

Machine Learning Specialization

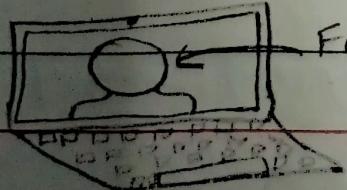
1). Supervised Machine Learning

Week -1

Overview of ML

Machine Learning is a field of study in AI concerned with the development & study of statistical algorithms that can use to teach machine to perform a task without human instruction.

for example, anytime you want to find out something like "how do i make a sushi roll?" You can do a web search on google to find out. And that work so well because their Machine Learning software has figured out how to rank web pages.



Face
recognition

Applications of Machine Learning

- Image Recognition
- Speech Recognition
- Automatic Language Translation
- Self driving car
- Recommendation systems, etc.

Supervised Vs Unsupervised ML

Machine Learning : "Field of study that gives computers the ability to learn without being explicitly programmed".

Machine Learning Algorithm

↓
Supervised Learning

↓
Unsupervised Learning

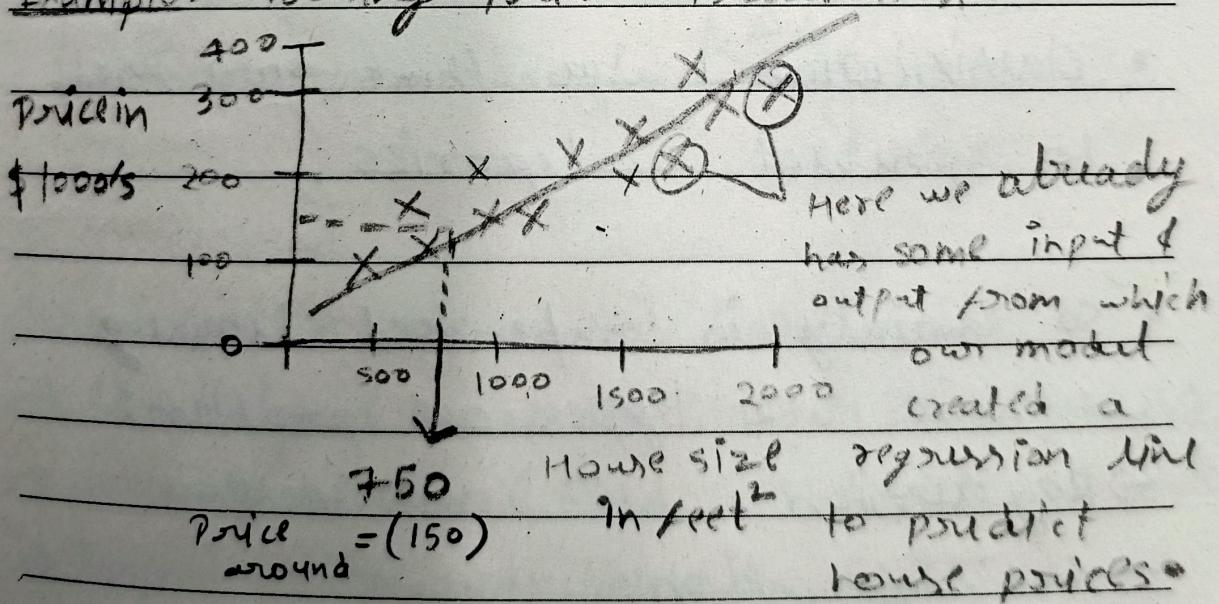
Supervised Learning

$X \rightarrow Y$
(Input) (Output) Label

Learn from being given "right answers"

- supervised Machine Learning or more commonly, supervised learning, refers to algorithm that learn X to Y or input to output mappings.
- the key characteristic of supervised learning is that you give your learning algorithm examples to learn from. That includes the right answer, whereby right answers, means that a correct label Y for a given input X .

