

Welcome to

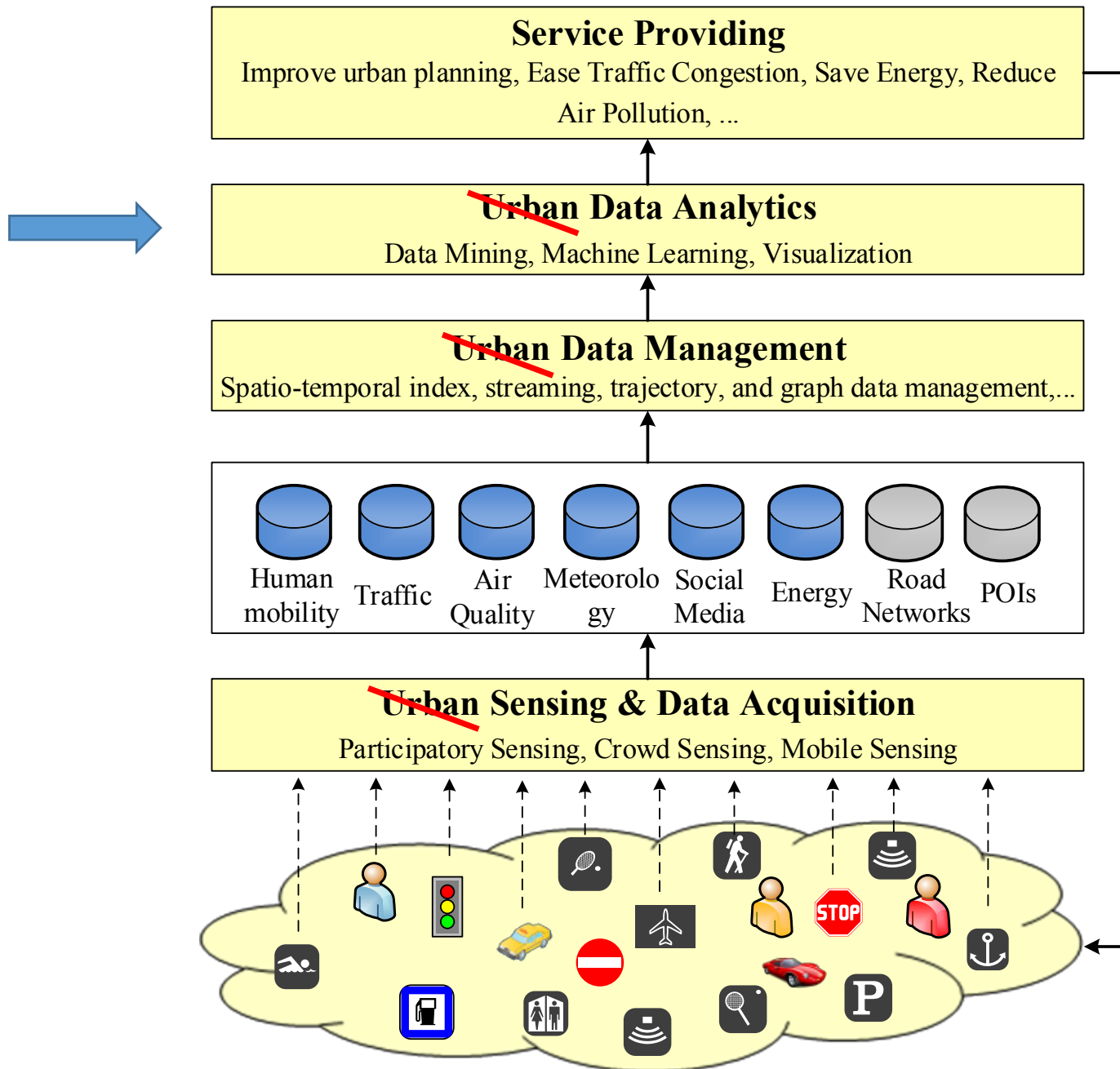
*DS3010:*  
*DS-III: Computational Data Intelligence*  
Gradient Descent  
Prof. Yanhua Li

