

Week 02

1] Multiple features

x_j = j^{th} feature

n = number of features

$\vec{x}^{(i)}$ = features of i^{th} training example

$x_j^{(i)}$ = value of feature j in i^{th} training example

Multiple Linear Regression Model

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

$\vec{w} = [w_1, w_2, \dots, w_n]$ parameters
 b = is a number

$\vec{x} = [x_1, x_2, \dots, x_n]$ vector of features

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

2) Vectorization

Parameters and features

*) $\vec{w} = [w_1, w_2, w_3]$

*) b is a number

*) $\vec{x} = [x_1, x_2, x_3]$

Linear algebra count from 1

$$w = np.array([1.0, 2.5, -3.3])$$

$$b = 4$$

$$x = np.array([10, 20, 30])$$

→ Code
Snippet

Code count from 0

Without vectorization

$$1) f(\vec{w}, b)(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$f = \begin{aligned} & w[0] * x[0] + \\ & w[1] * x[1] + \\ & w[2] * x[2] + b \end{aligned} \quad \left. \begin{array}{l} \text{code} \\ \text{snippet} \end{array} \right.$$

$$2) f(\vec{w}, b)(\vec{x}) = \left(\sum_{j=1}^n w_j x_j \right) + b \quad \sum_{j=1}^n \Rightarrow j=1 \dots n$$

range(0, n) / range() $\rightarrow j=0 \dots n-1$

$$f = 0$$

$$\text{for } j \text{ in range}(0, n) \quad \left. \begin{array}{l} \\ f = f + w[j] * x[j] \end{array} \right. \quad \begin{array}{l} \text{code} \\ \text{snippet} \end{array}$$

$$f = f + b$$

With Vectorization

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

$f = np.dot(w, x) + b$ } code

Snippet

Without Vectorization

```
for j in range(0, 16)
    f = f + w[j] * x[j]
```

t₀,
 $f + w[0] * x[0]$

t₁,
 $f + w[1] * x[1]$

t₁₅,
 $f + w[15] * x[15]$

Vectorization

$np.dot(w, x)$:

t₀,
 $w[0] | w[1] | \dots | w[15]$

* * * *
 $x[0] | x[1] | \dots | x[15]$

in parallel

$w[0]*x[0] + w[1]*x[1] + \dots + w[15]*x[15]$

- *) Vectorization is fast, shorter code and efficient specially when scaled to large datasets

3) Gradient descent for Multiple linear regression

	previous notation	Vector notation
Parameters	w_1, \dots, w_n b	$\vec{w} = [w_1, \dots, w_n]$ b
Model	$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + \dots + w_n x_n + b$	$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
Cost function	$J(w_1, \dots, w_n, b)$	$J(\vec{w}, b)$
Gradient descent	repeat { $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_1, \dots, w_n, b)$ $b = b - \alpha \frac{\partial}{\partial b} J(w_1, \dots, w_n, b)$ }	repeat { $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$ $b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$ }

One feature	n features ($n \geq 2$)
<p>repeat {</p> $\omega = \omega - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\omega, b}(x^{(i)}) - y^{(i)}) x^{(i)}$ <p style="text-align: center;">↓</p> $\frac{\partial}{\partial \omega} J(\omega, b)$ <p>$b = b - \alpha \frac{1}{m} \sum_{i=1}^n (f_{\omega, b}(x^{(i)}) - y^{(i)})$</p> <p>Simultaneously update ω, b</p> <p style="text-align: center;">}</p>	<p>repeat {</p> $\omega_j = \omega_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\omega, b}(\bar{x}^{(i)}) - y^{(i)}) x_{j, i}$ <p style="text-align: center;">↓</p> $\frac{\partial}{\partial \omega_j} J(\bar{\omega}, b)$ <p>$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\bar{\omega}, b}(\bar{x}^{(i)}) - y^{(i)})$</p> <p>Simultaneously update ω_j (for $j = 1, \dots, n$) and b</p> <p style="text-align: center;">}</p>