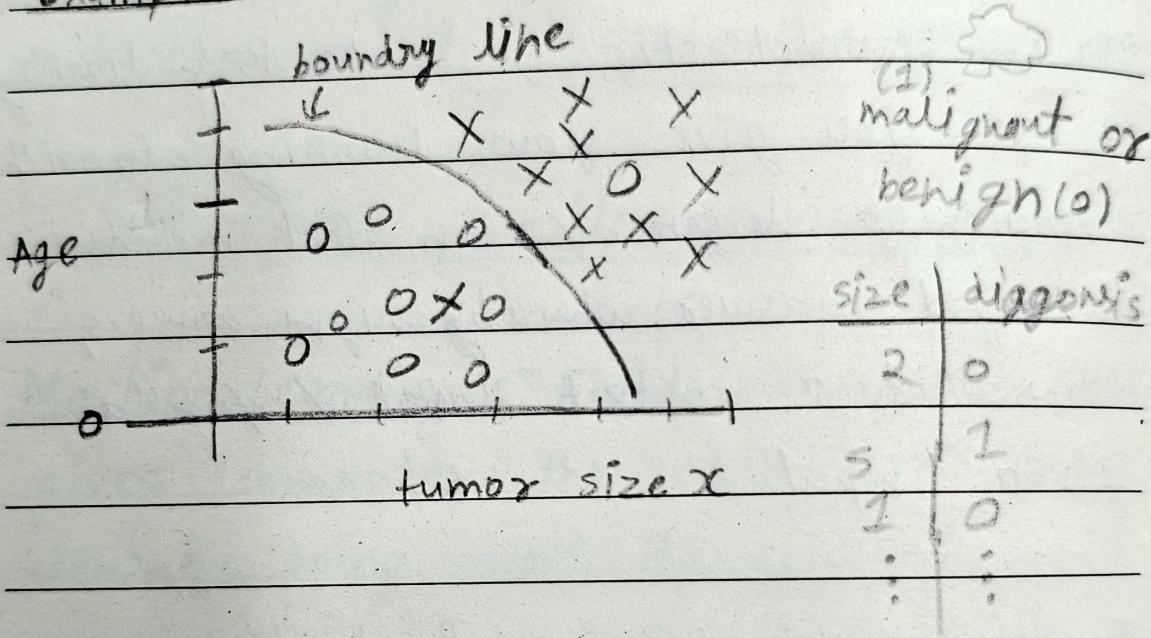
Example: Housing Price Prediction



Regression

- Another example of a supervised learning can be a "classification" problem.

Example: Breast cancer detection



- Classification algorithms are used to predict categories.
- So mainly in supervised learning we have 2 types of problems:
→ Regression
→ Classification.

Unsupervised Learning

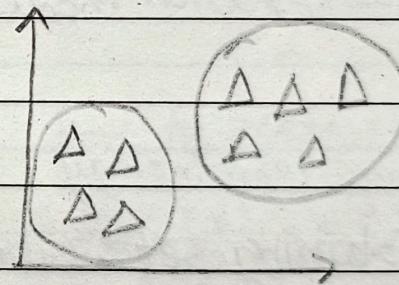
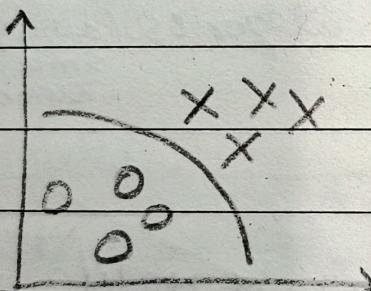
- Unlike supervised learning, unsupervised learning model are given unlabeled data & allowed to discover patterns & insights without any explicit guidance or instruction.

Supervised Learning

Learn from data labeled
with the "right answers"