Time: 11:00am – 12:50pm M & R

Location: HL 114

D-term 2022

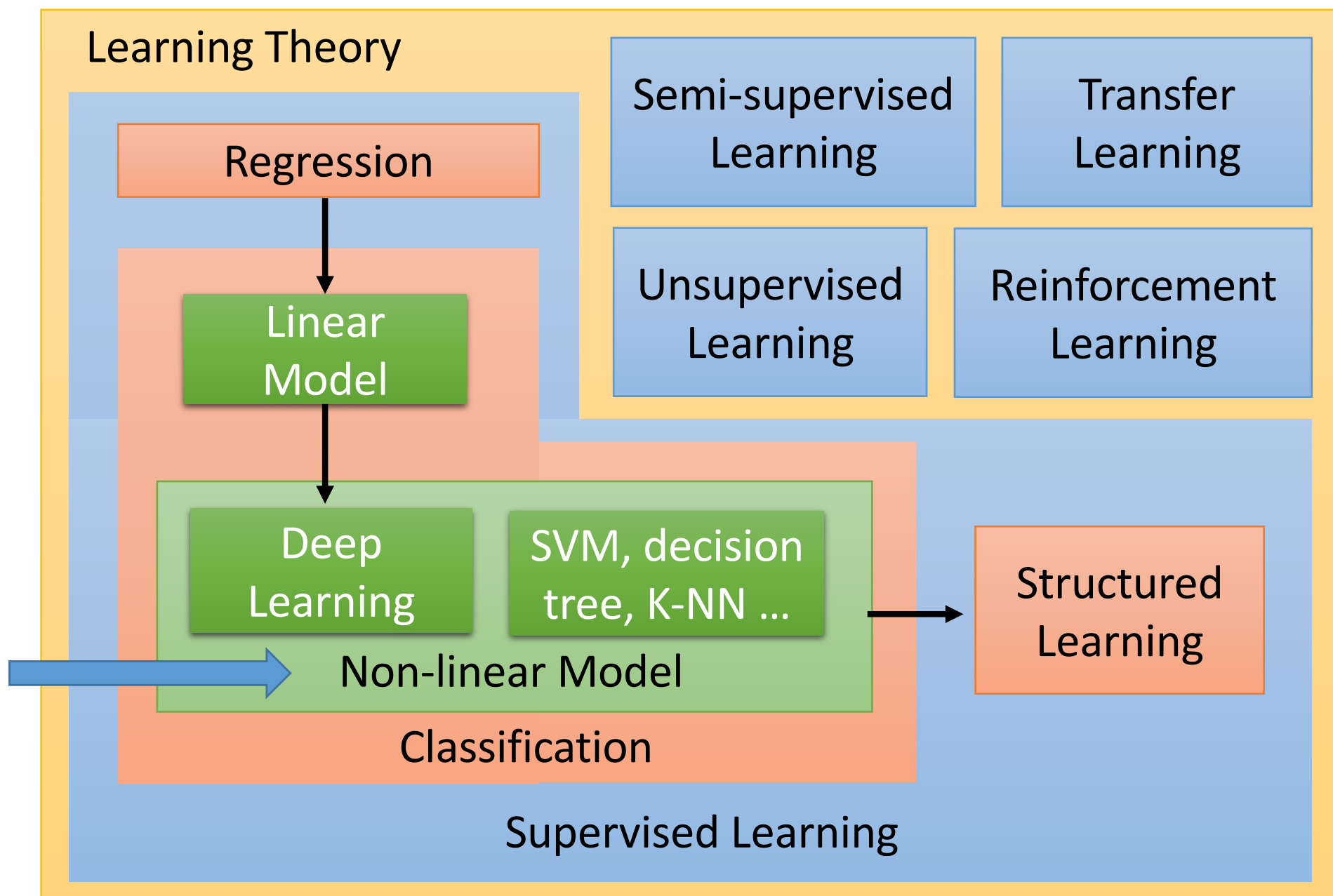
# Data pipeline



**Urban Computing: concepts, methodologies, and applications.**  
Zheng, Y., et al. *ACM transactions on Intelligent Systems and Technology*.

# Learning Map

scenario task method



# Gradient Descent

# Review: Gradient Descent

- In step 3, we have to solve the following optimization problem:

$$\theta^* = \arg \min_{\theta} L(\theta) \quad \text{L: loss function} \quad \theta: \text{parameters}$$

