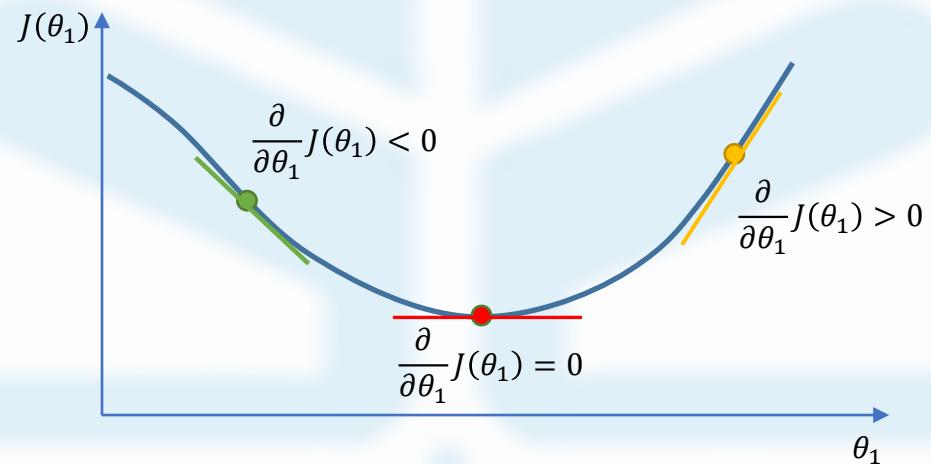


## Machine Learning

### Data Mining



Dr. Ali Valinejad

valinejad.ir

valinejad@umz.ac.ir

University of Mazandaran



## Outline:

- ❖ Supervised Learning
- ❖ Linear Regression
- ❖ Normal equations
- ❖ Gradient descent
  - Batch Gradient Descent
  - Stochastic Gradient Descent
  - Mini-Batch Gradient Descent
- ❖ Linear Regression: Probabilistic interpretation
- ❖ Locally weighted linear regression



# Supervised Learning

University of Mazandaran



# Supervised Learning

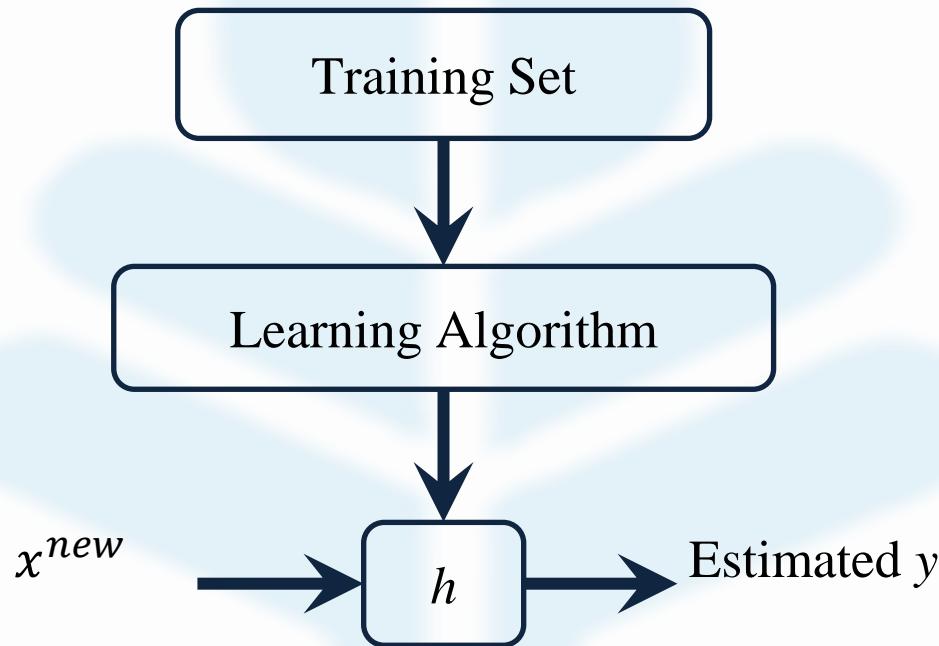
Price of a House



what do you think is the best guess for the price of the house  ?