Unsupervised Learning

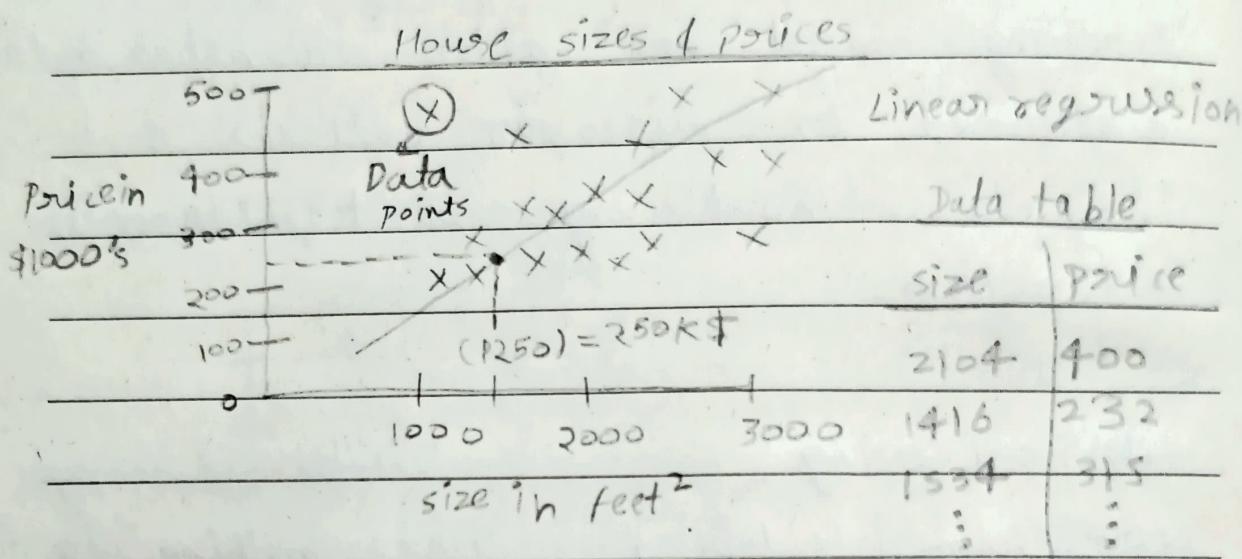
Find something interesting
in unlabeled data.



clustering

- Some type of unsupervised learning are
 - Clustering (group similar data points together).
 - Anomaly detection (Find unusual data points).
 - Dimensionality reduction (compress data with fewer numbers).

Regression Model



- Regression model predicts numbers. (Infinite)
- Classification model predicts categories. (discrete)
small numbers

Some terminology

- Training set : Data used to train the model

x size in foot ²	y price in \$1000's
(1) 2104	400
(2) 1416	232
(3) 1534	315
:	:
(47)	

- Notation : x = "input" variable feature

y = "output" variable

"target" variable

$$\Rightarrow x = 2104, y = 400$$

m = number of training examples

(x, y) = single training example

$$\Rightarrow (x, y) = (2104, 400)$$

$(x^{(i)}, y^{(i)})$ = i^{th} training example
 $(1^{\text{st}}, 2^{\text{nd}}, 3^{\text{rd}}, \dots)$

$$\Rightarrow x^{(1)} = 2104, y^{(1)} = 400$$

$$(x^{(1)}, y^{(1)}) = (2104, 400)$$

$$(x^{(2)}, y^{(2)}) = (1416, 232)$$

$$\Rightarrow x^{(2)} \neq x^2 \text{ not exponent}$$

f = function

\hat{y} = "y-hat" predicted value

Model

How to represent f ?

training set features
target

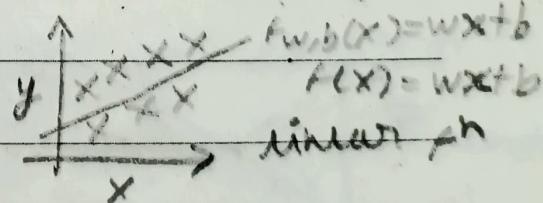
$$f_{w,b}(x) = wx + b$$

$$f(x)$$

learning algorithm

for now let w, b are
some numbers which help
us to find the \hat{y} .

$x \rightarrow [F] \rightarrow \hat{y}$



feature prediction
(estimated y)

target • Linear regression with
one variable.

size $\rightarrow [E] \rightarrow \text{price}$

(size 2)

(estimated)

• Univariate linear
regression

Cost Function

Training set:

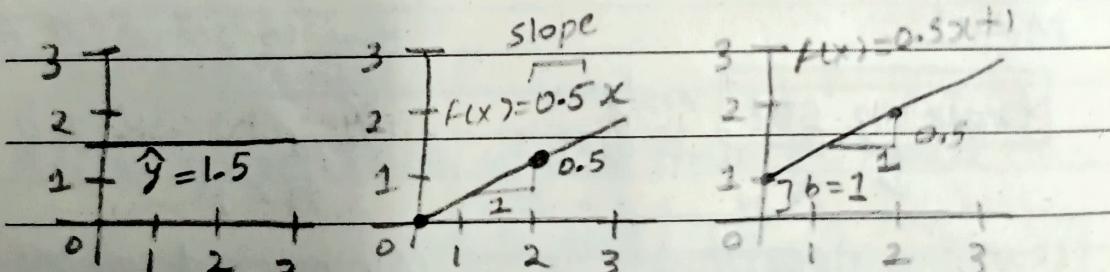
size in feet ² (x)	price \$1000's (y)
2104	460
1416	232
1534	315
852	178
:	:

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

w, b = parameters what do "w, b" do?