---

Suppose that  $\theta$  has two variables  $\{\theta_1, \theta_2\}$

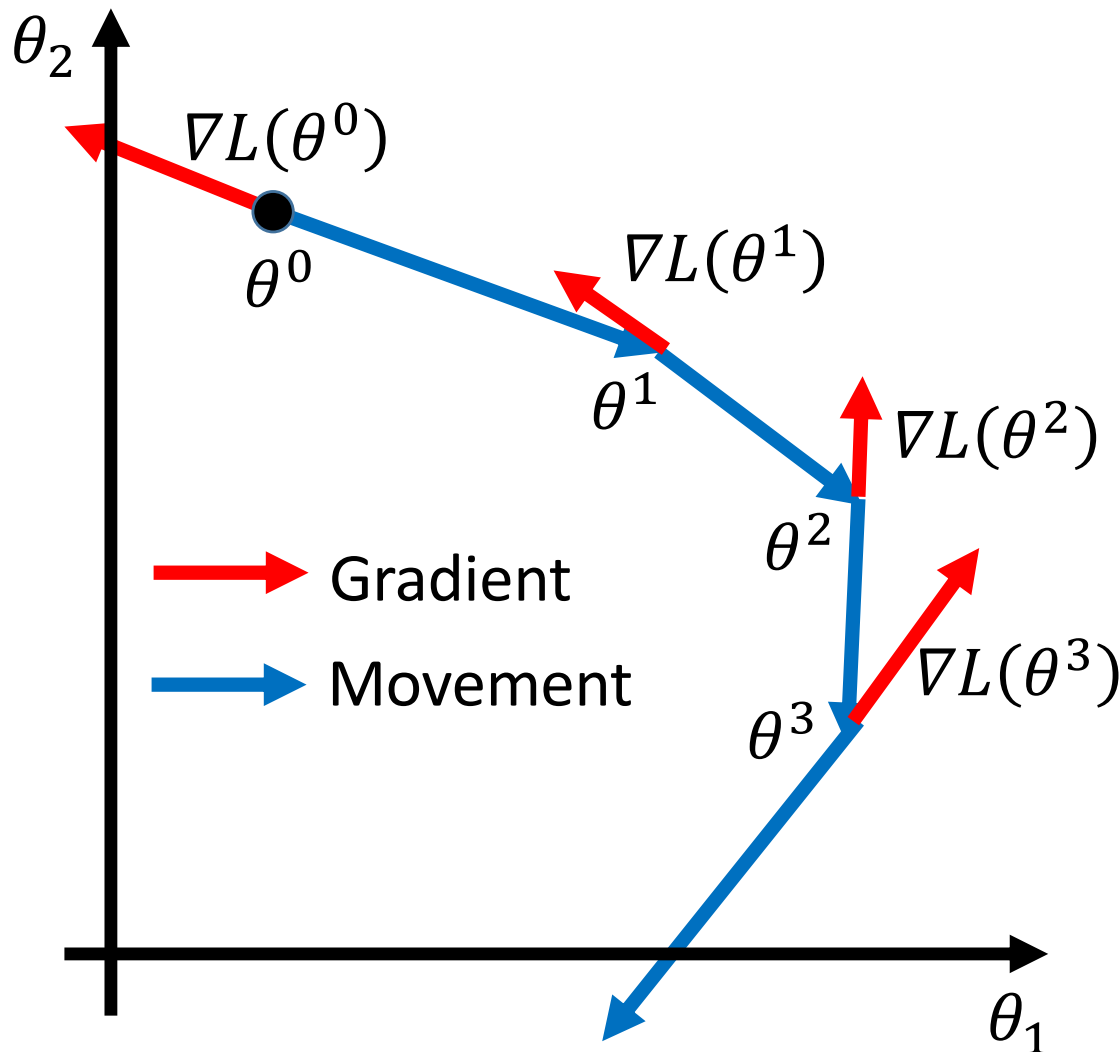
Randomly start at  $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta_1)/\partial \theta_1 \\ \partial L(\theta_2)/\partial \theta_2 \end{bmatrix}$$

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta_1^0)/\partial \theta_1 \\ \partial L(\theta_2^0)/\partial \theta_2 \end{bmatrix} \Rightarrow \theta^1 = \theta^0 - \eta \nabla L(\theta^0)$$

$$\begin{bmatrix} \theta_1^2 \\ \theta_2^2 \end{bmatrix} = \begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta_1^1)/\partial \theta_1 \\ \partial L(\theta_2^1)/\partial \theta_2 \end{bmatrix} \Rightarrow \theta^2 = \theta^1 - \eta \nabla L(\theta^1)$$

