## Supervised ML Dr Andrew Ng

week 1
Applications of ML

Machine Learning $\longrightarrow$ Field of Study that gives Computer
the ability to learn without being explicitly Programmed

Arthur Samuel (1959)

Machine Learning Algorithms

Supervised Learning ⟶ unsupervised learning

Supervised Learning ↑

used Most in real world
Applications

## Supervised Learning

$X \longrightarrow Y$
input       output Label

$\longrightarrow$ Learns from being given "right Answers"

example

| Input (X) | output (Y) | Applications |
|---|---|---|
| email $\rightarrow$ | spam (o/1) | spam filtering |
| audio $\rightarrow$ | text transcripts | speech recognition |
| English $\rightarrow$ | Spanish | Machine translation |
| ad, userinfo $\rightarrow$ | click? (o/1) | online Advertising |
| image, radar info $\rightarrow$ | position of other cars | self-driving car |
| image of Phone $\rightarrow$ | defect? (o/1) $\rightarrow$ | visual inspection |

Regression: Housing price Prediction

Q) Choose Most appropriate line to fit the data

Regression: Predict a number infinitely many possible outputs

Supervised learning Classification
example Breast Cancer detection

size
2
5
7

diagnosis
c
1
c

⊗ We can use two or more inputs

To Recap

Supervised Learning: Learns from being given "right Answers"

Regression → predict a number infinitely many possible
outputs

Classification → predict Categories small number of possible
outputs

unsupervised learning

Supervised Learning: Learn from data labeled with the "right

Answers"

unsupervised Learning: find something interesting in unlabeled
data

example: Clustering → we use it in Google news

→ DNA microarray

→ Grouping Customers

unsupervised learning: Data only comes within puts X

but not o/p label Y.
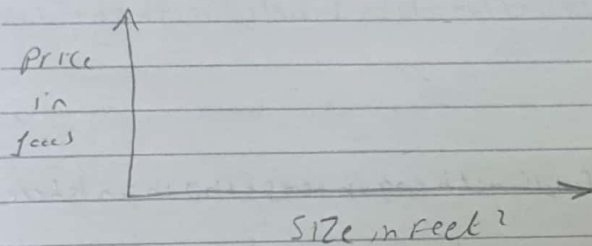
Algorithm has to find structure in the data

1) Clustering: Group similar data points together

2) Anomaly detection: Find unusual data points

3) Diemensionality reduction: Compress data using fewer
   number

## Regression model

### Linear regression



Regression model predicts number
Classification model predicts categories

## Terminology

Training set: Data used to train the model

| Size in feet$^2$ | Price in $ 1000's |
|---|---|
| 2104 | 400 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 572 |
| ... | ... |
| 327 | 87 |

## Notation

$x \to$ "input" variable "feature"

$y \to$ "output" variable
   "target variable"

$m \to$ number of training ex...

$(x, y) \to$ Single trai...

$(x^{(c)}, y^{(c)}) \to c$th t...



$x \to f \to \hat{y}$
Feature     hypothesis/function

Q) How to represe...

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



Linear regr...
equivalent

Notation

x → "input" variable "feature"