# Supervised Learning



**In order to design a learning algorithm, we have to answer the following question:**

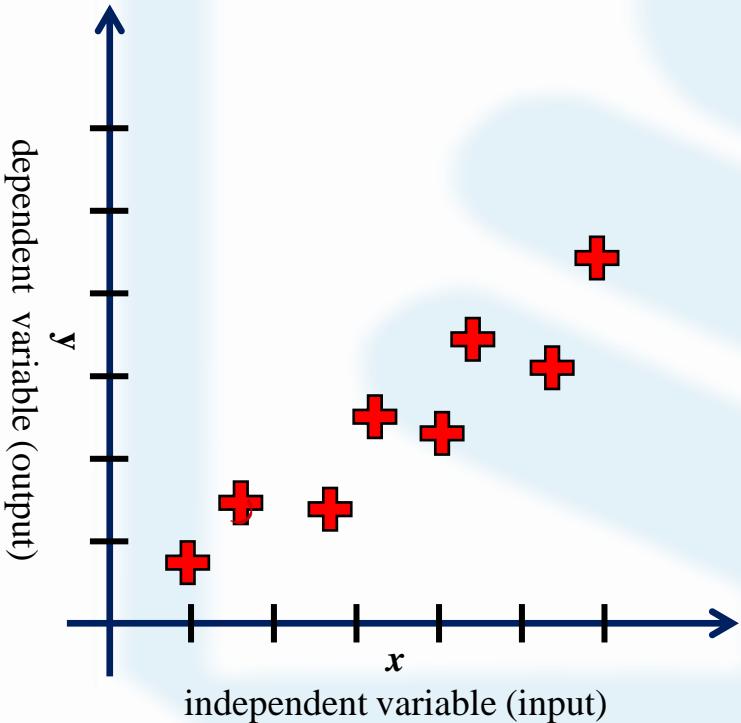
How we want to represent Hypothesis?

# Linear Regression

University of Mazandaran



# Linear Regression (One Variable)



- **One feature:**  $x$
- **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:** minimize  $J(\theta_0, \theta_1)$

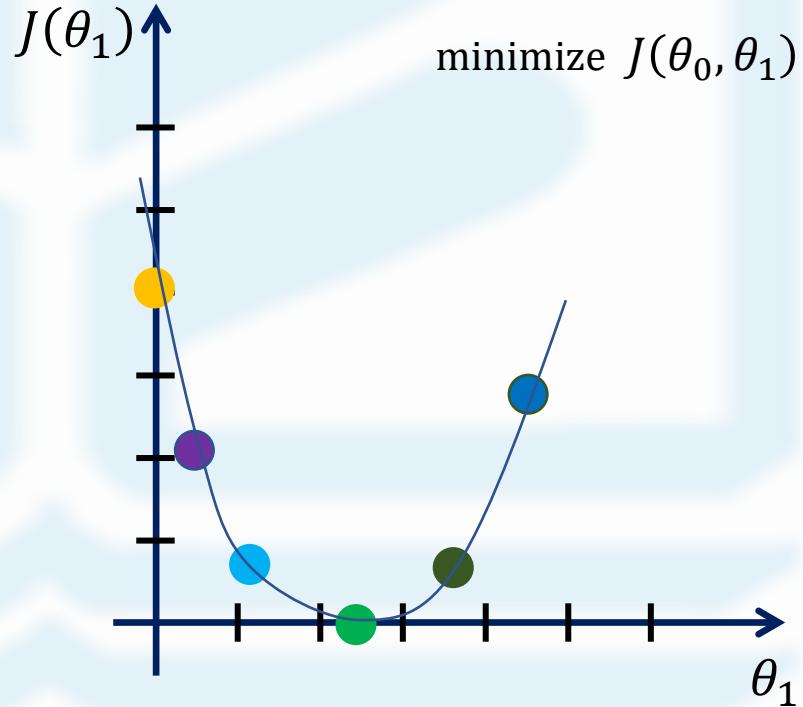
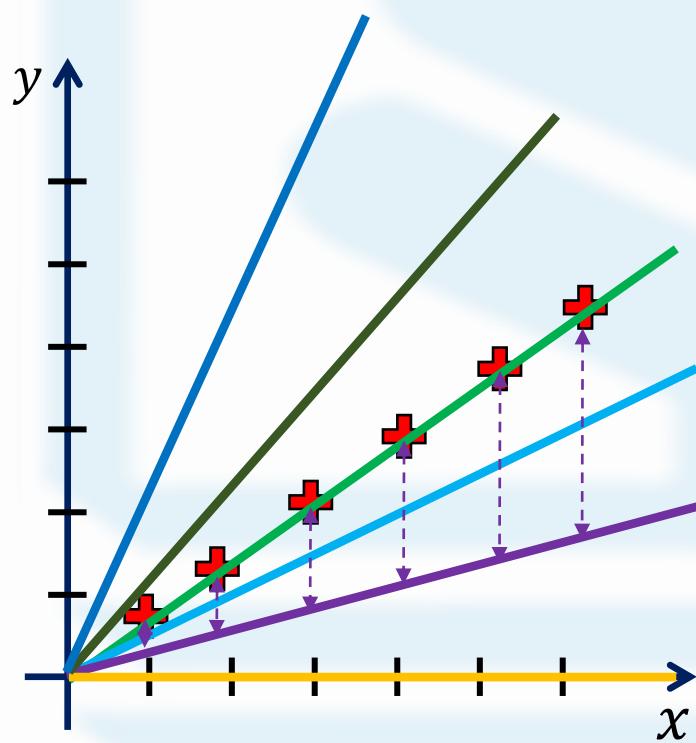
**Idea:** Choose  $\theta_0, \theta_1$  so that  $h_{\theta}(x)$  is close to  $y$  for our training examples  $(x, y)$ .