# Review: Gradient Descent



Start at position  $\theta^0$

Compute gradient at  $\theta^0$

Move to  $\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

Compute gradient at  $\theta^1$

Move to  $\theta^2 = \theta^1 - \eta \nabla L(\theta^1)$

⋮

# Gradient Descent

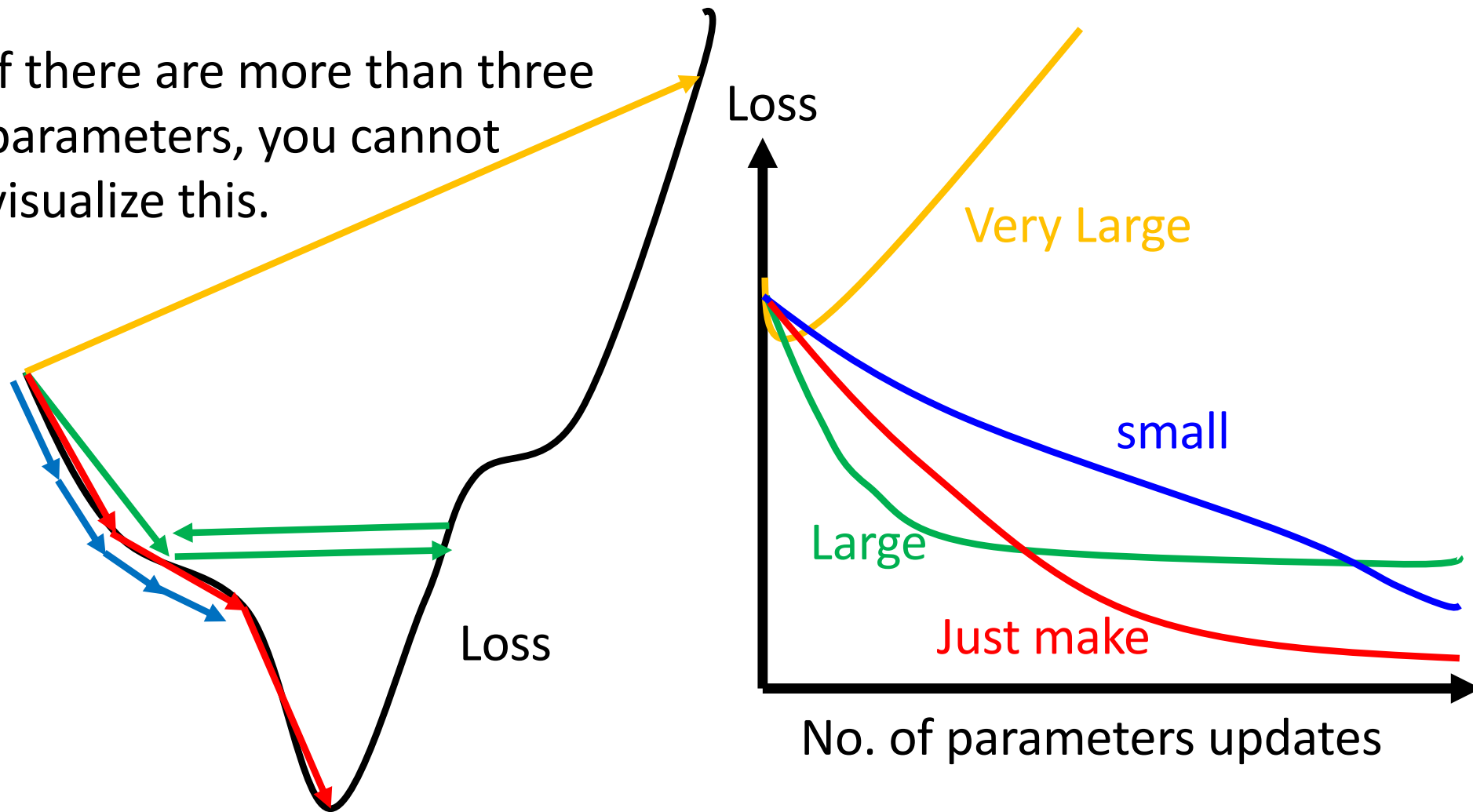
Tip 1: Tuning your  
learning rates

# Learning Rate

$$\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$$

Set the learning rate  $\eta$  carefully

If there are more than three parameters, you cannot visualize this.



But you can always visualize this.



# Adaptive Learning Rates

- Popular & Simple Idea: Reduce the learning rate by some factor every few epochs.
  - At the beginning, we are far from the destination, so we use larger learning rate
  - After several epochs, we are close to the destination, so we reduce the learning rate
  - E.g. 1/t decay:  $\eta^t = \eta / \sqrt{t + 1}$
- Learning rate cannot be one-size-fits-all
  - Giving different parameters different learning rates

$\eta$ : Base learning rate

Adagrad

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad g^t = \frac{\partial L(\theta^t)}{\partial w}$$

- Divide the learning rate of each parameter by the ***root mean square of its previous derivatives***

Vanilla Gradient descent

$$w^{t+1} \leftarrow w^t - \eta^t g^t$$

w is one parameters

Adagrad

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

$\sigma^t$ : ***root mean square*** of the previous derivatives of parameter w

Parameter dependent

# Adagrad

$\sigma^t$ : *root mean square* of the previous derivatives of parameter  $w$

$$w^1 \leftarrow w^0 - \frac{\eta^0}{\sigma^0} g^0$$

$$w^2 \leftarrow w^1 - \frac{\eta^1}{\sigma^1} g^1$$

$$w^3 \leftarrow w^2 - \frac{\eta^2}{\sigma^2} g^2$$

$\vdots$

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

$$\sigma^0 = \sqrt{(g^0)^2}$$

$$\sigma^1 = \sqrt{\frac{1}{2} [(g^0)^2 + (g^1)^2]}$$

$$\sigma^2 = \sqrt{\frac{1}{3} [(g^0)^2 + (g^1)^2 + (g^2)^2]}$$

$$\sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$$

# Adagrad

- Divide the learning rate of each parameter by the ***root mean square of its previous derivatives***

The diagram illustrates the Adagrad update rule. It shows the update equation  $w^{t+1} \leftarrow w^t - \eta^t \sigma^t g^t$  with  $\eta^t$  highlighted in an orange box and  $\sigma^t$  in a blue box. A red arrow points from the orange box to the definition  $\eta^t = \frac{\eta}{\sqrt{t+1}}$  with the text "1/t decay" in red. A blue arrow points from the blue box to the definition  $\sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$ , where the  $\frac{1}{t+1}$  term is also crossed out with a red line. A large blue arrow points from the first equation to the simplified version below.

$$w^{t+1} \leftarrow w^t - \eta^t \sigma^t g^t$$
$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad \text{1/t decay}$$
$$\sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$$
$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

Contradiction?  $\eta^t = \frac{\eta}{\sqrt{t+1}}$   $g^t = \frac{\partial L(\theta^t)}{\partial w}$

### Vanilla Gradient descent

$$w^{t+1} \leftarrow w^t - \eta^t \underline{g^t} \longrightarrow \text{Larger gradient, larger step}$$

### Adagrad

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (\underline{g^i})^2}} \underline{g^t}$$

Larger gradient, larger step

Larger gradient, smaller step

# Intuitive Reason

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad g^t = \frac{\partial L(\theta^t)}{\partial w}$$

- How surprise it is Contrast

$g^0$	$g^1$	$g^2$	$g^3$	$g^4$	.....
0.001	0.001	0.003	0.002	0.1	.....
$g^0$	$g^1$	$g^2$	$g^3$	$g^4$	.....
10.8	20.9	31.7	12.1	0.1	.....

