# MACHINE LEARNING COURSE

PRESENTED BY ABDEL RAHMAN ALSABBAGH

LECTURE #3 – SAT - 17.5.2023

In the name of Allah, the most gracious, the most merciful, we start :)

# Today's Quote

"Growth happens without requiring much effort"

- Chris Gardener.
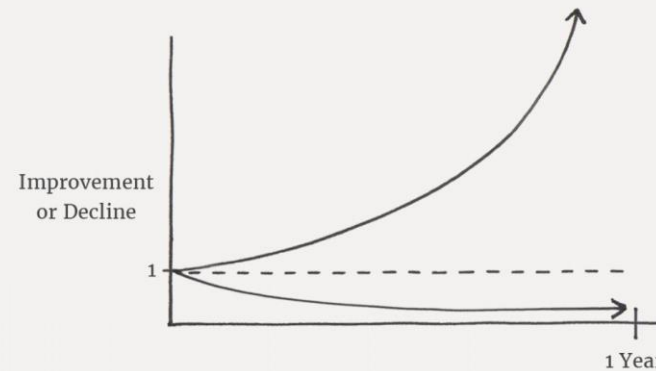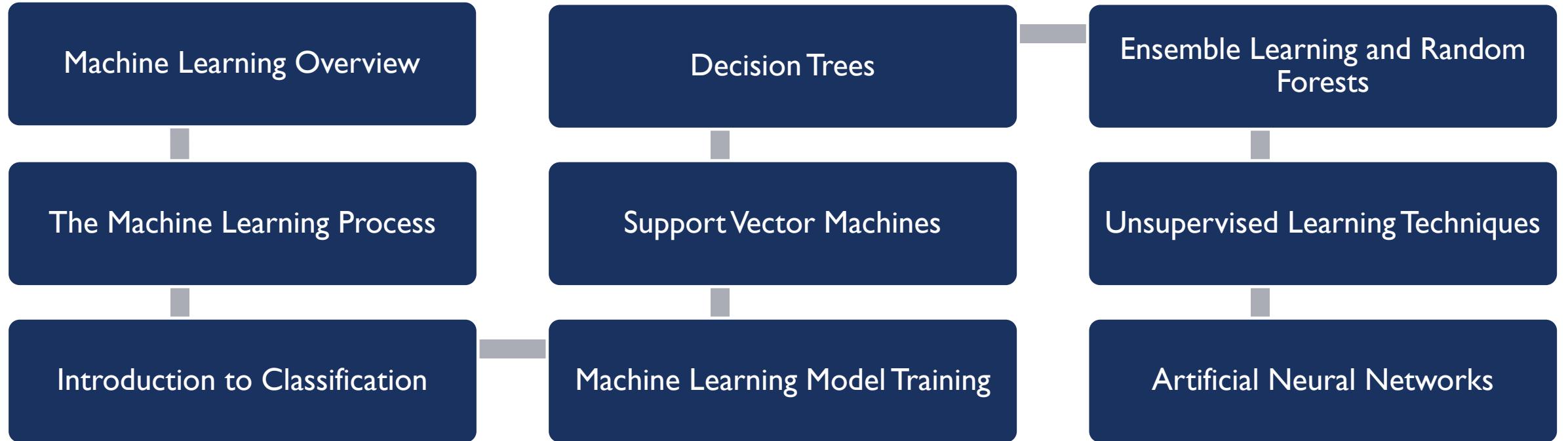
# Today's Quote

# Course Outline

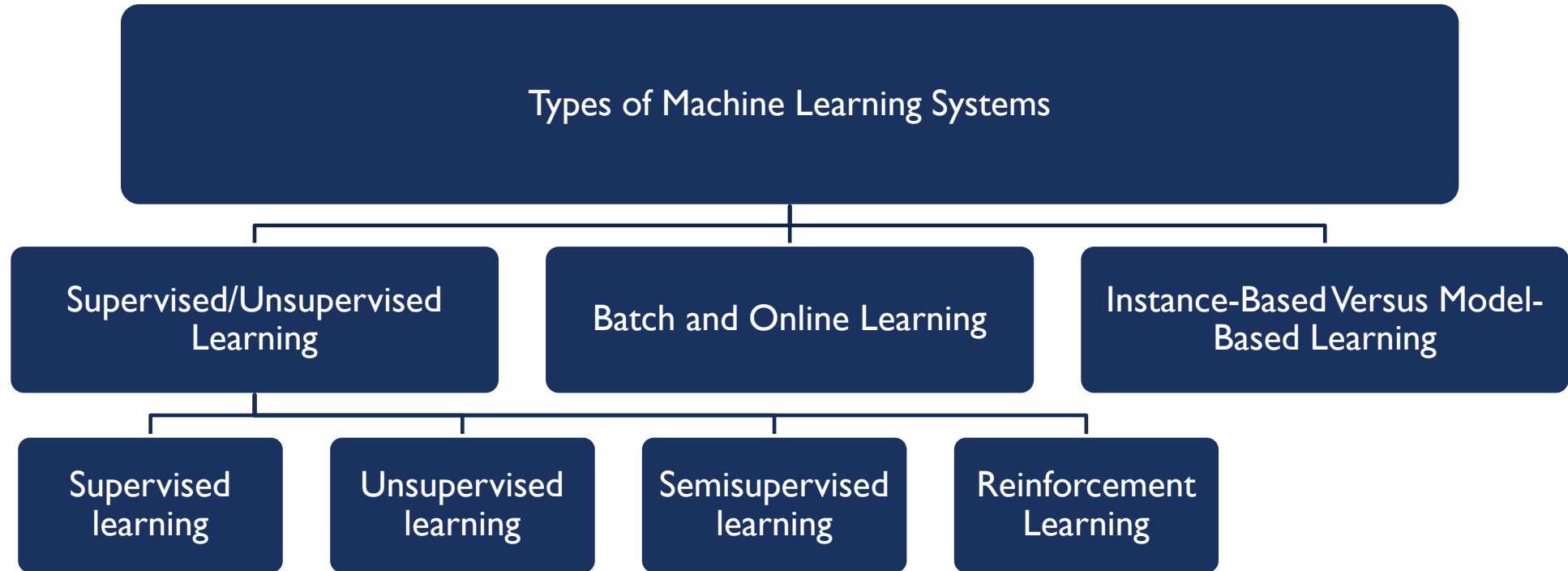| | | |
|---|---|---|
| Machine Learning Overview | Decision Trees | Ensemble Learning and Random Forests |
| The Machine Learning Process | Support Vector Machines | Unsupervised Learning Techniques |
| Introduction to Classification | Machine Learning Model Training | Artificial Neural Networks |

Resources used:
- Hands-on Machine Learning with Scikit-Learn, Keras, and TensorFlow by Aurélien Géron
- Machine Learning Specialization by Andrew Ng and Stanford Online

Machine Learning Course
Abdel Rahman AlSabbagh

# Types of Machine Learning Systems

# Recap

**Main Challenges of Machine Learning**

**Bad Data**

**Bad Algorithm**

| Insufficient Quantity of Training Data | Nonrepresentative Training Data | Poor-Quality Data | Irrelevant Features | Overfitting the Training Data | Underfitting the Training Data |