# Linear Regression (One Variable)

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

Assume  $\theta_0 = 0$ , so  $h_\theta(x) = \theta_1 x$  and  $J(\theta_0, \theta_1) = J(\theta_1)$



# Linear Regression (Multiple Variables)

- **Multiple features:**  $x_1, x_2, \dots, x_n$
- **Hypothesis:**  $x_0 := 1 \Rightarrow h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \boldsymbol{\theta}^T x$
- **Parameters:**  $\theta_0, \theta_1, \dots, \theta_n$
- **Cost function:**
$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$
- **Goal:** minimize  $J(\boldsymbol{\theta}) := J(\theta_0, \theta_1, \dots, \theta_n)$

- **Learning**
  - ❖ Solving normal equation  $\boldsymbol{\theta} = (X^T X)^{-1} X^T y$
  - ❖ Gradient descent

- **Inference**

$$\hat{y} = h_{\theta}(x^{\text{test}}) = \boldsymbol{\theta}^T x^{\text{test}}$$



# Normal equations to minimize $J(\theta)$



# Normal equations to minimize $J(\theta)$

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

$$\boldsymbol{\theta}^* = \operatorname*{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$



# Normal equations to minimize $J(\theta)$

Given a training set

$x_1^{(i)}$	...	$x_n^{(i)}$	$y^{(i)}$
$x_1^{(1)}$	...	$x_n^{(1)}$	$y^{(1)}$
$x_1^{(2)}$	...	$x_n^{(2)}$	$y^{(2)}$
$\vdots$	$\ddots$	$\vdots$	$\vdots$
$x_1^{(m)}$	...	$x_n^{(m)}$	$y^{(m)}$

set  $x_0^{(i)} = 1$ ,  $i = 1, 2, \dots, m$

$x_0^{(i)}$	$x_1^{(i)}$	...	$x_n^{(i)}$	$y^{(i)}$
1	$x_1^{(1)}$	...	$x_n^{(1)}$	$y^{(1)}$
1	$x_1^{(2)}$	...	$x_n^{(2)}$	$y^{(2)}$
$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
1	$x_1^{(m)}$	...	$x_n^{(m)}$	$y^{(m)}$



# Normal equations to minimize $J(\theta)$

Given a training set, define the **design matrix**  $X$

$$X := \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(m)})^T - \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times (n+1)}$$

where

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

set:

$$y := \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m \quad \theta := \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$



# Normal equations to minimize $J(\theta)$

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad X := \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad y := \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \theta := \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \sum_{j=0}^n \theta_j x_j$$

$$\Rightarrow h_{\theta}(x) = x^T \theta$$

$$\Rightarrow h_{\theta}(x^{(i)}) = (x^{(i)})^T \theta$$

$$\Rightarrow h_{\theta}(x^{(i)}) - y^{(i)} = (x^{(i)})^T \theta - y^{(i)}$$

$$\Rightarrow X^T \theta - y = \begin{bmatrix} (x^{(1)})^T \theta - y^{(1)} \\ (x^{(2)})^T \theta - y^{(2)} \\ \vdots \\ (x^{(m)})^T \theta - y^{(m)} \end{bmatrix}$$