Extremely large

Extremely small

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

→ Make the contrast effect

# Gradient Descent

## Tip 2: Stochastic Gradient Descent

Make the training faster

# Stochastic Gradient Descent

$$L = \sum_n \left( \hat{y}^n - \left( b + \sum w_i x_i^n \right) \right)^2$$

Loss is the summation over all training examples

◆ **Gradient Descent**  $\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$

◆ **Stochastic Gradient Descent**

Faster!

Pick an example  $x^n$

$$L^n = \left( \hat{y}^n - \left( b + \sum w_i x_i^n \right) \right)^2 \quad \theta^i = \theta^{i-1} - \eta \nabla L^n(\theta^{i-1})$$

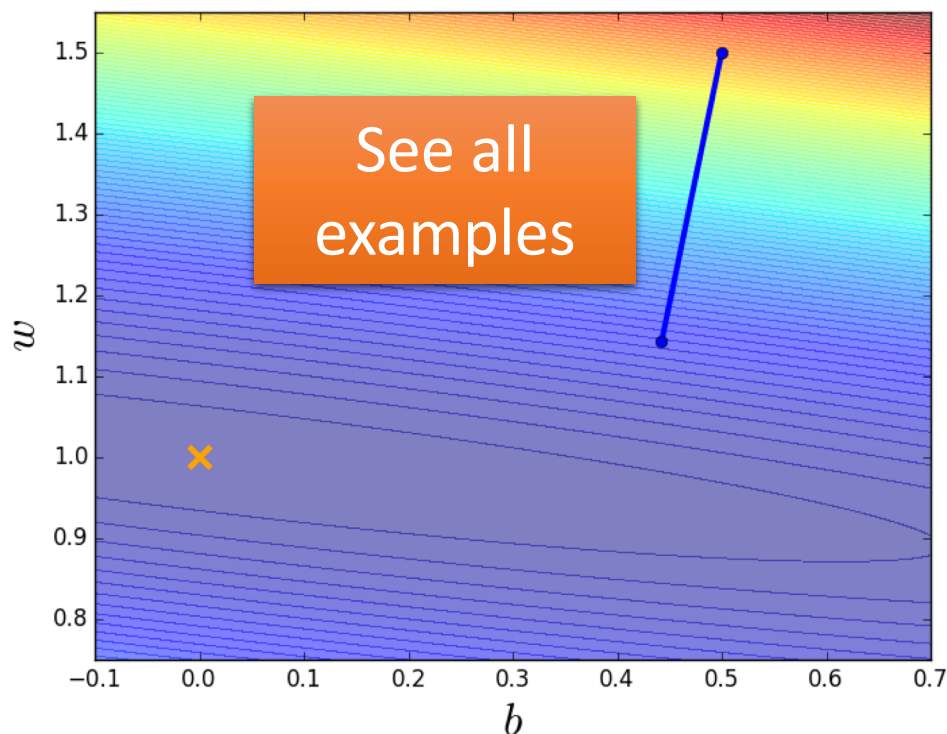
Loss for only one example



# Stochastic Gradient Descent

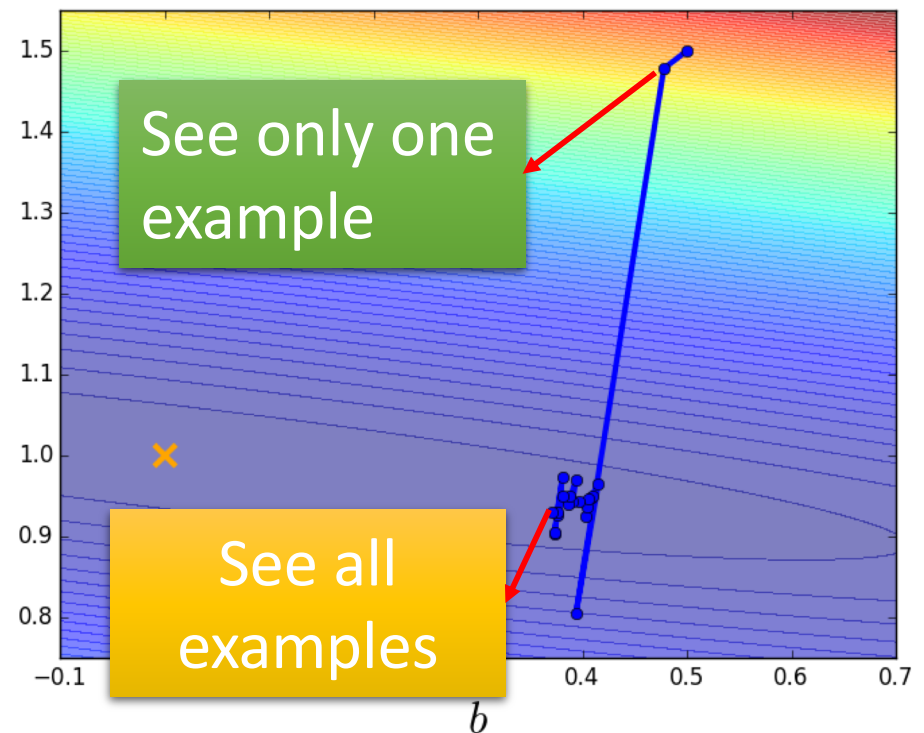
## Gradient Descent

Update after seeing all examples



## Stochastic Gradient Descent

Update for each example  
If there are 20 examples,  
20 times faster.