# Recap

Testing and Validating

Hyperparameter Tuning and Model Selection

Data Mismatch

Machine Learning Course
Abdel Rahman AlSabbagh

# The Machine Learning Process

- Linear regression.
- Loss and cost functions.
- Visualizing the cost function.
- Gradient descent.
- Learning rate.
- Gradient descent for linear regression.
- Evaluation metrics.
- Linear regression for multiple features.
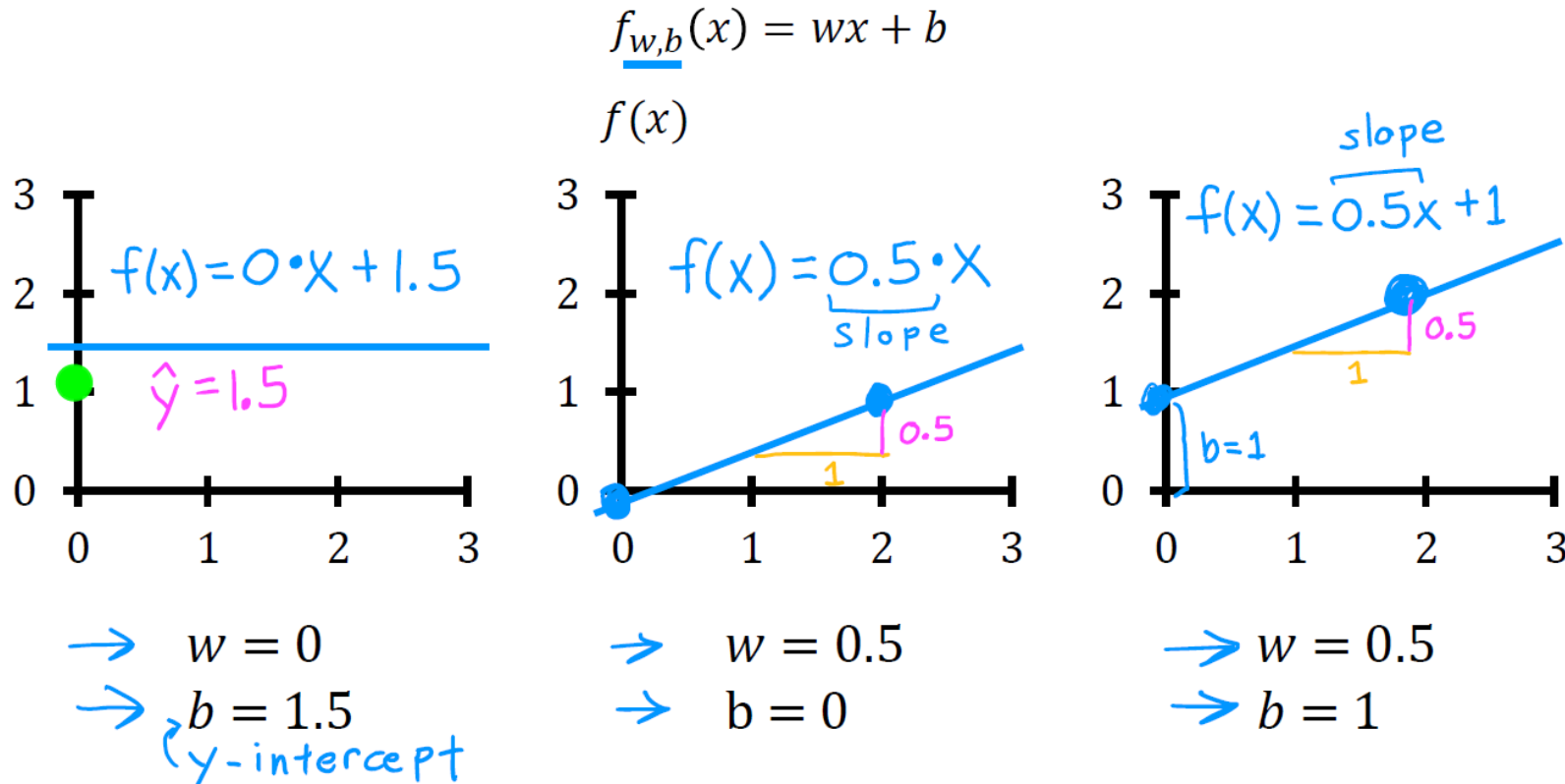- Gradient descent for multiple features.
- Feature scaling.

Source: Machine Learning Specialization by Andrew Ng and Stanford Online.
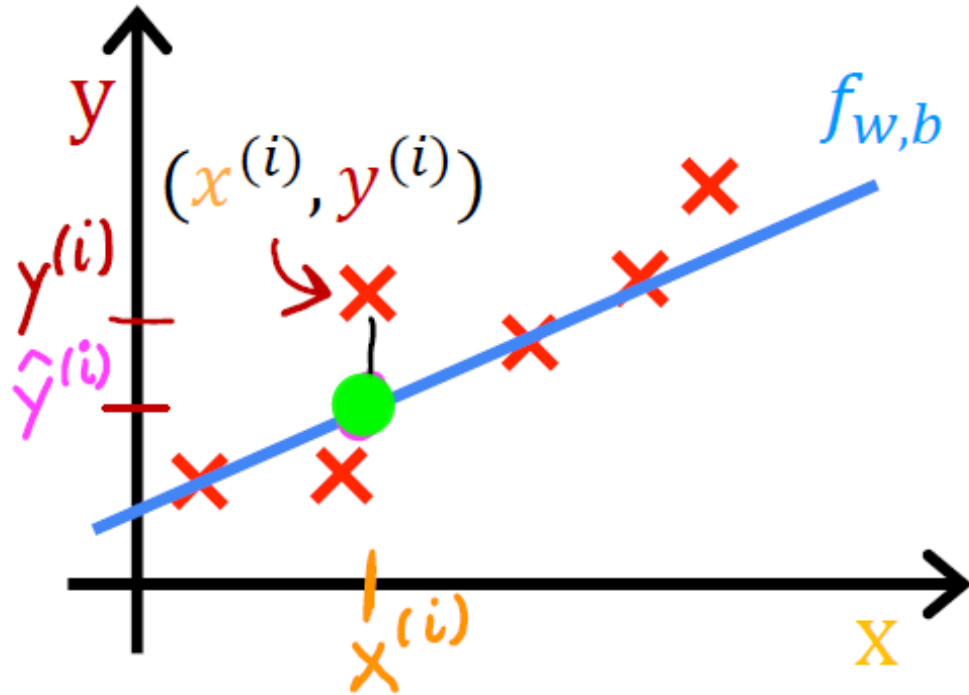
# Linear Regression

$$\text{Model: } f_{w,b}(x) = wx + b$$

$w, b$ : parameters

coefficients
weights

IEEE
Computational
Intelligence
Society®

University of Jordan Chapter

# Linear Regression

$$f_{w,b}(x) = wx + b$$

$f(x)$



Left graph:
$f(x) = 0 \cdot x + 1.5$
$\hat{y} = 1.5$
$\to w = 0$
$\to b = 1.5$ (y-intercept)

Middle graph:
$f(x) = 0.5 \cdot x$ (slope)
$\to w = 0.5$
$\to b = 0$

Right graph:
slope
$f(x) = 0.5x + 1$
$b = 1$
$\to w = 0.5$
$\to b = 1$

IEEE Computational Intelligence Society®
University of Jordan Chapter

# Linear Regression



$$\hat{y}^{(i)} = f_{w,b}\left(x^{(i)}\right) \Leftarrow$$

$$f_{w,b}\left(x^{(i)}\right) = wx^{(i)} + b$$

# Loss and Cost Functions

$$\hat{y}^{(i)} = f_{w,b}(x^{(i)}) \quad \leftarrow$$

$$f_{w,b}(x^{(i)}) = wx^{(i)} + b$$

Find $w, b$:
$\hat{y}^{(i)}$ is close to $y^{(i)}$ for all $(x^{(i)}, y^{(i)})$.

Loss function

$$\left(f_{w,b}(x^{(i)}) - y^{(i)}\right)^2$$

Cost function

m = number of training examples

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^{m} \left(f_{w,b}(x^{(i)}) - y^{(i)}\right)^2$$

intuition (next!)

IEEE
Computational
Intelligence
Society®
University of Jordan Chapter

Machine Learning Course
Abdel Rahman AlSabbagh

# Cost Function

model:

$$f_{w,b}(x) = wx + b$$

parameters:

$$w, b$$



cost function:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)})^2$$
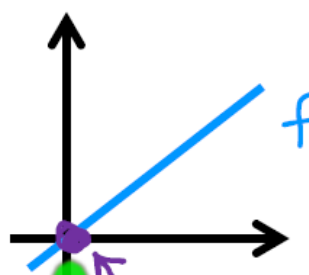
goal:

$$\underset{w,b}{\text{minimize}} \, J(w, b)$$

IEEE
Computational
Intelligence
Society®
University of Jordan Chapter

Machine Learning Course
Abdel Rahman AlSabbagh

# Cost Function

simplified

$$f_w(x) = wx \qquad b = \emptyset$$

$w$



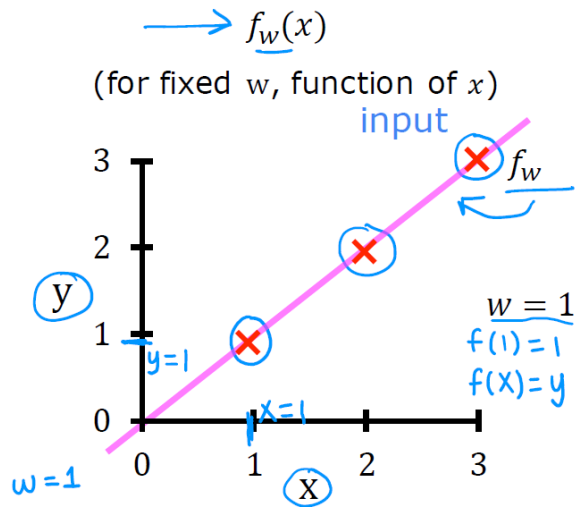$$J(w) = \frac{1}{2m} \sum_{i=1}^{m} (f_w(x^{(i)}) - y^{(i)})^2$$

$$wx^{(i)}$$

$$\underset{w}{\text{minimize}}\, J(w)$$

# Cost Function

# Cost Function



$f_w(x)$

(for fixed $w$, function of $x$)
input

$y$

$x$

$w = 1$

$f_w$

$w = 1$
$f(1) = 1$
$f(x) = y$

$y = 1$

$x = 1$

(function of $w$)
parameter

$J(w)$

$J(1) = 0$

?