An alternative for gradient descent

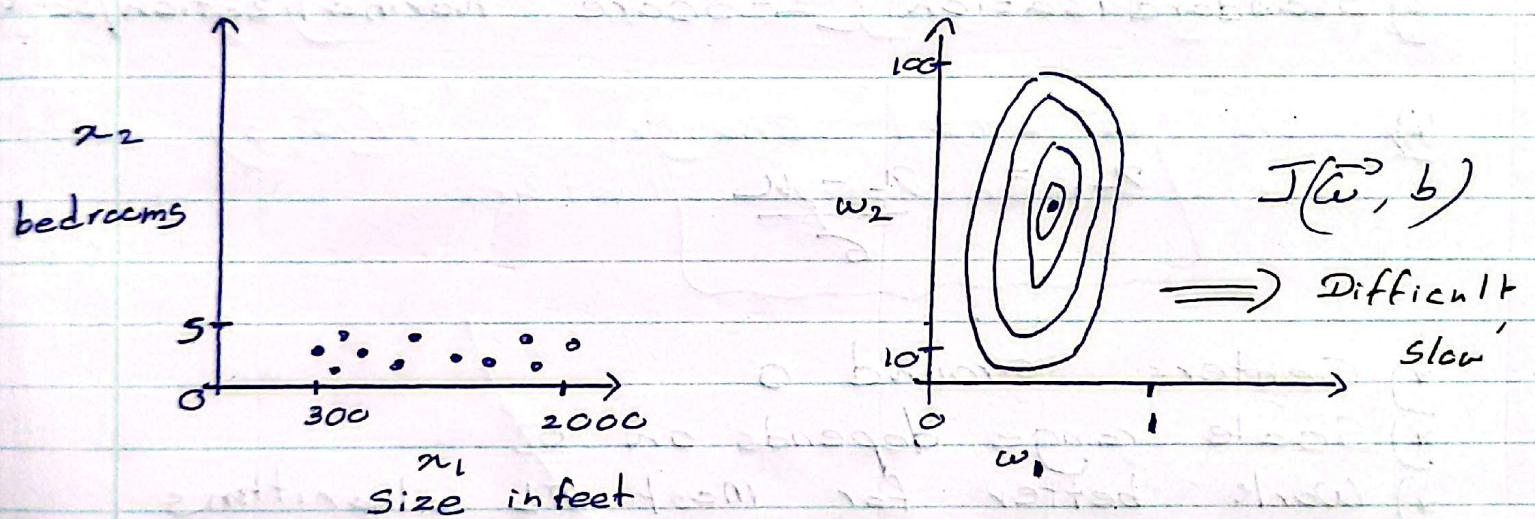
- *) while gradient descent is the recommended method for finding parameters ω, b Normal equation (solve ω, b without iterations) may be used in ML libraries that implement linear regression)

4] Feature Scaling

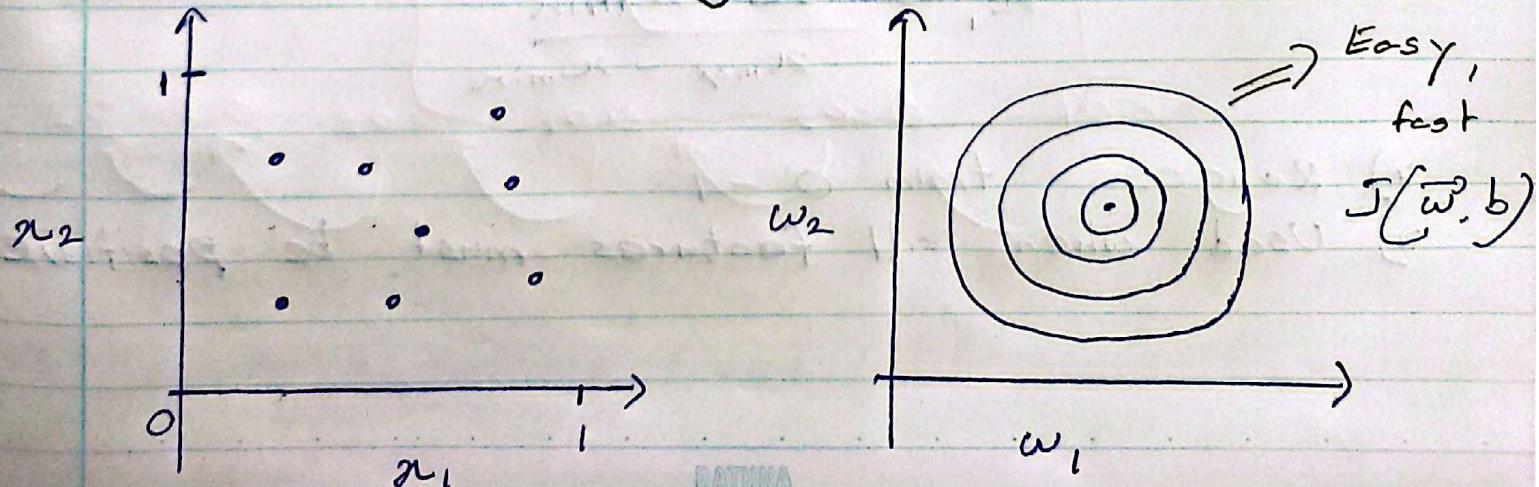
* Feature scaling is a preprocessing step in ML where you adjust the range of input features(variables) so that they are on a similar scale.