y → "output" variable
   "target variable

m → number of training examples

(x, y) → Single training example
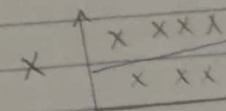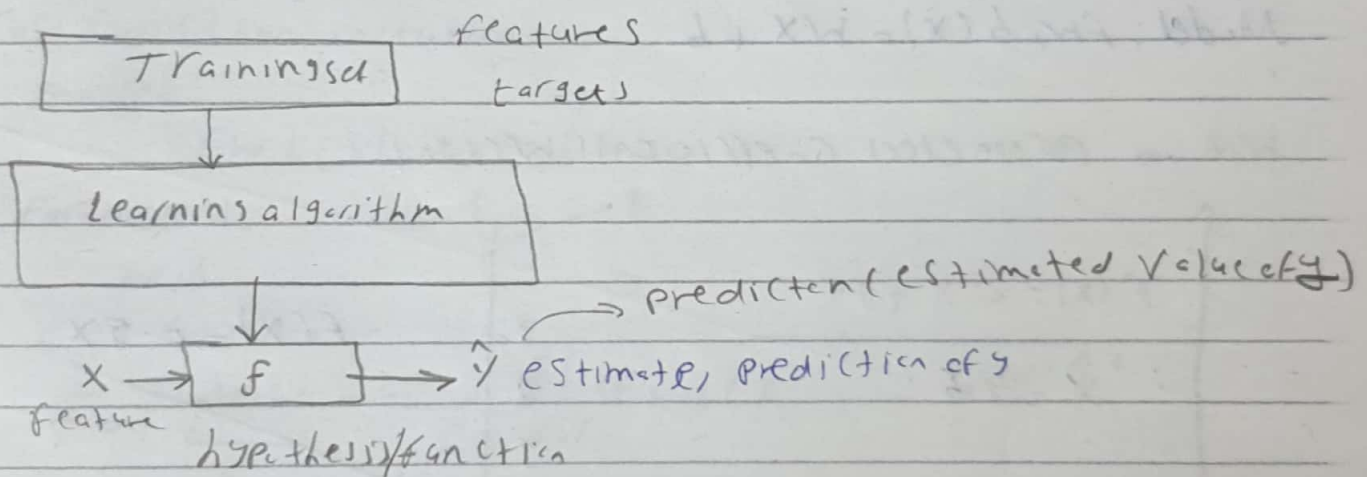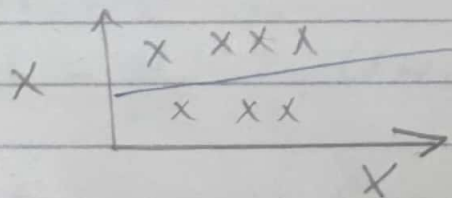
$(x^{(i)}, y^{(i)})$ → i'th training example



Training set          features
                      targets

learning algorithm

x → f → ŷ estimate, prediction of y
feature
    hypothesis/function

→ prediction (estimated value of y)

② | How to represent f

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



→ Single feature x

⌐ Linear regression with one variable
⌐ equivalent to univariate linear regression

## Cost function

Training Set
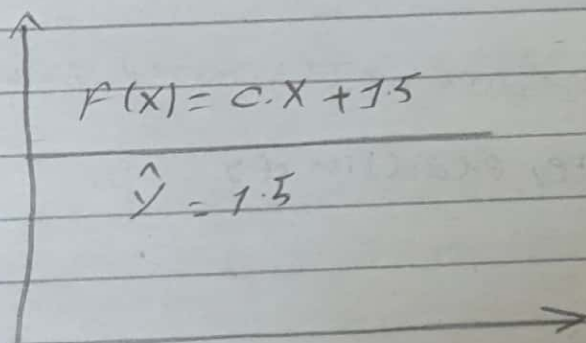
features Size in feet² (X) | targets price $ ccc's (Y)
--- | ---
2104 | 460
1416 | 232
1534 | 315
852 | 178

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

$w,b \rightarrow$ Parameters Coefficients Weight



$f(x) = 0.x + 1.5$

$\hat{y} = 1.5$

$w = 0$
$b = 1.5$
y intercept

$f(x) = 0.5x$

$w = 0.5$
$b = 0$

$f(x) = 0.5x + 1$



$(x^{(i)}, y^{(i)})$

$f_{w,b}$

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

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

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

Cost function : Squared error cost function

$$\frac{1}{2m} \sum_{c=1}^{m} (\hat{y}^{(i)} - y^{(i)})^2 \quad\Rightarrow m \Rightarrow \text{number of training example}$$

error

→ to make the Calculation neater

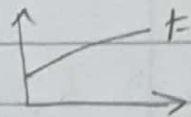$$J(W, b) = \frac{1}{2m} \sum_{c=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Cost function initution