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



# Normal equations to minimize $J(\theta)$

$$J(\theta) = \frac{1}{2m} (X^T \theta - y)^T (X^T \theta - y)$$

$$\Rightarrow \nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2m} (X^T \theta - y)^T (X^T \theta - y)$$

$$\Rightarrow \nabla_{\theta} J(\theta) = \frac{1}{m} (X^T X \theta - X^T y)$$

if  $\nabla_{\theta} J(\theta^*) = \mathbf{0}$  then  $\theta^* = (X^T X)^{-1} X^T y$

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



# Normal equations to minimize $J(\theta)$

**Examples:**  $m = 5$ .

	Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178
1	3000	4	1	38	540

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \\ 1 & 3000 & 4 & 1 & 38 \end{bmatrix} \quad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \\ 540 \end{bmatrix}$$

$$\theta^* = (X^T X)^{-1} X^T y$$



# Normal equations to minimize $J(\theta)$

$$\theta^* = (X^T X)^{-1} X^T y$$

What if  $X^T X$  is non-invertible?

❖ *Redundant features (linearly dependent).*

e.g. if  $x_1 = \text{size(in feet}^2)$  and  $x_2 = \text{size(in m}^2)$  then  
 $x_1$  and  $x_2$  are linearly independent.  $x_1$  can be removed.

❖ *Too many features (e.g. #training example  $\leq$  #features).*

Delete some features, or use regularization.

# Gradient descent to minimize $J(\theta)$

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

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



# Gradient Descent

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

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

For a single training example ( $m = 1$ ),

$$\theta_j := \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

This rule is called the **LMS** (*least mean squares*) update rule, and is also known as the **Widrow-Hoff** learning rule.



# Gradient descent to minimize $J(\theta)$

## Three types of gradient descents

- **Batch Gradient Descent:**

Parameters are updated after computing the gradient of error with respect to the *entire training set*

- **Stochastic Gradient Descent(SGD):**

Parameters are updated after computing the gradient of error with respect to a *single training example*

- **Mini Batch Gradient Descent:**

Parameters are updated after computing the gradient of error with respect to a *subset of the training set*



# **Gradient descent to minimize $J(\theta)$**

**Batch Gradient Descent**



# Batch Gradient Descent

```
Repeat {  
     $\theta := \theta - \alpha \nabla_{\theta} J(\theta)$   
}
```

$\alpha$ : Learning rate



# Batch Gradient Descent

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n = \sum_{j=0}^n \theta_j x_j$$

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

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_0} \\ \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix}$$

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

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}, \quad j = 1, 2, \dots, n$$



# Batch Gradient Descent

**Repeat** {

$$\theta := \theta - \alpha \nabla_{\theta} J(\theta)$$

}

$\alpha$ : Learning rate

**Repeat** {

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

(simultaneously update  $\theta_j$  for every  $j = 0, 1, \dots, n$ )

}



# Batch Gradient Descent

**Repeat** {

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{m} \left[ \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

(simultaneously update for every  $j = 0, 1, \dots, n$ )

}

This rule is called the **Least Mean Squares (LMS)** update rule for a training set of  $m$  data points, which is also known as the **Widrow-Hoff learning** rule.



# Batch Gradient Descent

```
def gradient_descent(X, y, theta, learning_rate, iterations):
    m, n = X.shape
    n = n-1
    cost_history = np.zeros(iterations)
    theta_history = np.zeros((iterations, n))
    for it in range(iterations):
        prediction = np.dot(X,theta)
        theta = theta -(1/m)*learning_rate*( X.T.dot((prediction - y)))
        theta_history[it,:] = theta.T
        cost_history[it] = cal_cost(theta,X, y)
    return theta, cost_history, theta_history
# -----
def cal_cost(theta, X, y):
    """ Calculates the cost for given theta, X and Y. """
    m = len(y)
    predictions = X.dot(theta)
    cost = 1/(2*m) * np.sum(np.square(predictions-y))
    return cost
```