* Many ML algorithms work best when all the features have similar range.

Feature size and Gradient Descent



↓
After
Rescaling



Feature Scaling Method

1) Mean Normalization

$$x_i = \frac{x_i - \mu}{x_{\max} - x_{\min}}$$

- *) Centers around 0
- *) Scale range usually $[-1 \Rightarrow 1]$
- *) Sensitive to outliers

2) Standardization (Z-score normalization/Scaling)

$$x_i = \frac{x_i - \mu}{\sigma}$$

- *) Centers around 0
- *) Scale range depends on σ
- *) Work better for most ML algorithms

3) Min max Scaling

$$x_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

- *) Ranges from 0 - 1
- *) Used when all features must be positive

$$0 \leq w_1 \leq 3$$

(No rescale needed)

$$-2 \leq w_2 \leq 0.5$$

(No rescale needed)

$$-100 \leq w_3 \leq 100$$

(Too large, rescale)

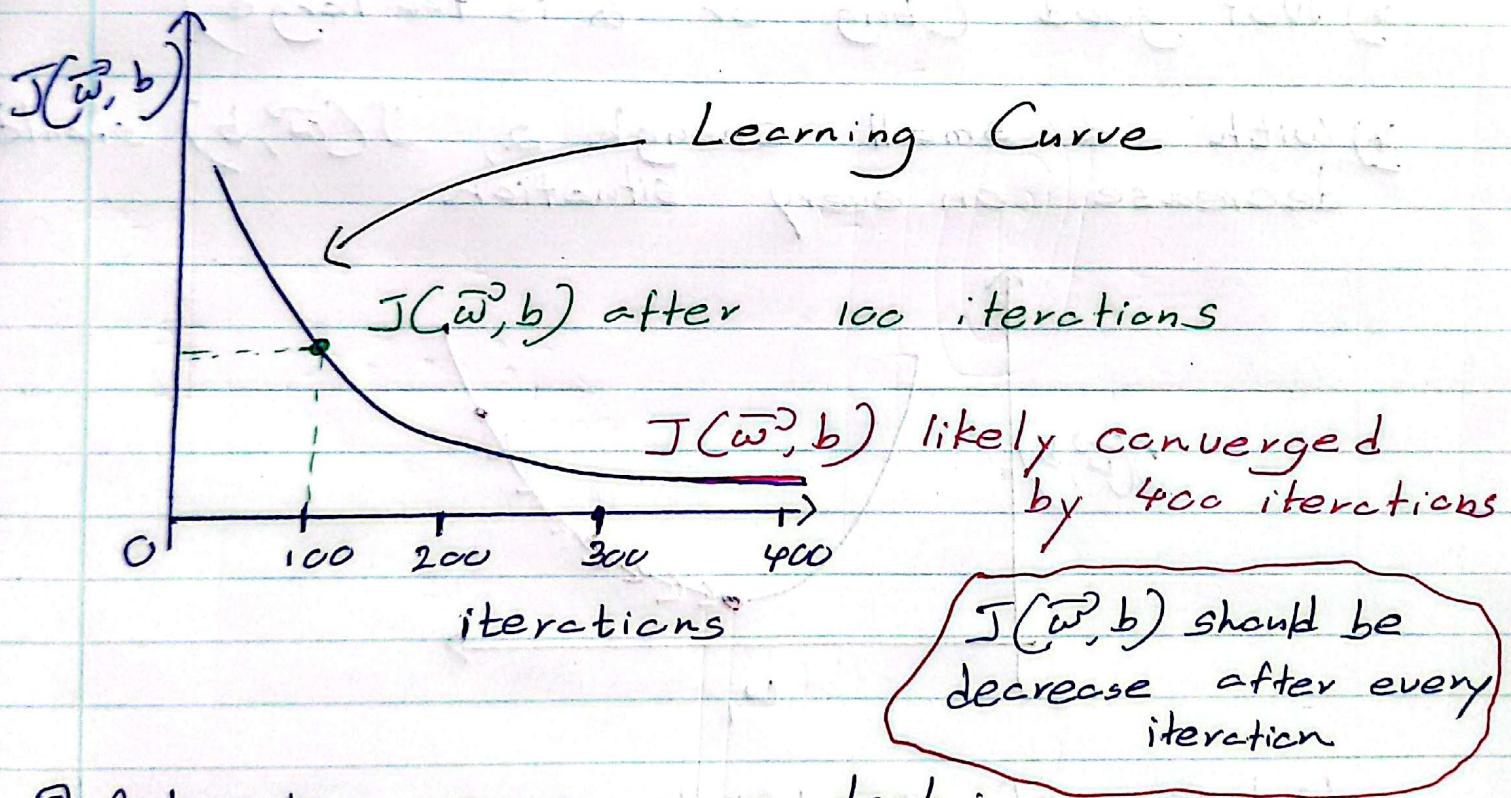
$$-0.001 \leq w_4 \leq 0.001$$

(Too small, rescale)

$$98.6 \leq w_5 \leq 105$$

(Too large → rescale)

5) Checking gradient descent and convergence



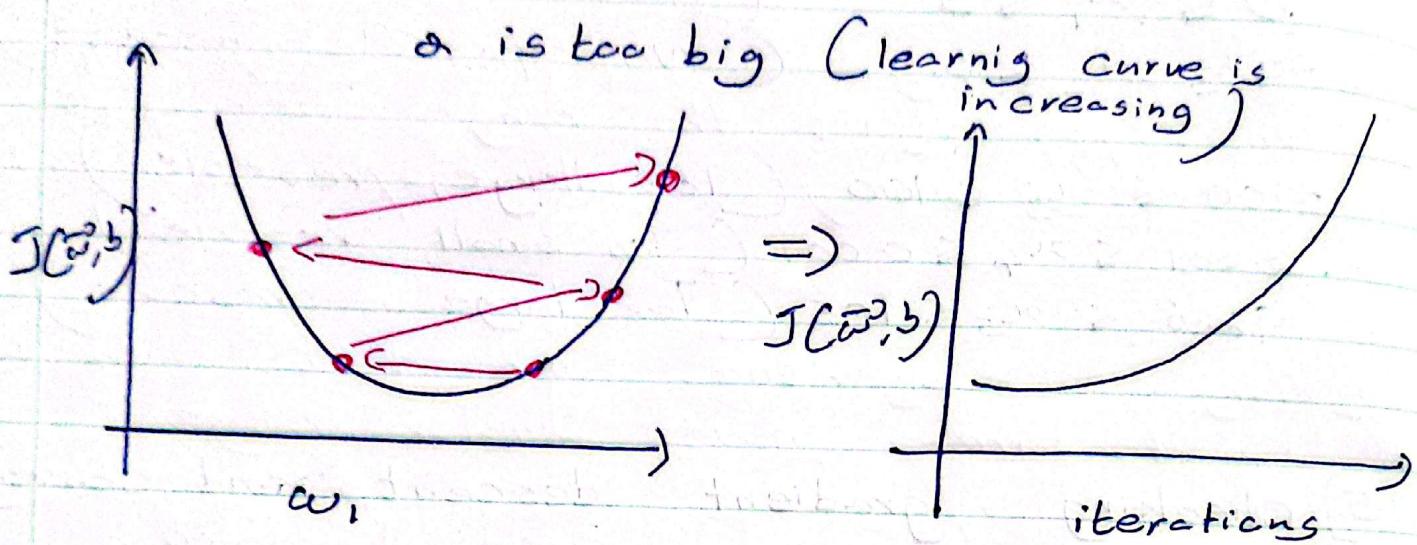
④ Automatic convergence test :

$$\text{Let } \epsilon = 10^{-3}$$

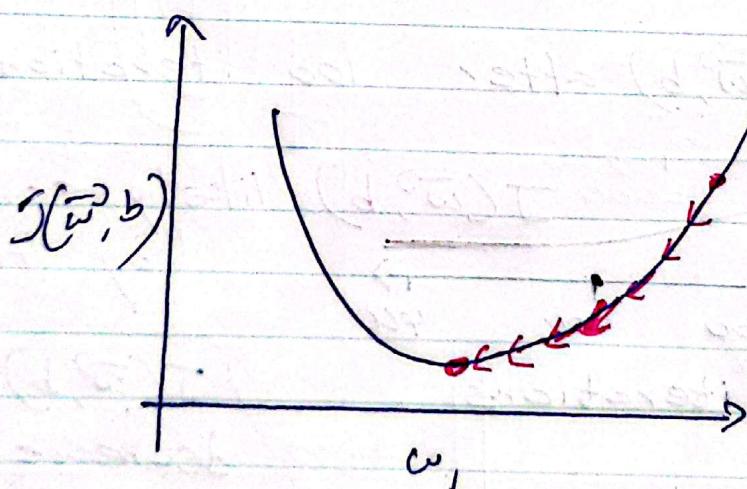
IF $J(\vec{w}, b)$ decreases by $\leq \epsilon$ in one iteration, declare convergence