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 :
Minimize $J(W, b)$

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

goal :
minimizes $J(w, b)$

Simplified

$$f_w(x) = wx$$

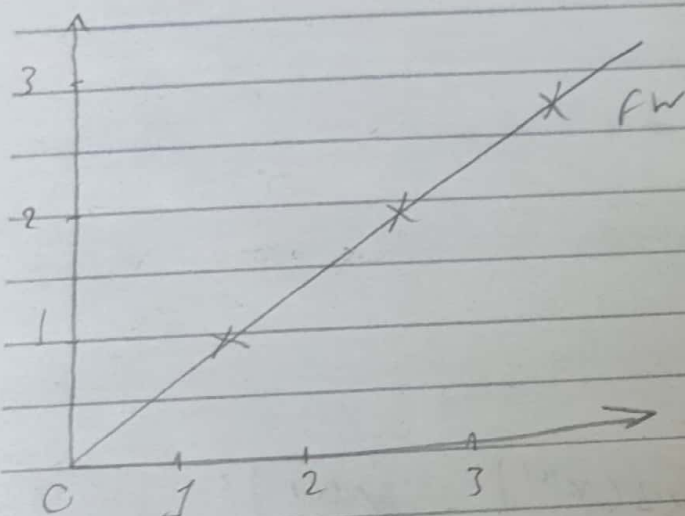$$J(w) = \frac{1}{2m} \sum_{i=1}^{m} (f_w(x^{(i)}) - y^{(i)})^2$$

minimize $J(w)$
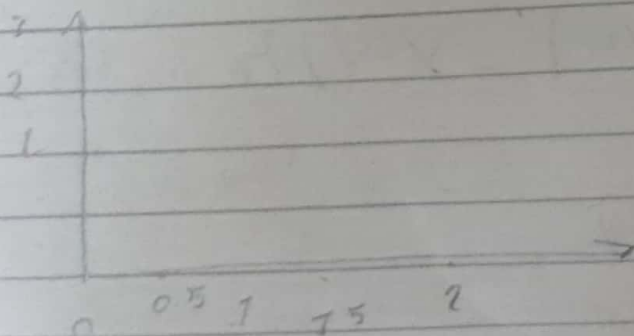
$f_w(x)$
(for fixed $w$, function of $x$)

$J(w)$
(function of $w$)



3

2

1

C    1    2    3

$w=1$

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

$J(w)$ (function of $w$)

3

2

1

0    0.5    1    1.5    2

Simplified $= \frac{1}{2m} ((0.5-1)^2 + (1-2)^2 + (1.5-3)^2) = \frac{3.5}{6}$

$\approx 5$

Goal of linear regression:- minimize $J(w)$

general case :
minimize $J(w, b)$

_____

visualizing cost function

gradient descent $\Rightarrow$ we can use it to train deep
learning model!

Have some function $J(w, b)$ for

want min $J(w, b)$
$\quad w, b$

outline :- start with some $w, b$     set $w = c, b = c$)

Keep Changing $w, b$ to reduce $J(w, b)$
until we settle at or near a minimum

_____

implementing gradient descent

$\qquad$ Assignment $\qquad\qquad$ learning rate

$w = w - \alpha \left[\dfrac{\partial}{\partial w} \left(J(w, b)\right)\right]$

| | Assignment | Truth assertion |
|---|---|---|
| old value | Cost function | |
| | $a = c$ | $a = c$ |
| | | $G = G + 1$ |
| | Code | math |
| $b = b - \alpha \dfrac{\partial}{\partial b} J(w, b)$ | | $a = c$ |