# **Gradient descent to minimize $J(\theta)$**

**Stochastic Gradient Descent(SGD)**



# Stochastic Gradient Descent

*Randomly shuffle(reorder) examples in training set*

**Repeat {**

**for**  $i=1$  to  $m$  {

$$\theta_j := \theta_j - \alpha \frac{1}{m} [(h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}]$$

(simultaneously update for every  $j = 0, 1, \dots, n$ )

}

}



# Stochastic Gradient Descent

*Randomly shuffle(reorder) examples in training set*

**Repeat** {

**for**  $i=1$  to  $m$  {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} [(h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}]$$

        (simultaneously update for every  $j = 0, 1, \dots, n$ )

    }

}



# Stochastic Gradient Descent

```
def stochastic_gradient_descent(X, y, theta, learning_rate, iterations):
    m = len(y)
    cost_history = np.zeros(iterations)
    for it in range(iterations):
        cost = 0.0
        for i in range(m):
            rand_ind = np.random.randint(0, m)
            X_i = X[rand_ind,:].reshape(1, X.shape[1])
            y_i = y[rand_ind].reshape(1, 1)
            prediction = np.dot(X_i, theta)
            theta = theta -(1/m)*learning_rate*( X_i.T.dot((prediction - y_i)))
            cost += cal_cost(theta, X_i, y_i)
        cost_history[it] = cost
    return theta, cost_history

# -----
def cal_cost(theta, X, y):
    """ Calculates the cost for given theta, X and Y. """
    m = len(y)
    predictions = X.dot(theta)
    cost = 1/(2*m) * np.sum(np.square(predictions-y))
    return cost
```



# **Gradient descent to minimize $J(\theta)$**

**Mini Batch Gradient Descent(SGD)**



# Mini-Batch Gradient Descent

Given training set:

$x^{(1)}$	$x^{(2)}$			$\dots$		$x^{(i-1)}$	$x^{(i)}$	$x^{(i+1)}$				$\dots$		$x^{(m)}$
$y^{(1)}$	$y^{(2)}$			$\dots$		$y^{(i-1)}$	$y^{(i)}$	$y^{(i+1)}$				$\dots$		$y^{(m)}$

first mini-batch                    second mini-batch                     $s^{th}$  mini-batch

$x^{(1)}$	$x^{(2)}$	$\dots$	$x^{(b)}$	$x^{(b+1)}$	$x^{(b+2)}$	$\dots$	$x^{(2b)}$	...	$x^{(1)}$	$x^{(2)}$	$\dots$	$x^{(b)}$	$y^{(1)}$	$y^{(2)}$	$\dots$	$y^{(b)}$	$y^{(b+1)}$	$y^{(b+2)}$	$\dots$	$y^{(2b)}$	...	$x^{(m)}$	$y^{(m)}$
-----------	-----------	---------	-----------	-------------	-------------	---------	------------	-----	-----------	-----------	---------	-----------	-----------	-----------	---------	-----------	-------------	-------------	---------	------------	-----	-----------	-----------

# training example in each mini-batch:  $b$

$$\# \text{mini-batch: } s = \left[ \frac{m}{b} \right]$$

first mini-batch:  $x^{(1)}, x^{(2)}, \dots, x^{(b)}$

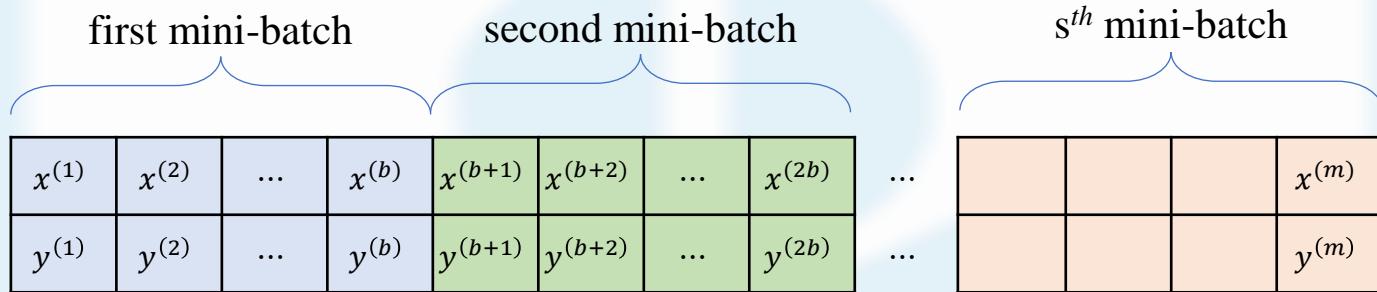
second mini-batch:  $x^{(b+1)}, x^{(b+2)}, \dots, x^{(2b)}$