$w$

$w = 1$
$\downarrow$

$J(w) = \dfrac{1}{2m} \sum_{i=1}^{m} (f_w(x^{(i)}) - y^{(i)})^2$

$w x^{(i)}$

$= \dfrac{1}{2m} \sum_{i=1}^{m} (wx^{(i)} - y^{(i)})^2$

$\varnothing$

$= \dfrac{1}{2m} ( 0^2 + 0^2 + 0^2 ) = 0$

# Cost Function

$f_w(x)$

(function of $x$)



$J(w)$

(function of $w$)



$$J(0.5) = \frac{1}{2m}\left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2\right] = \frac{1}{2\times3}[3.5] = \frac{3.5}{6} \simeq 0.58$$

Machine Learning Course

Abdel Rahman AlSabbagh

# Cost Function



$f_w(x)$

(function of $x$)

$\omega=1$

$f_w$

$f$

$w=0$

$f(x)=-0.5x$

$$J(0) = \frac{1}{2m}\left(1^2+2^2+3^2\right)=\frac{1}{6}[14]\approx 2.3$$

$J(w)$

(function of $w$)

5.25

$\omega=1$

$J$

$J(w)$

$w$

-0.5  0  0.5  1  1.5  2  2.5

$w$

how to choose $w$?
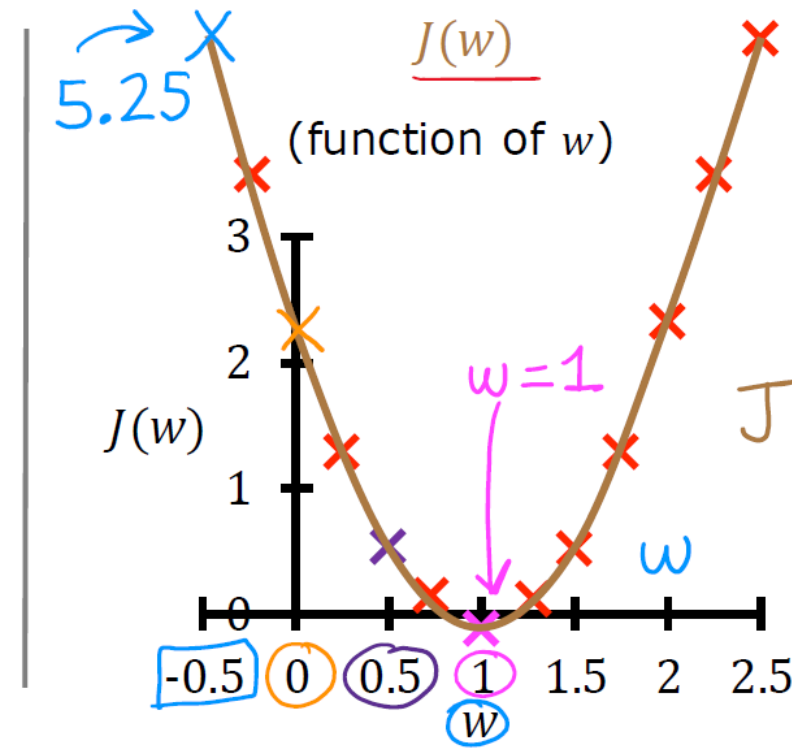
Machine Learning Course

Abdel Rahman AlSabbagh

# Cost Function

goal of linear regression:

$$\underset{w}{\text{minimize}}\, J(w)$$

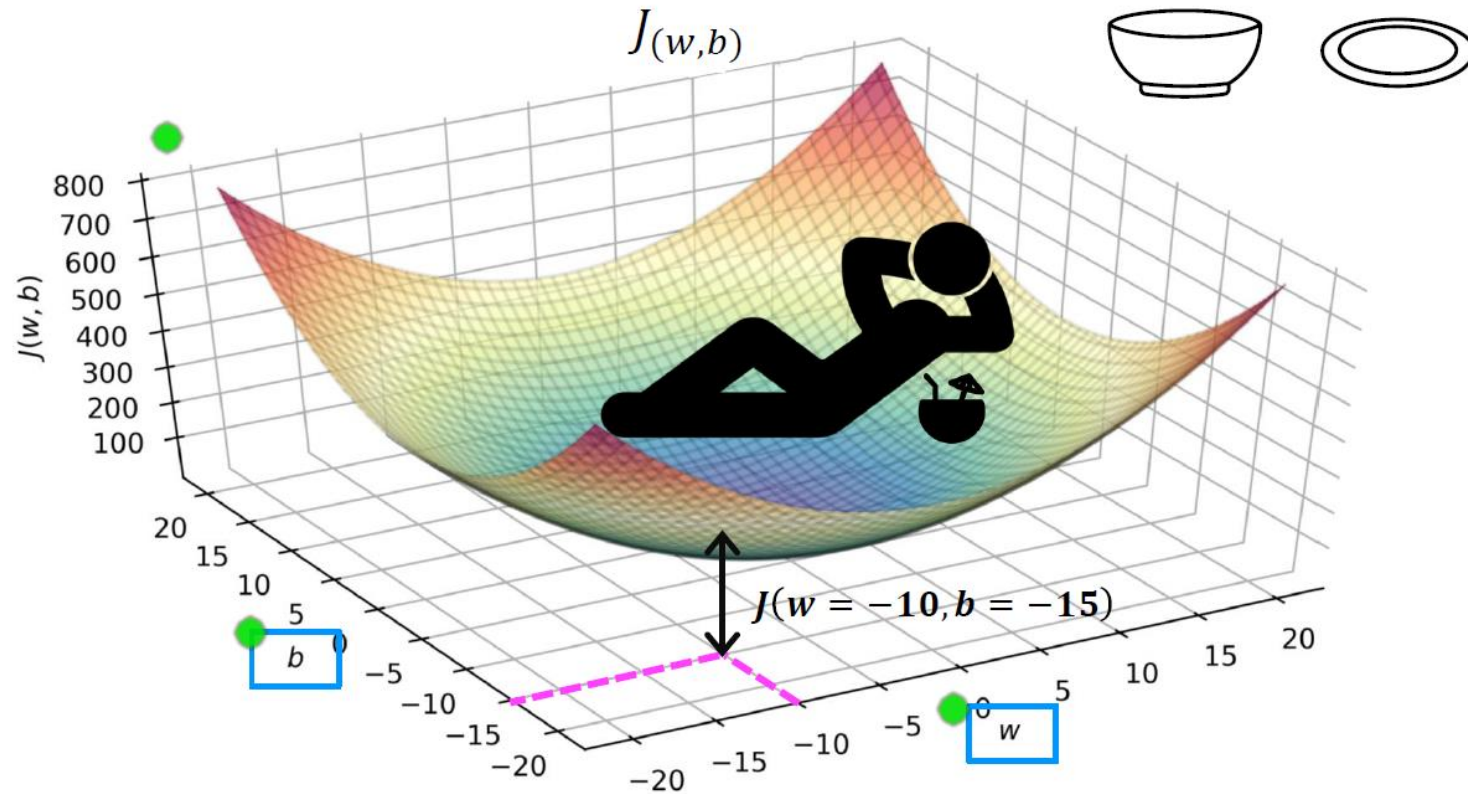general case:

$$\underset{w,b}{\text{minimize}}\, J(w,b)$$



$J(w)$

(function of $w$)

$J(w)$

$w=1$

choose w to minimize J(w)

# Visualizing the Cost Function

$f_{w,b}$

(function of $x$)