Repeat until convergence

$\qquad\qquad$ Simultaneously
$\qquad\qquad$ update $w$ and $b$

Correct Simultaneous update

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

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

$$w = tmp\_w$$
$$b = tmp\_b$$

Gradient descent Algorithm
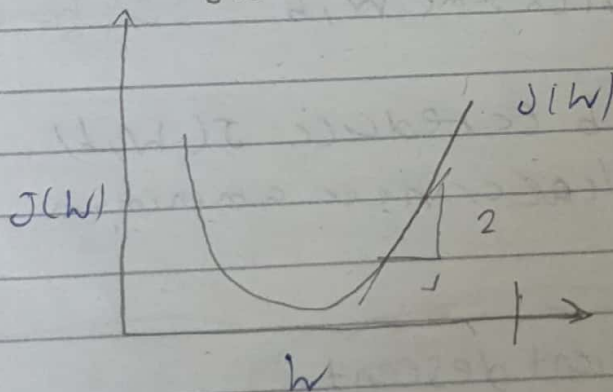
repeat until convergence {

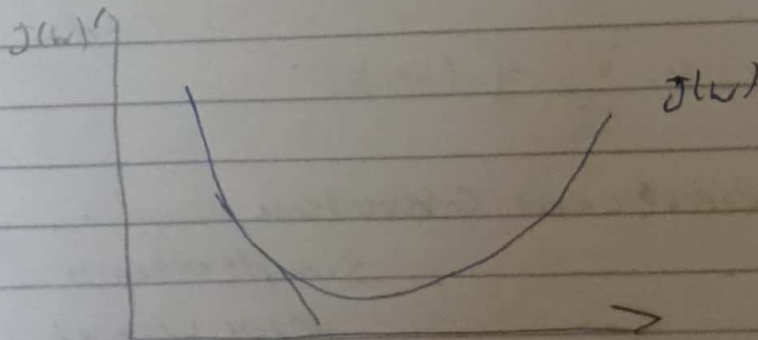$$w = w - \underset{\text{Learning rate}}{\alpha} \underbrace{\frac{\partial}{\partial w} J(w,b)}_{\text{derivative}}$$

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

J(w)

J(w)

2

1

w

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

$$w = w - \alpha \cdot (\text{Positive number})$$

J(w)

J(w)
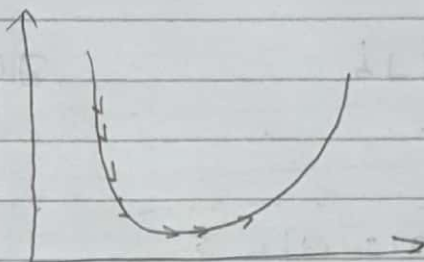
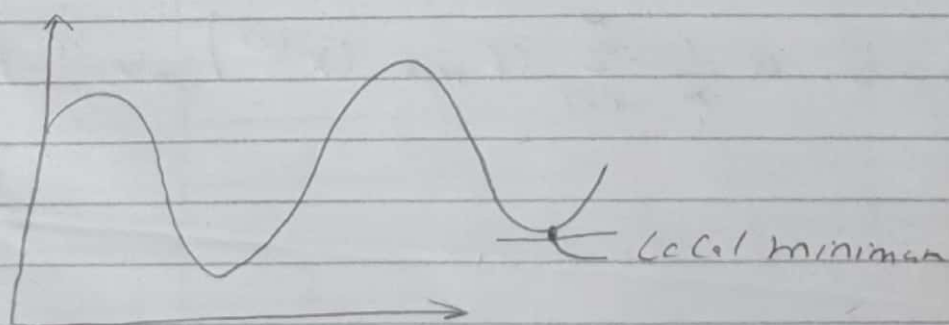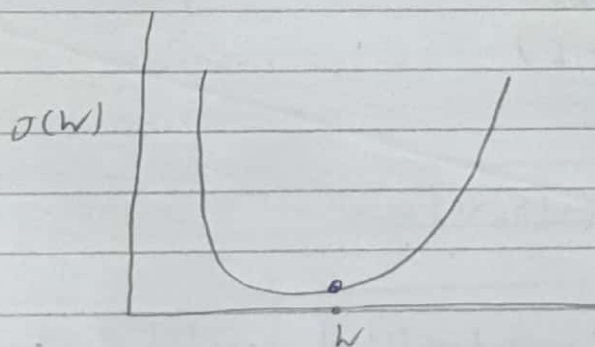$$\frac{d}{dw} J(w)$$

$$w = w - \alpha \cdot (\text{negative number})$$

Learningrate

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

⌐ if $\alpha$ is too Small → very small b.by step
└→ Gradient descent maybe slow



if $\alpha$ is too large

J(w)



w



Local minimu

$$w = w - \alpha \frac{\partial}{\partial w} J(w) \Rightarrow w = w - \alpha \cdot c$$