$\vdots$

$k^{th}$  mini-batch:  $x^{((k-1)b+1)}, x^{((k-1)b+2)}, \dots, x^{(kb)}$



# Mini-Batch Gradient Descent



**Repeat {**

*Randomly shuffle(reorder) examples in training set*

*partition new training set into  $s = \lceil \frac{m}{b} \rceil$  mini-batches of size  $b$*

**for**  $k=1$  to  $s$  {

$$\text{set } X^{\{k\}} = \{x^{((k-1)b+1)}, x^{((k-1)b+2)}, \dots, x^{(kb)}\}$$

$$\text{set } Y^{\{k\}} = \{y^{((k-1)b+1)}, y^{((k-1)b+2)}, \dots, y^{(kb)}\}$$

$$\theta_0 := \theta_0 - \alpha \frac{1}{b} \sum_{i=1}^b (h_\theta(X^{\{k\}}(i)) - Y^{\{k\}}(i))$$

$$\theta_j := \theta_0 - \alpha \frac{1}{b} \sum_{i=1}^b (h_\theta(X^{\{k\}}(i)) - Y^{\{k\}}(i)) X^{\{k\}}(i, j)$$

(simultaneously update for every  $j = 0, 1, \dots, n$ )

}

}

$$X^{\{k\}}(i, j) := x_j^{((k-1)b+i)}$$
$$Y^{\{k\}}(i) := y^{((k-1)b+i)}$$



# Mini-Batch Gradient Descent

```
def mini_batch_gradient_descent(X, y, theta, learning_rate, iterations, batch_size):
    m = len(y)
    cost_history = np.zeros(iterations)
    n_batches = int(m/batch_size)
    for it in range(iterations):
        cost = 0.0
        indices = np.random.permutation(m)
        X = X[indices]
        y = y[indices]
        for i in range(0, m, batch_size):
            X_i = X[i:i+batch_size]
            Y_i = y[i:i+batch_size]
            X_i = np.c_[np.ones(len(X_i)), X_i]
            prediction = np.dot(X_i,theta)
            theta = theta -(1/m)*learning_rate*( X_i.T.dot((prediction - Y_i)))
            cost += cal_cost(theta,X_i,Y_i)
        cost_history[it] = cost
    return theta, cost_history

# -----#
def cal_cost(theta, X, y):
    """ Calculates the cost for given theta, X and Y. """
    m = len(y)
    predictions = X.dot(theta)
    cost = 1/(2*m) * np.sum(np.square(predictions-y))
    return cost
```

# Gradient Descent

Batch Gradient Descent	Stochastic Gradient Descent	Mini-Batch Gradient Descent
Since entire training data is considered before taking a step in the direction of gradient, therefore it takes a lot of time for making a single update.	Since only a single training example is considered before taking a step in the direction of gradient, we are forced to loop over the training set and thus cannot exploit the speed associated with vectorizing the code.	Since a subset of training examples is considered, it can make quick updates in the model parameters and can also exploit the speed associated with vectorizing the code.
It makes <i>smooth</i> updates in the model parameters	It makes <i>very noisy</i> updates in the parameters	Depending upon the batch size, the updates can be made <i>less noisy</i> – greater the batch size less noisy is the update

Thus, mini-batch gradient descent makes a compromise between the speedy convergence and the noise associated with gradient update which makes it a more flexible and robust algorithm.

# References

Andrew Ng, <https://www.coursera.org/learn/machine-learning>

<https://www.geeksforgeeks.org/gradient-descent-algorithm-and-its-variants/>

<https://www.geeksforgeeks.org/ml-mini-batch-gradient-descent-with-python/>

[https://www.holehouse.org/mlclass/04\\_Linear\\_Regression\\_with\\_multiple\\_variables.html](https://www.holehouse.org/mlclass/04_Linear_Regression_with_multiple_variables.html)

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

$$J^* = \min_{\theta} J(\theta)$$

$$J^* = J(\theta^*)$$