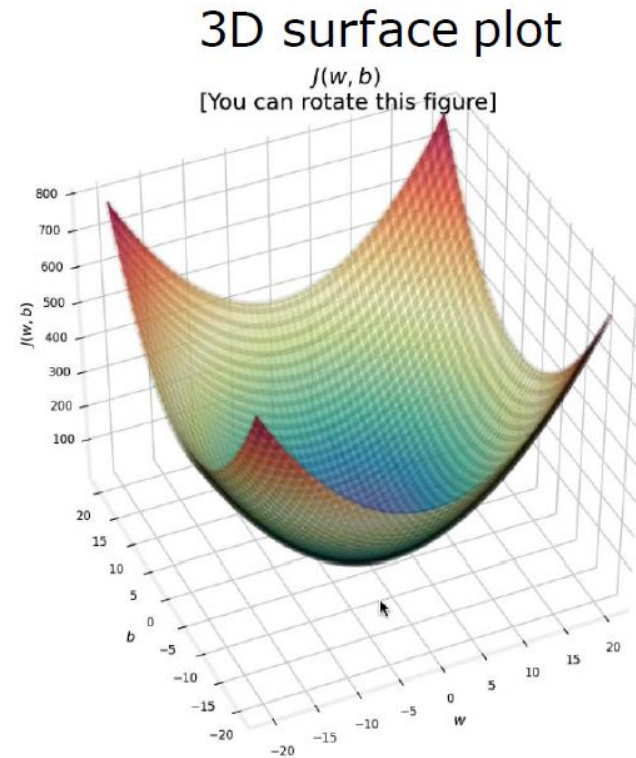$f_{w,b}(x) = 0.06x + 50$

$J$

(function of $w, b$)
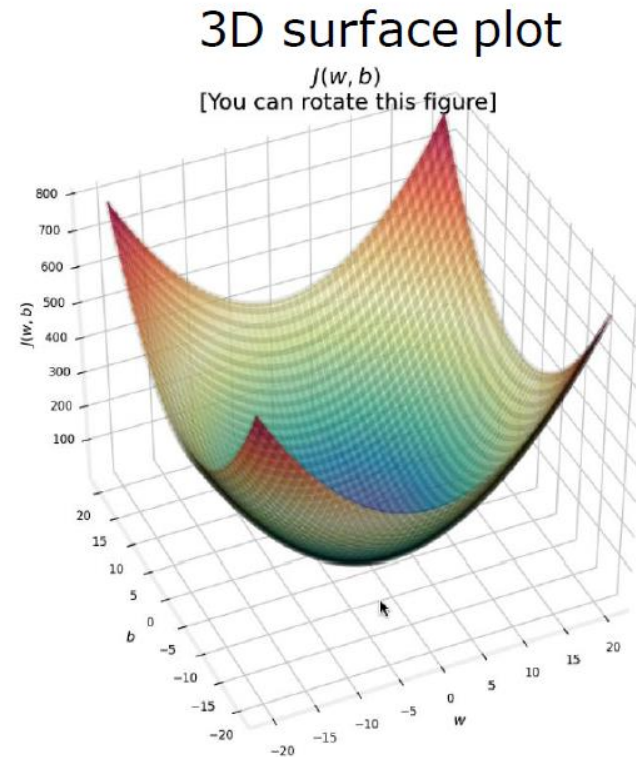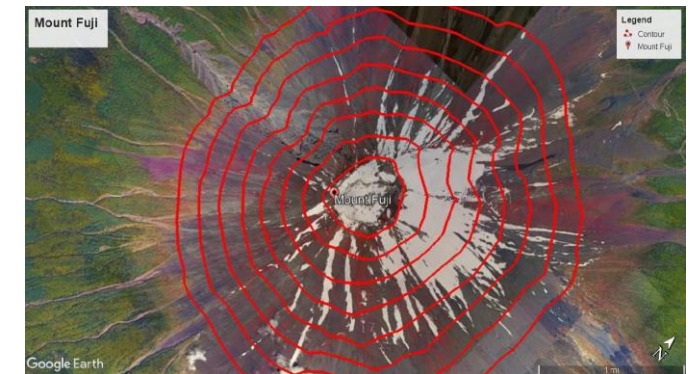


$w, b$

# Visualizing the Cost Function

# Visualizing the Cost Function

# Visualizing the Cost Function



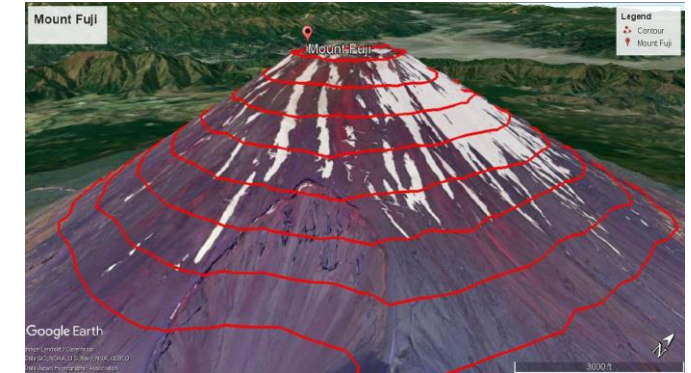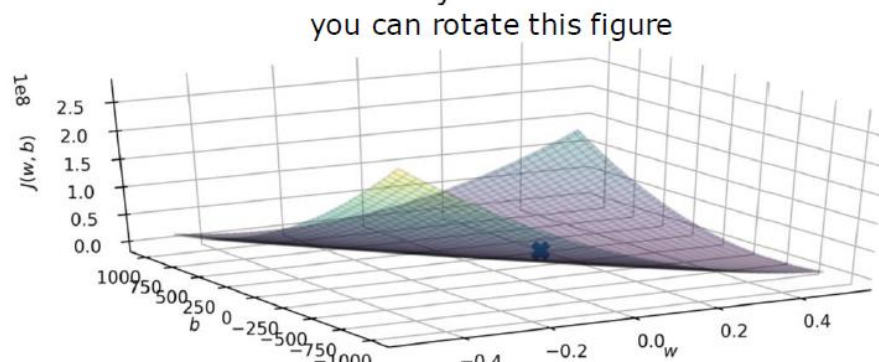3D surface plot

$J(w, b)$
[You can rotate this figure]

IEEE
Computational
Intelligence
Society®

University of Jordan Chapter

# Visualizing the Cost Function

# Gradient Descent

Have some function  $J(w, b)$  *for linear regression or any function*
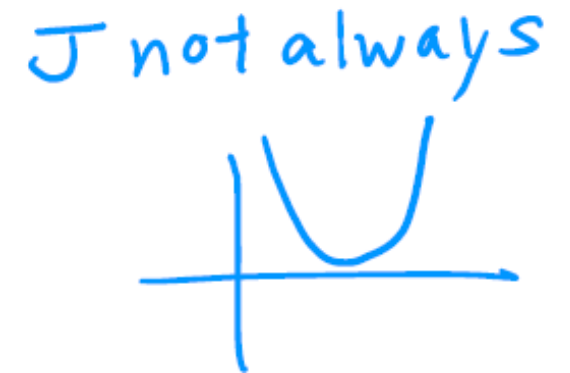
Want  $\min_{w,b} J(w, b)$    $\min_{w_1, \ldots, w_n, b} J(w_1, w_2, \ldots, w_n, b)$
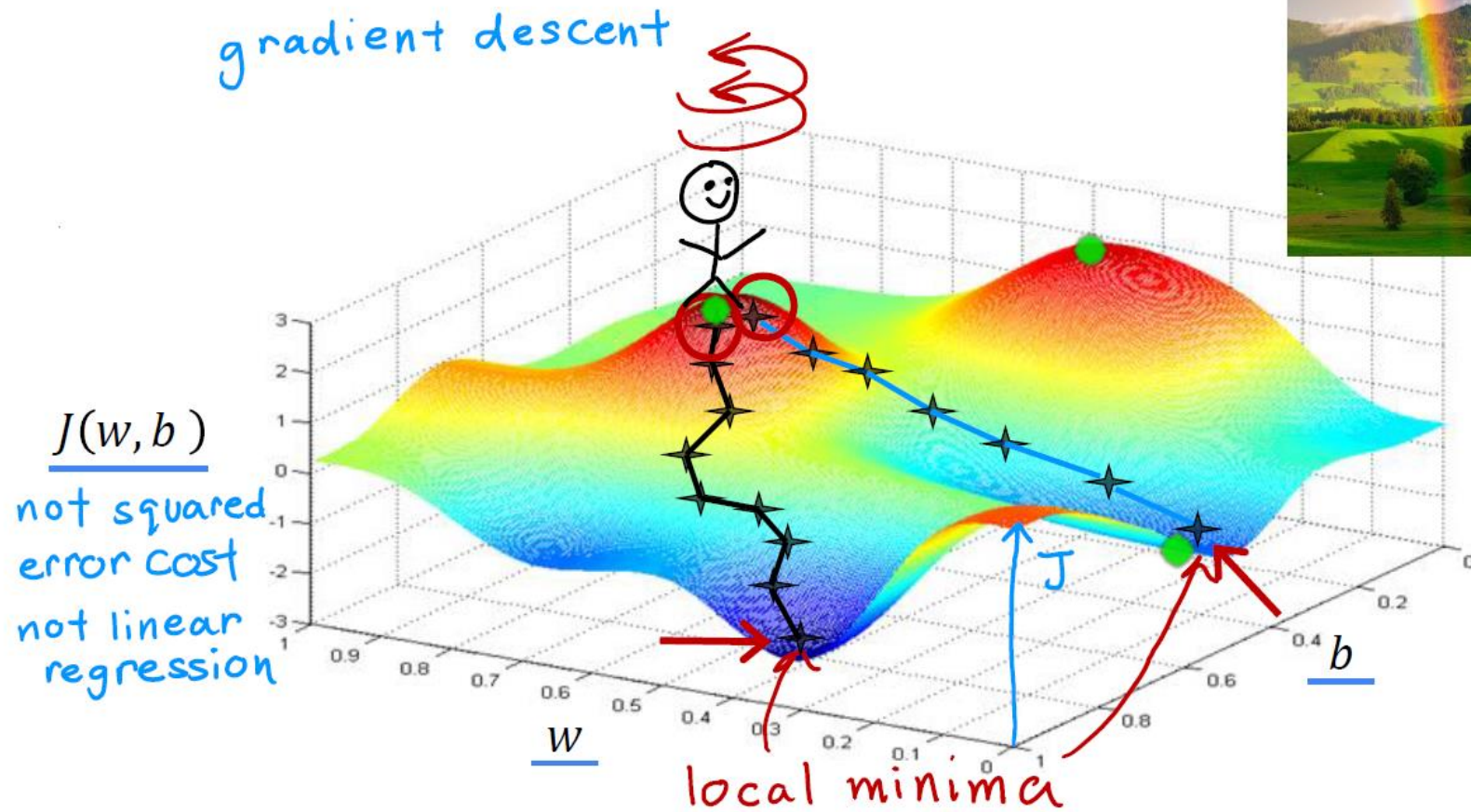
Outline:

    Start with some $w, b$    *(set $w=0, b=0$)*

    Keep changing $w, b$ to reduce $J(w, b)$

    Until we settle at or near a minimum

*$J$ not always*

# Gradient Descent

# Gradient Descent Algorithm

Repeat until convergence

$$w = w - \alpha \boxed{\frac{d}{dw} J(w,b)}$$

$$b = b - \alpha \frac{d}{db} J(w,b)$$

Learning rate
Derivative

Simultaneously update w and b

$$tmp\_w = w - \alpha \frac{\partial}{\partial w} J(w,b)$$

$$tmp\_b = b - \alpha \frac{\partial}{\partial b} J(w,b)$$

$$w = tmp\_w$$

$$b = tmp\_b$$

# Gradient Descent Algorithm



$$w = w - \alpha \frac{d}{dw} J(w)$$
$$> 0$$

$$w = w - \underline{\alpha} \cdot (positive\ number)$$

$$\frac{d}{dw} J(w) \quad < 0$$

$$w = w - \alpha \cdot (negative\ number)$$

J(w)

J(w)

2

negative slope

-2

1

# Gradient Descent Algorithm

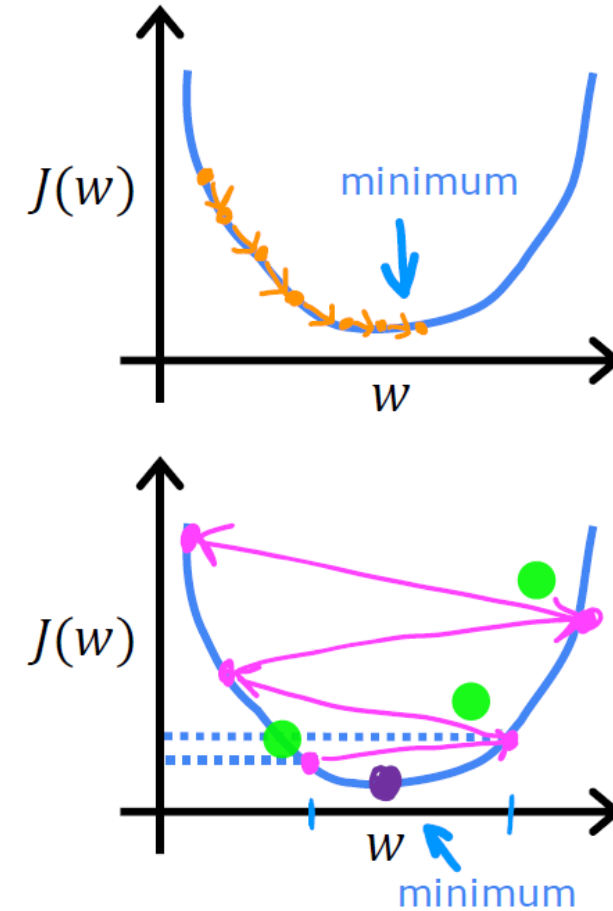$$w = w - \textcircled{\alpha} \frac{d}{dw} J(w)$$

🟢
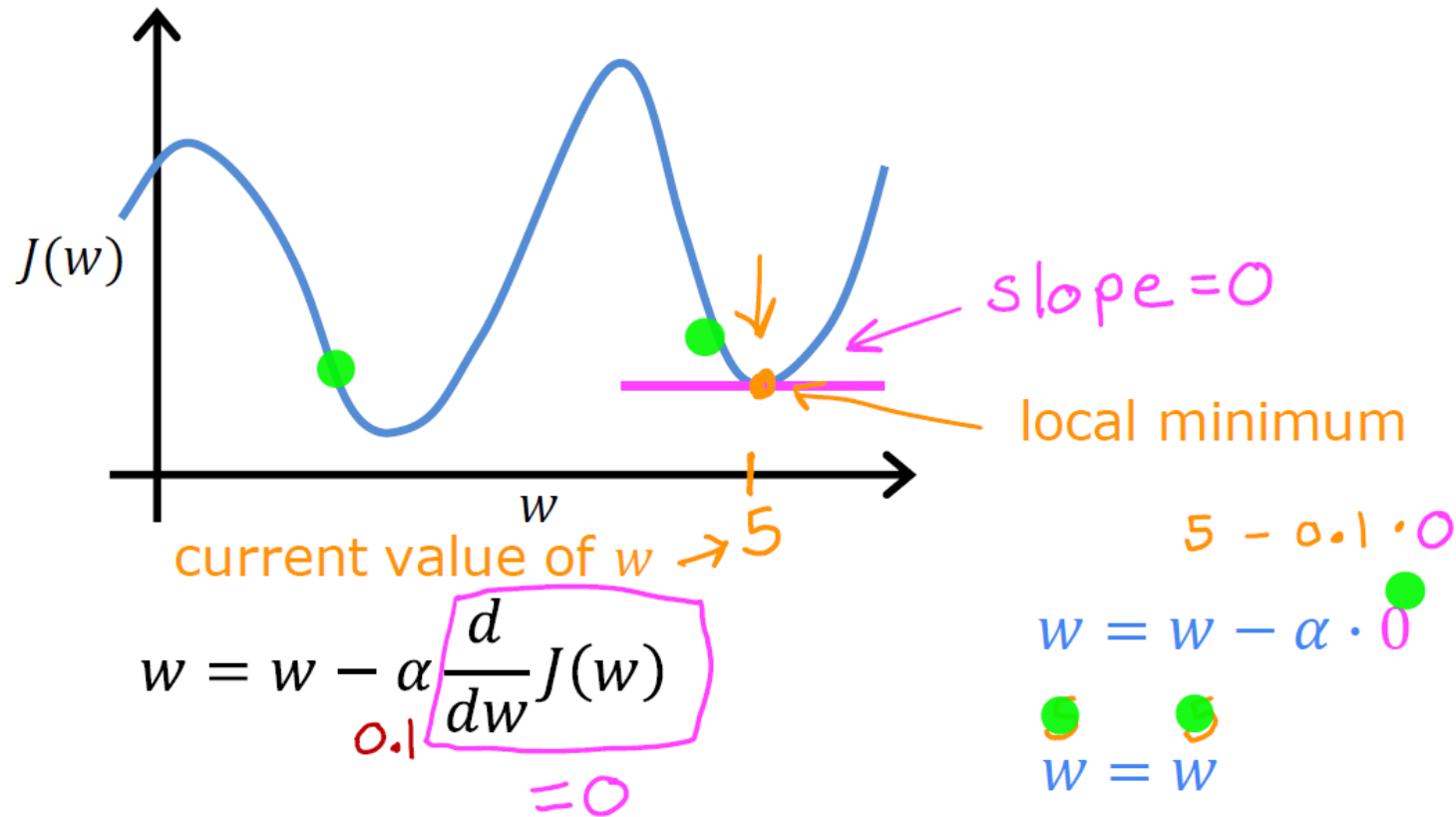
If $\alpha$ is too small...
Gradient descent may be slow.

If $\alpha$ is too large...