Can reach Local minimum with fixed learning rate

$$W = W - \alpha \frac{\partial}{\partial W} J(W)$$

Gradient descent for linear regression

$$f_{w,b}(X) = WX + b$$

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

Gradient descent Algorithm {

repeat until Convergence {

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

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

}

after applying derivative

repeat until Convergence {

$$W = W - \alpha \frac{1}{m} \sum_{c=1}^{m} (f_{w,b} (x^{(c)}) - y^{(c)}) x^{(c)}$$

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

Week2

## Multiple features (Variables)

Single feature

| Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 400 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |

## Multiple features

| Size in feet 2 | Number of bedrooms | Number of floors | Age of home in years | Price($) in 1000's |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | |
| | | 2 | 40 | |
| 1416 | 3 | | | |
| | | . | | |
| 1534 | 3 | 2 | 30 | |
| | | | | |
| 852 | 2 | 1 | 36 | |

$x_j = j$ th feature

$n =$ number of features

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

## Model

Previously:- $f_{w,b}(X) = Wx + b$

$f_{w,b}(X) = 0.1 X_1 + 4 X_2 + 10 X_3$

$f_{w,b}(X) = W_1 X_1 + W_2 X_2 + \cdots + W_n X_n + b$

$f_{\vec{w},b}(\vec{X}) = \vec{W} \cdot \vec{X} + b = W_1 X_1 + W_2 X_2 + \cdots + W_n X_n + b$

---

### Parameters and features

$\vec{W} = [W_1, W_2, W_3]$    $n = 3$

$b$ is a number

$\vec{X} = [X_1 \ X_2 \ X_3]$

linear algebra count from 1

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

$b = 4$

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

$f = np.dot(W, X) + b$

without victorization                                    Vectorization

for J in range (o, 16);                                $np.dot(w/x)$
   $F = F + W[J] * X[J]$

Gradient descent $\vec{W} = (W_1, W_2, \ldots, W_{16})$ b Parameters

$\vec{d} = (d_1, d_2, \ldots d_{16})$

$W = np.array([o.5, 1.3, \ldots, 3.4])$

$d = np.array([o.3, c.2, \ldots, c.4])$

$W_J = W_J - c.1 d_J$ for $J=1 \ldots 16$

with vectorization

$\vec{W} = \vec{W} - c.1 \vec{d}$

$W = W - c.1 \times d$

$np.random.random\_sample(4)$                    random vector اس١ لیں

$G = np.arange(4)$

## Gradient descent

repeat $\{$

$$WJ = WJ - \alpha \frac{\partial}{\partial WJ} J(W_1, \dots, W_n, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(W_1, \dots, W_n, b)$$

$$WJ = WJ - \alpha \boxed{\frac{1}{m} \sum_{c=1}^{m} (f_{\vec{W}, b}(\vec{X}^{(c)}) - y^{(c)}) X_J^{(c)}}$$

$$\frac{\partial}{\partial W_j} J(\vec{W} + b)$$

## An Alternative to gradient descent

### Normal equation
$\rightarrow$ only for linear regression
$\rightarrow$ Solve for W, b without
$\rightarrow$ iterations

### Disadvantages
$\rightarrow$ Doesn't generalize to other
$\rightarrow$ learning Algorithms.

Gradient descent Practice

$$Price = W_1 X_1 + W_2 X_2 + b$$

$\downarrow$ size

$\downarrow$ #bedrooms

X1: Size (feet²)
range: 300 – 2,000

X2: # bedrooms
rang: 0:5

House: $X_1 = 2000$, $X_2 = 5$, Price = 500 k   one training
Size of the Parameters      example
$W_1, W_2$