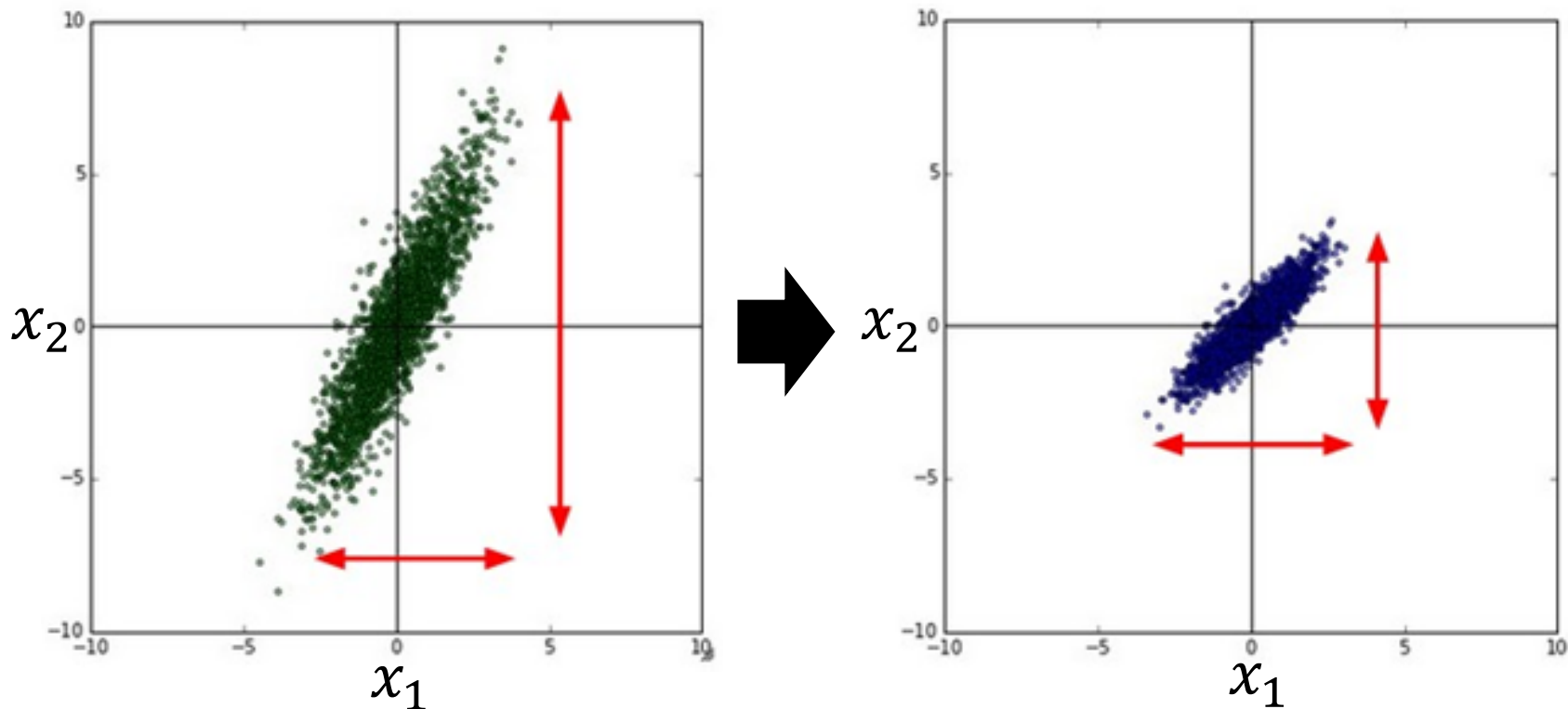
# Gradient Descent

## Tip 3: Feature Scaling

# Feature Scaling

Source of figure:  
<http://cs231n.github.io/neural-networks-2/>