Gradient descent may:
 - Overshoot, never reach minimum
 - Fail to converge, diverge

# Gradient Descent Algorithm



$J(w)$

slope = 0

local minimum

current value of $w$ → 5

$$w = w - \alpha \frac{d}{dw} J(w)$$

0.1 = 0

$5 - 0.1 \cdot 0$

$w = w - \alpha \cdot 0$

$w = w$

# Gradient Descent for Linear Regression

Linear regression model

$$f_{w,b}(x) = wx + b$$

Cost function

$$J(w,b) = \frac{1}{2m}\sum_{i=1}^{m}(f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Gradient descent algorithm

repeat until convergence {

$$w = w - \alpha\frac{\partial}{\partial w}J(w,b) \longrightarrow \frac{1}{m}\sum_{i=1}^{m}(f_{w,b}(x^{(i)}) - y^{(i)})x^{(i)}$$

$$b = b - \alpha\frac{\partial}{\partial b}J(w,b) \longrightarrow \frac{1}{m}\sum_{i=1}^{m}(f_{w,b}(x^{(i)}) - y^{(i)})$$

}

# Evaluation Metrics for Linear Regression

Mean Absolute Error

$$\text{MAE} = \frac{1}{n} \sum_{j=1}^{n} |y_j - \hat{y}_j|$$

Mean Squared Error

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

Root Mean Squared Error

$$RMSE = \sqrt{\sum_{i=1}^{n} \frac{(\hat{y}_i - y_i)^2}{n}}$$

R Squared

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \overline{y})^2}$$

# Linear Regression for Multiple Features

Model:

Previously: $f_{w,b}(x) = wx + b$

$$f_{w,b}(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b$$

example

$$f_{w,b}(x) = 0.1 x_1 + 4 x_2 + 10 x_3 + -2 x_4 + 80$$

↑ size   ↑ #bedrooms   ↑ #floors   ↑ years   ↑ base price

$$f_{w,b}(x) = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b$$

# Linear Regression for Multiple Features

$$\text{vector} \quad \vec{x} = \begin{bmatrix} x_1 & x_2 & x_3 \dots x_n \end{bmatrix}$$

$$\vec{w} = \begin{bmatrix} w_1 & w_2 & w_3 \dots w_n \end{bmatrix}$$

$$b$$

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + b$$

dot product

multiple linear regression

Machine Learning Course

Abdel Rahman AlSabbagh

# Vectorization

```
w = np.array([1.0,2.5,-3.3])
b = 4              X[0] X[1] X[2]
x = np.array([10,20,30])
```

```
f = 0          range(n)
for j in range(0,n):
        f = f + w[j] * x[j]
f = f + b
```

$$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

```
f = w[0] * x[0] +
    w[1] * x[1] +
    w[2] * x[2] + b
```

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

```
f = np.dot(w,x) + b
```

Machine Learning Course

Abdel Rahman AlSabbagh

IEEE
Computational
Intelligence
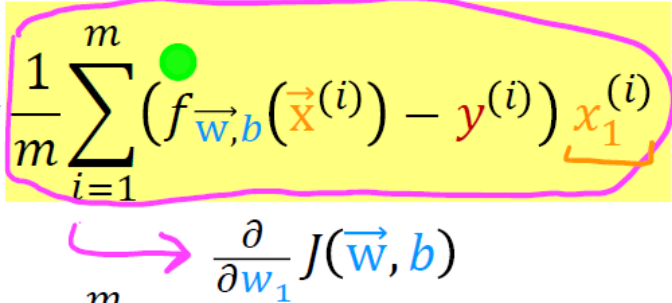Society®

University of Jordan Chapter

# Gradient Descent for Multiple Features

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

}

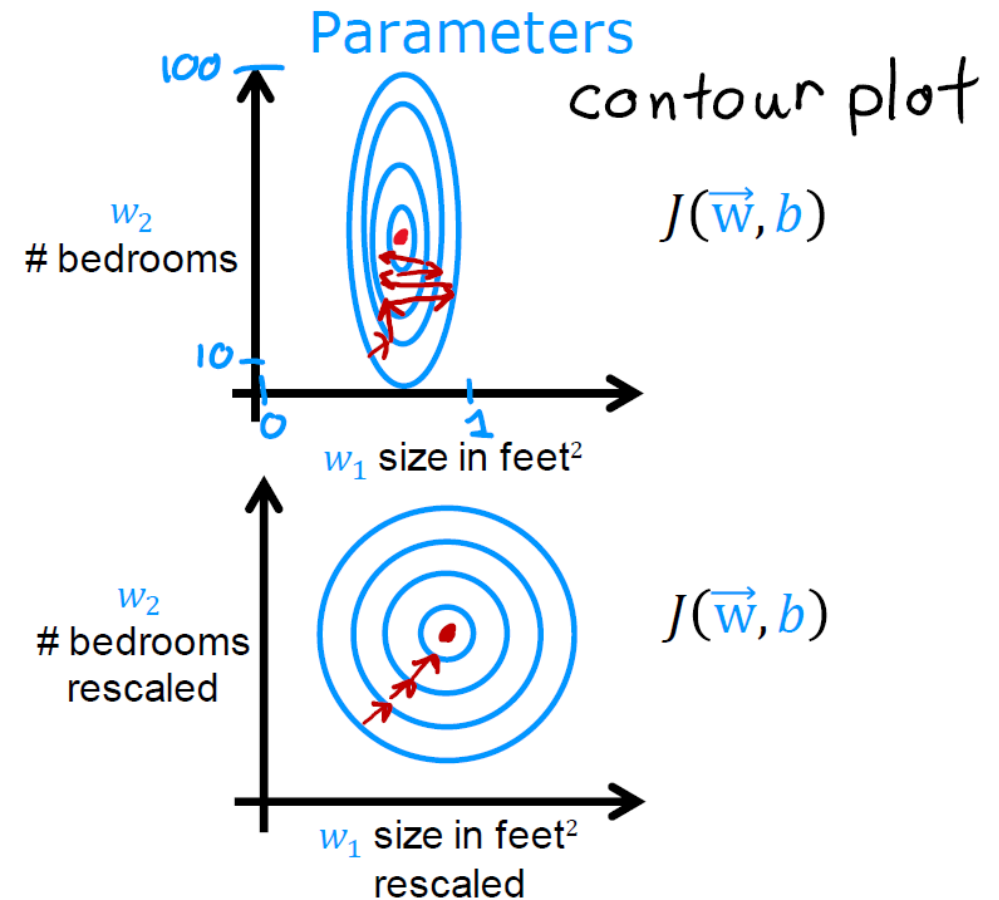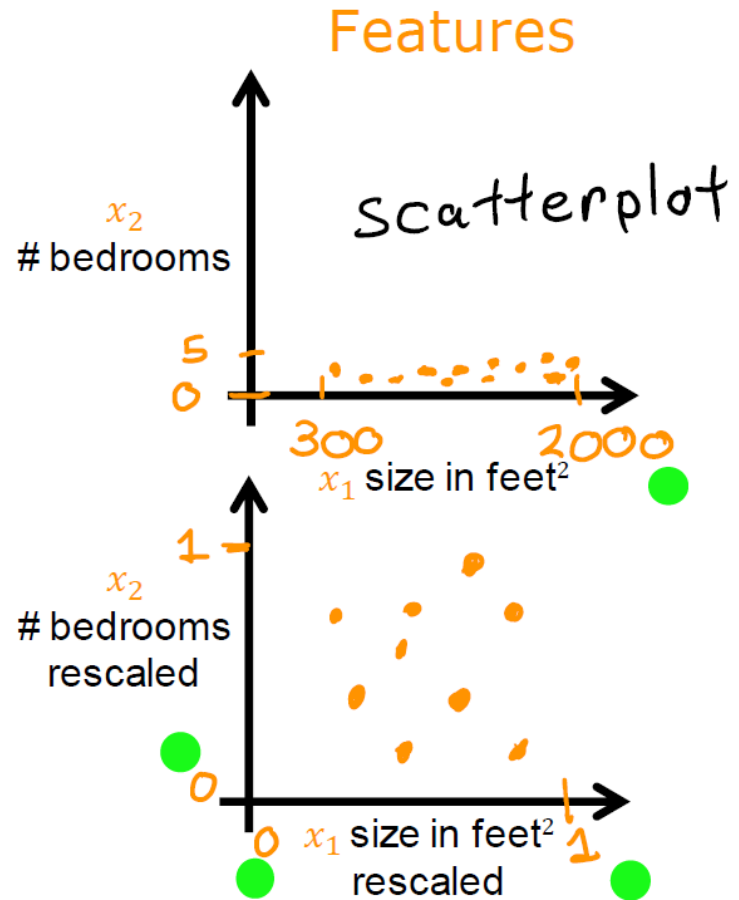# Gradient Descent for Multiple Features

$n$ features $(n \geq 2)$

repeat {

$j=1$

$w_1 = w_1 - \alpha \dfrac{1}{m} \displaystyle\sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right) x_1^{(i)}$

$\underbrace{\qquad}$

$\vdots$

$\dfrac{\partial}{\partial w_1} J(\vec{w}, b)$

$j=n$

$w_n = w_n - \alpha \dfrac{1}{m} \displaystyle\sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right) x_n^{(i)}$

$b = b - \alpha \dfrac{1}{m} \displaystyle\sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right)$