→ Convergence \Rightarrow (Found parameters \vec{w}, b to get close to global minimum)

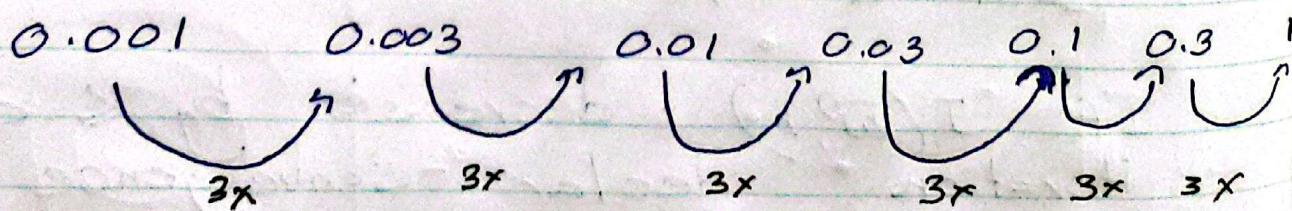
b) Choosing the learning rate

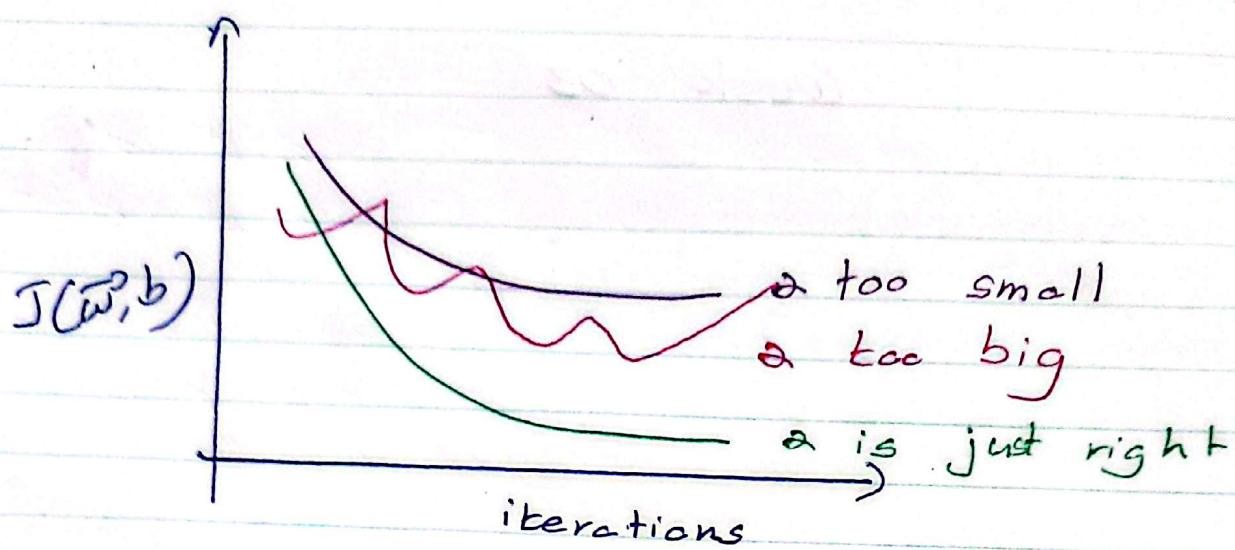


- ✗ Not good (bug or α is too large)
- * with a small enough α , $J(\vec{w}, b)$ should decrease on every situation



Values of α to try:





7) Feature Engineering

- * Using intuition to design new features, by combining or transforming original features
- * It improves model's performance as the algorithm generalize better to new data and learn faster

Turning raw data \rightarrow Useful form

8) Polynomial Regression

- * Polynomial regression is a type of linear regression where the relationship between input x and output y is modeled as an n th degree polynomial (curve) instead of a straight line

$$y = \omega_0 + \omega_1 x + \omega_2 x^2 + \omega_3 x^3 + \dots + \omega_n x^n$$