$W_1 = 50$, $W_2 = 0.1$, $b = 50$

$Price = 50 × 2000 + 0.1×5 + 50$

100,000 K      0.5k  5.k

$Price = \$ 100, 050.5 K = \$ 100, 050, 500$

$W_1 = 0.1$         $W_2 = 50$      $b = 50$
Small            large

$Price = 0.1 × 2000 K + 50×5+50$          $= 500 \$ k$

200 K          250 k  50 k                 Reasonable

Features N scaling     تعمل المعادلات بدون يوثر قيمة اكبر

global
اقرب كتلة Minimum للوصول       عشان تبحث عن اقرب
bouncing                        بحث عن

elipse       بدل circle   مش  Contourplt      ويصبح

# How do We Actually scale features

$$3cc \leq X1 \leq 2ccc \qquad\qquad c \leq X_2 \leq 5$$

$$X_1, scaled = \frac{X1}{2000} \qquad\qquad X_2, scaled = \frac{X_2}{5}$$

$$0.15 \leq X_1, scaled \leq 1 \qquad\qquad 0 \leq X_2, scaled \leq 1$$

قسمة على Maximum ae divide   قسمة على Mean normlize the.

Mean normalize the.

$$3cc \leq X1 \leq 2ccc \qquad\qquad c \leq X_2 \leq 5$$



$X_2$
#bedrooms
5
c
3cc   2ccc   → 2ccc

$X_2$
#bedrooms
normalized
→ size in feet2
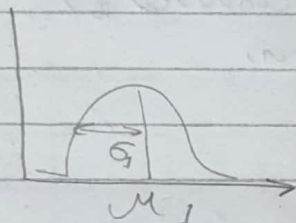
using Mean normalization

$$3cc \leq X_1 \leq 2000$$

Let $M_1 = 6cc$

$$X_1 = \frac{3cc - 6cc}{2000 - 3cc} \leq X_1 \leq \frac{2000 - 6cc}{2000 - 3cc}$$

$$-0.18 \leq X_1 \leq 0.82$$

Standard      مراد بها Z Score normalization (معدل) حارة ح
deviation



$$X_1 = \frac{X_1 - M_1}{\sigma_1} \qquad \text{Let } \sigma = 45c, \; M_1 = 6c$$

$$-0.67$$

$$\boxed{\frac{3cc - 6cc}{45c}} \leq X \leq 3.1$$

قبلها الأشكال تطبيق على ال second
feature

→ Note for feature scalling

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

$$-3 \leq X_j \leq 3 \qquad \left. \right\} \text{ Acceptable ranges}$$
$$-0.3 \leq X_j \leq 0.3$$

$$0 \leq X_j \leq 3 \qquad \text{okay, no rescaling}$$
$$-2 \leq X_2 \leq 0.5 \qquad \text{okay, no rescaling}$$
$$-100 \leq X_3 \leq 100 \qquad \text{too large} \rightarrow \text{rescale}$$
$$-0.001 \leq X_4 \leq 0.001 \rightarrow \text{too small} \rightarrow \text{rescale}$$

Checking gradient descent for convergence

objective $\min\limits_{\vec{w}, b} J(w, b)$



Learning curve

$J(\vec{w}, b)$ after 1 iterations

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

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

# iterations

Automatic convergence test let $\varepsilon$ epsilon be $10^{-3}$

0.001

if $J(\vec{w}, b)$ decrease $b \leq \varepsilon$ in one iteration declare

Convergence

Choosing good learning rate



# iterations

كول دايا دايق ف كا
large  او  busen code
learning
rate

Notice $w_1 = w_1 + \alpha d_1$ :)

descent $w_1 - \alpha d_1$ :)

derivative  minus :)
term

Feature engineering

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

frontage   depth

area = frontage × depth

$x_3 = x_1 x_2 \longrightarrow$ new feature

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

Polynomial regression



price
Y

Size

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

Size      Size²       Size³

1-1J      1-1c⁶      1-1⁹

We have choice of each feature we want to build most
suitable model