}

simultaneously update
$w_j$ (for $j = 1, \cdots, n$) and $b$

# The Effect of Feature Scaling

# The Effect of Feature Scaling

aim for about $-1 \leq x_j \leq 1$ for each feature $x_j$

$$-3 \leq x_j \leq 3$$
$$-0.3 \leq x_j \leq 0.3$$

$\Big\}$ acceptable ranges

$0 \leq x_1 \leq 3$     okay, no rescaling

$-2 \leq x_2 \leq 0.5$     okay, no rescaling
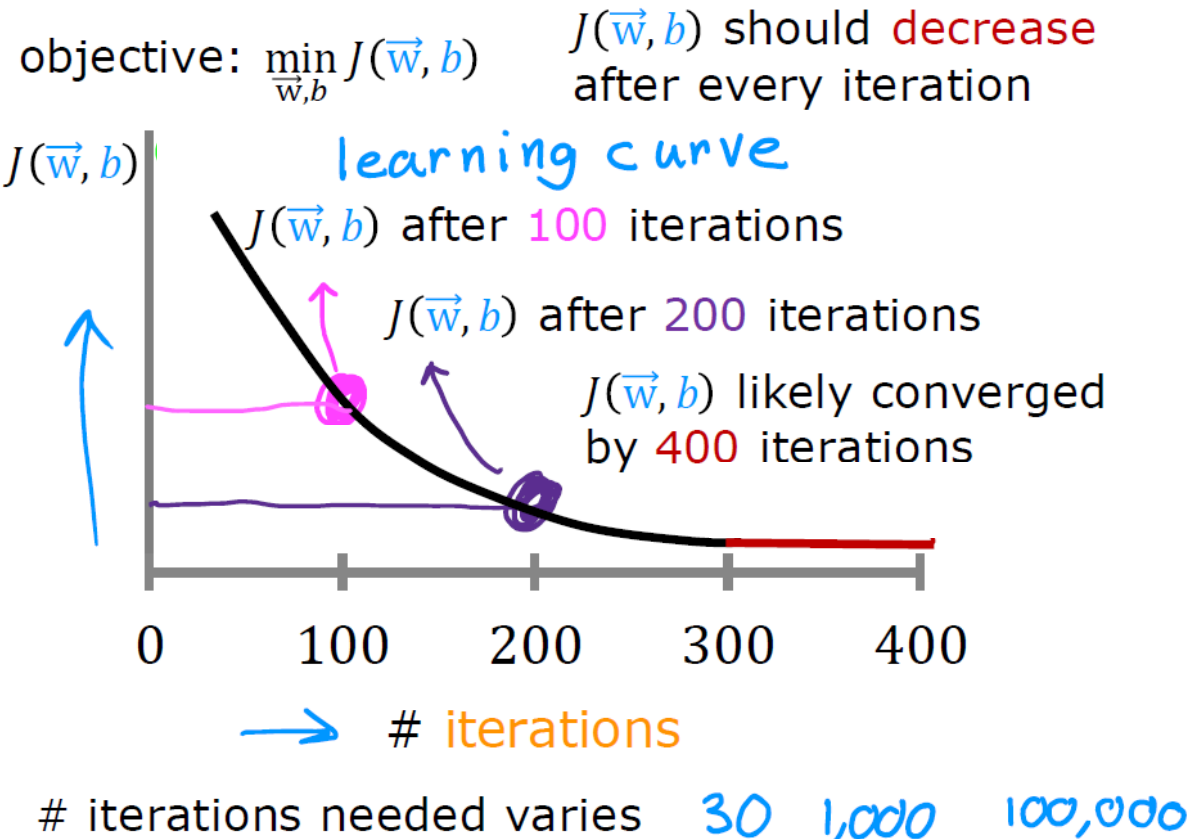
$-100 \leq x_3 \leq 100$     too large → rescale

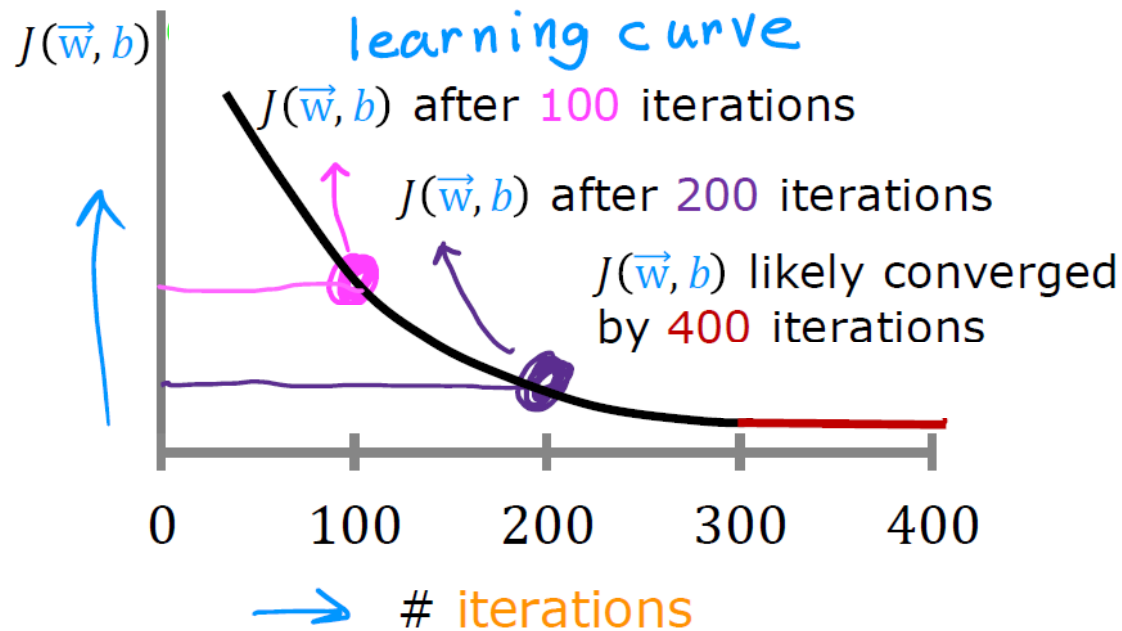$-0.001 \leq x_4 \leq 0.001$     too small → rescale

$98.6 \leq x_5 \leq 105$     too large → rescale

# Making Sure Gradient Descent is Working

# Checking Gradient Descent for Convergence

objective: $\min_{\vec{w},b} J(\vec{w}, b)$

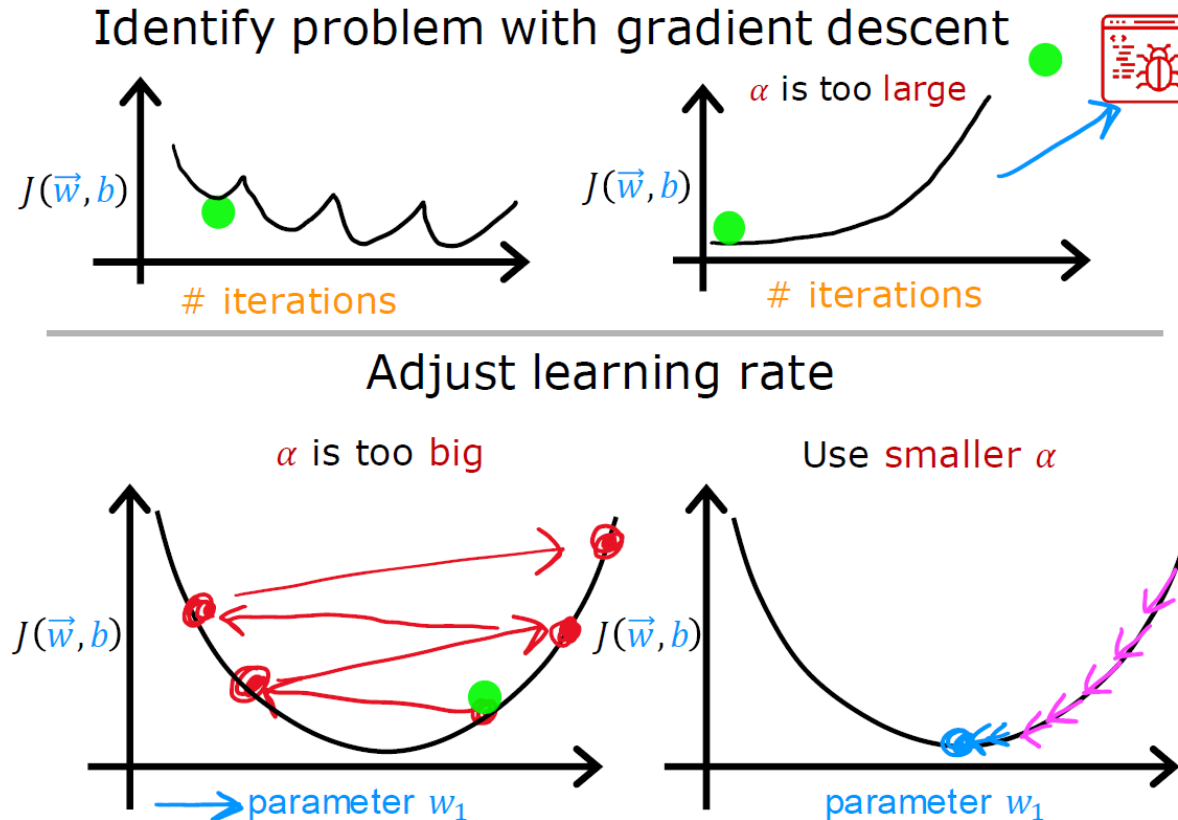$J(\vec{w}, b)$ should **decrease** after every iteration

**learning curve**

$J(\vec{w}, b)$ after **100** iterations

$J(\vec{w}, b)$ after **200** iterations

$J(\vec{w}, b)$ likely converged by **400** iterations

# iterations needed varies  30  1,000  100,000

Automatic convergence test

Let $\varepsilon$ "epsilon" be $10^{-3}$.