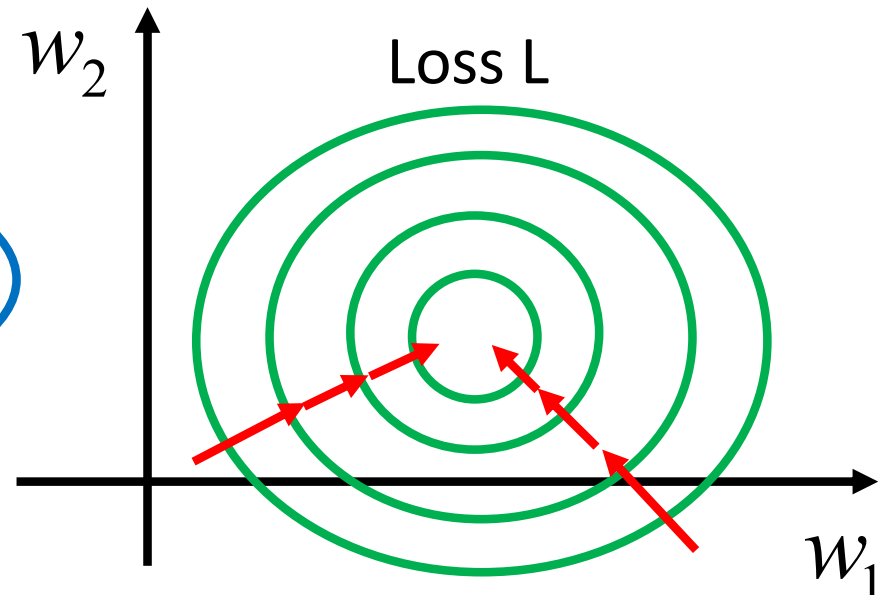
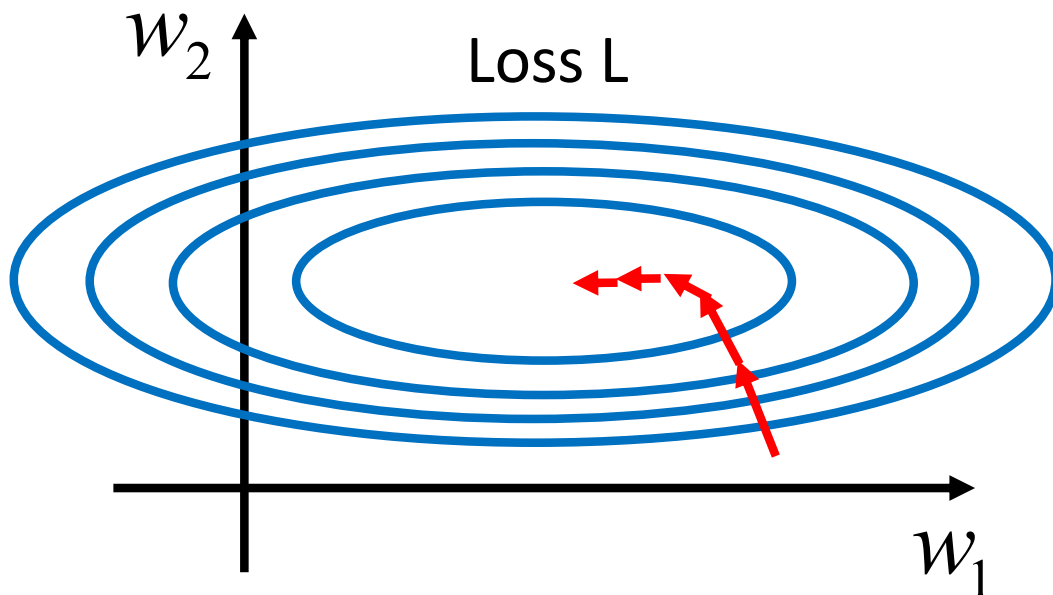
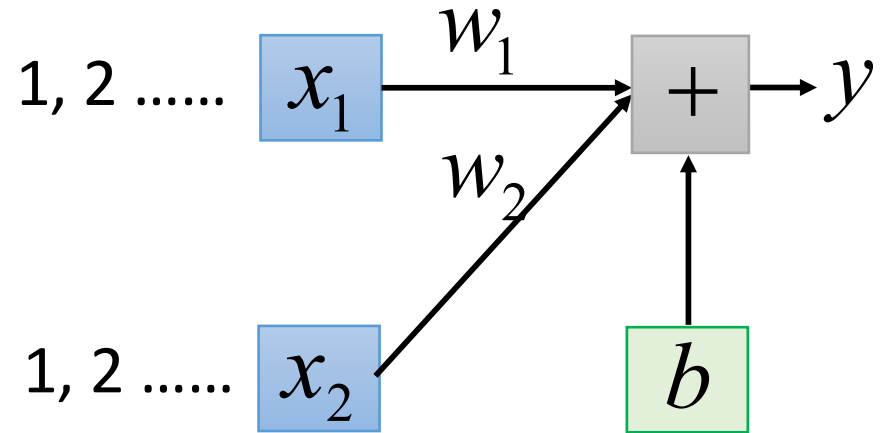
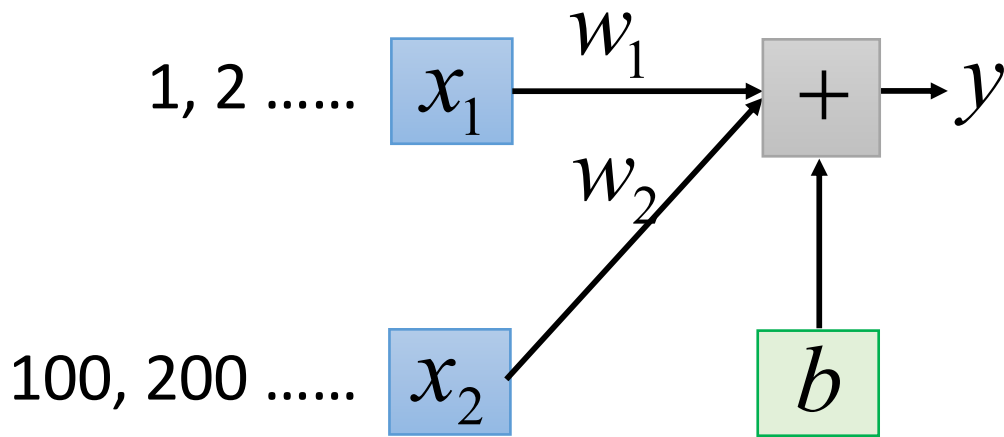
$$y = b + w_1x_1 + w_2x_2$$



