

- Neural Networks in Practice:
  - Mini-batches
  - Mini-batch Gradient Descent
  - Understanding Mini-batch Gradient Descent
  - Exponentially Weighted Averages