0.001

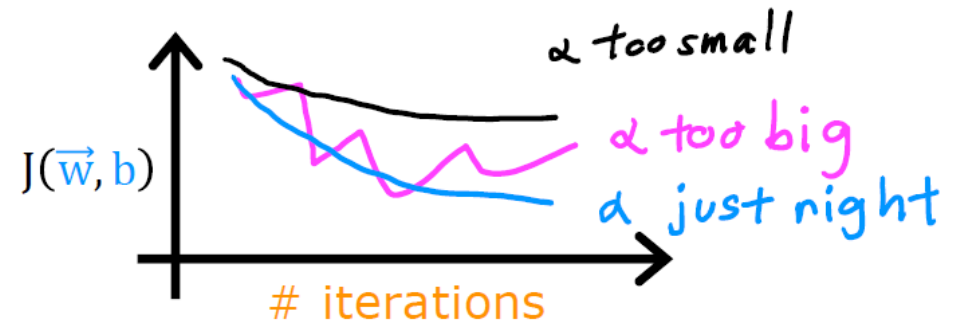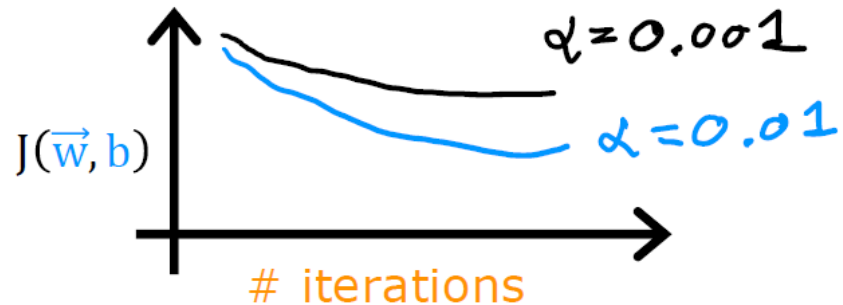If $J(\vec{w}, b)$ decreases by $\leq \varepsilon$ in one iteration, declare **convergence**.
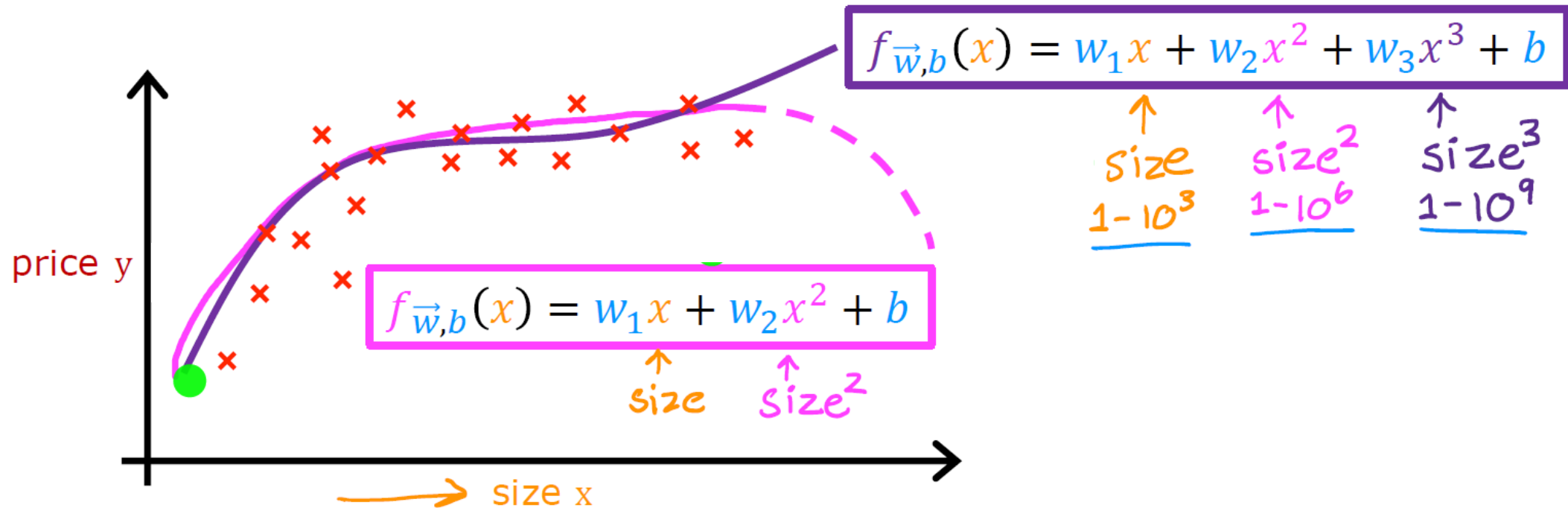
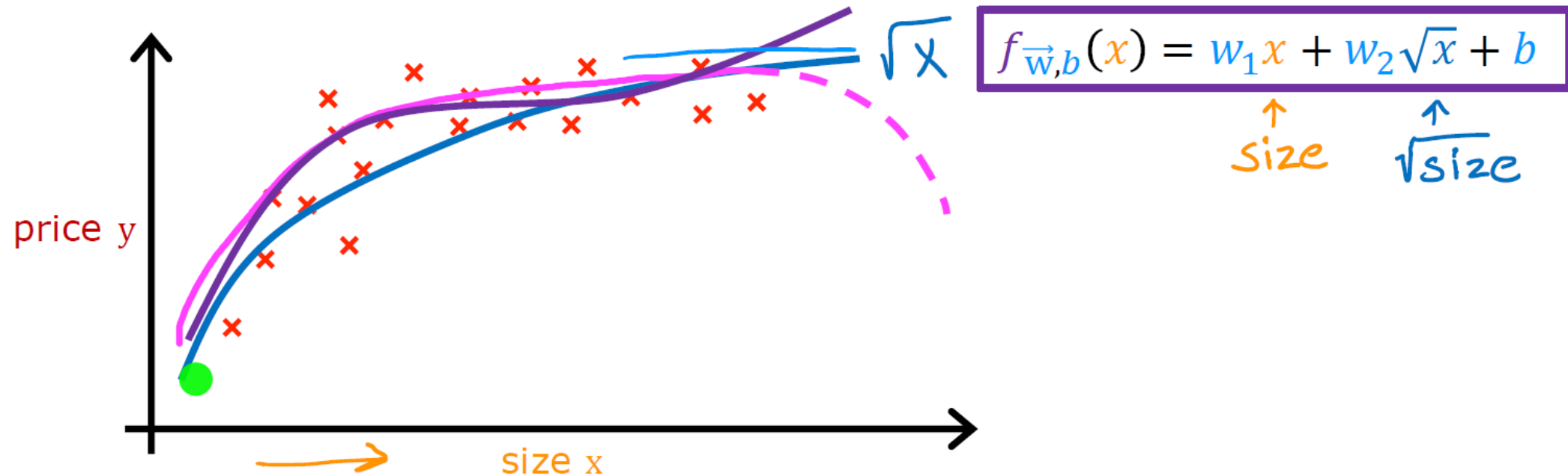# Choosing the Learning Rate

# Choosing the Learning Rate

Values of $\alpha$ to try:

... 0.001  0.003   0.01  0.03   0.1   0.3   1 ...

3X  ≈3X  3X  ≈3X  3X  ≈3X

$J(\vec{w}, b)$

$\alpha = 0.001$

$\alpha = 0.01$

# iterations

$J(\vec{w}, b)$

$\alpha$ too small

$\alpha$ too big

$\alpha$ just right

# iterations

# Polynomial Regression



$$f_{\vec{w},b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + b$$

size     $size^2$     $size^3$

$1-10^3$    $1-10^6$    $1-10^9$

$$f_{\vec{w},b}(x) = w_1 x + w_2 x^2 + b$$

size     $size^2$

price y

size x

# Polynomial Regression



$$f_{\vec{w},b}(x) = w_1 x + w_2 \sqrt{x} + b$$

size    √size

price y

size x

# THANK YOU

NEXT LECTURE WILL BE ONLINE
ON MON, 15.5.2023, IN SHAA ALLAH!

SABBAGH@IEEE.ORG

CONNECT ON LINKEDIN