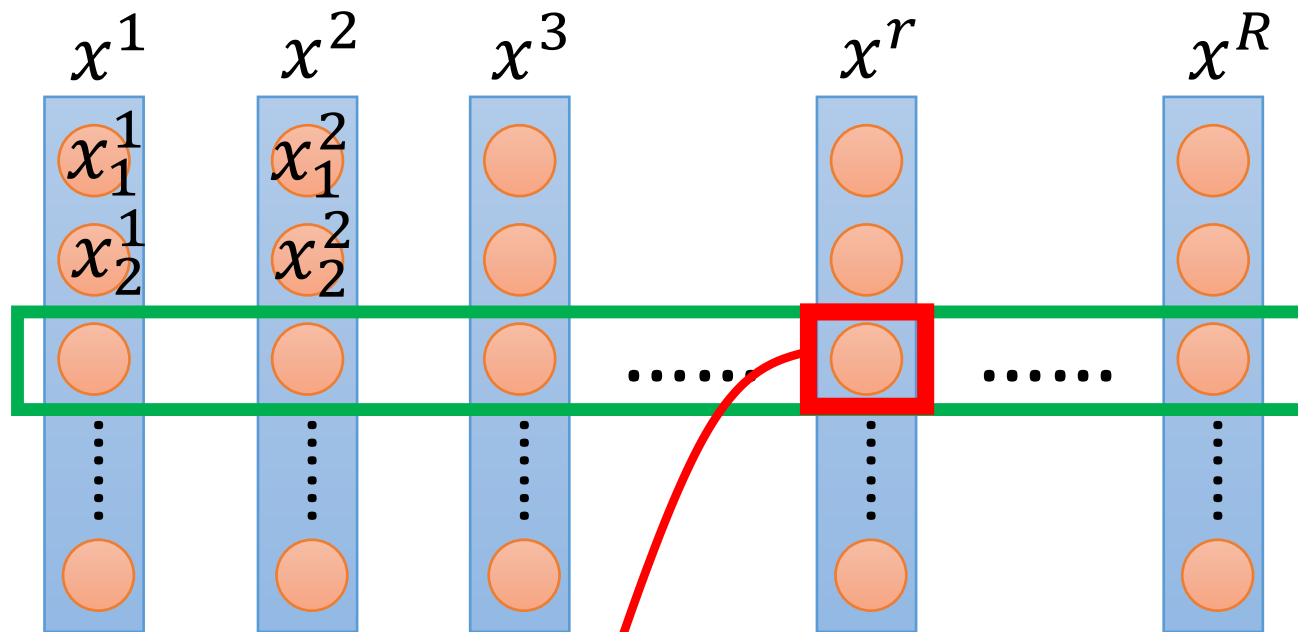
Make different features have the same scaling

# Feature Scaling

$$y = b + w_1x_1 + w_2x_2$$



# Feature Scaling



For each  
dimension  $i$ :

mean:  $m_i$

standard

deviation:  $\sigma_i$

$$x_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

The means of all dimensions are 0,  
and the variances are all 1

# Feature scaling in Python

```
from sklearn import preprocessing
```

```
scaler = preprocessing.StandardScaler().fit(x_train)
```

```
xtrain_scaled = scaler.transform(x_train)
```

```
xtest_scaled = scaler.transform(x_test)
```

- <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html>

# More on Gradient Descent Methods

- Vanilla Gradient Descent
- Adagrad
- RMSProp and Adam (with RMSProp and Momentum)
  - <https://blog.paperspace.com/intro-to-optimization-momentum-rmsprop-adam/>

# Questions