[coefficients]
[weights]

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



$$\Rightarrow w = 0$$

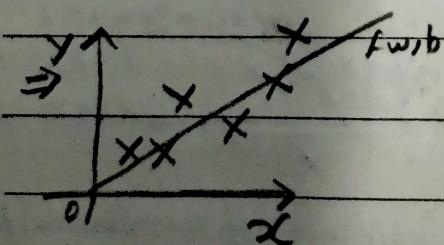
$$b = 1.5$$

$$w = 0.5$$

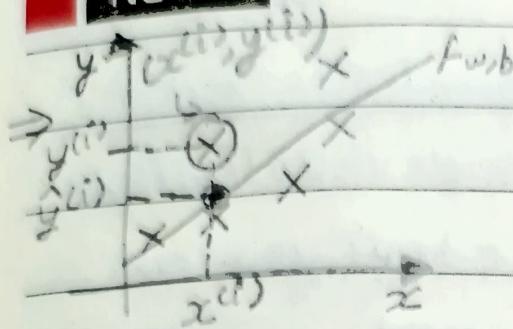
$$b = 0$$

$$w = 0.5$$

$$b = 1$$



• so w & b are some parameters which helps in finding a $f_{w,b}$ line by changing these parameters again & again



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

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

★ So, how should we know the correct value of w & b for the linear regression model?

⇒ So, the ans is "Cost function".

• Cost function helps in finding the difference b/w the actual & predicted value of the minimum the cost fn implise more accurate model.

⇒ Cost function : Squared Error Cost Function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2, m = \text{no. of train example}$$

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

Train the Model with Gradient Descent

* In last part we saw visualizations of the cost function J & how you can try different choices of the parameters " w " & " b ".

Gradient Descent

- It would be nice if we had a more systematic way to find the values of " w " & " b ", that results in the smallest possible cost, J of w, b . It turns out that there's an algorithm called "gradient descent" that you can use to do that.

* Gradient Descent is a machine learning algorithm which helps in minimizing the cost function as much as possible to find the best value of ' w ' & ' b '.

Outline of gradient descent

- 1). start with some w, b . (set $w=0, b=0$)
- 2). keep changing w, b to reduce $J(w, b)$ [by e^{2^n}]
- 3). Until we settle at or near a minimum (~ 0)

Gradient descent algorithm.→ Learning rate (α)

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

↓

Assignment operator

→ Derivative term of
cost fn $J(w, b)$

• Repeat
until

convergence
• simultaneously
update $w \& b$

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

* Correct way to implement: simultaneously update

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

we do this because in updating value

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

of 'b' we have to use old value of 'w'.

$$w = \text{temp_}w$$

$$b = \text{temp_}b$$

updating the
value of 'w' & 'b'

• More intuition about gradient descent?

repeat until convergence {

$$w = w - \alpha \left[\frac{d J(w, b)}{dw} \right] \text{"derivative" term}$$

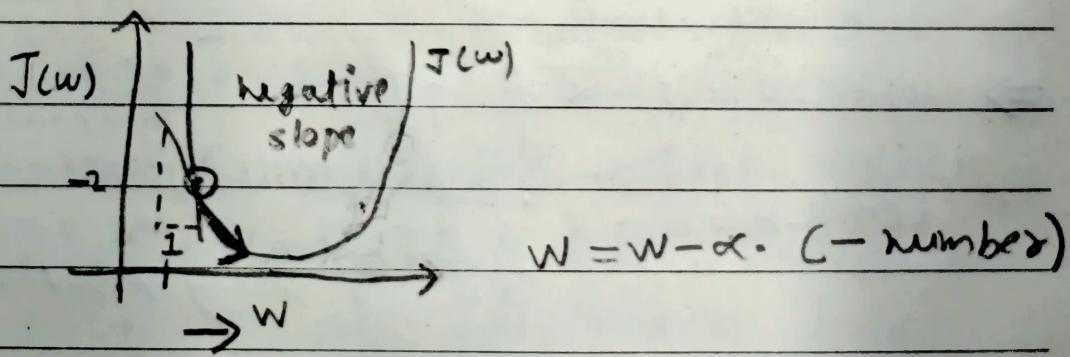
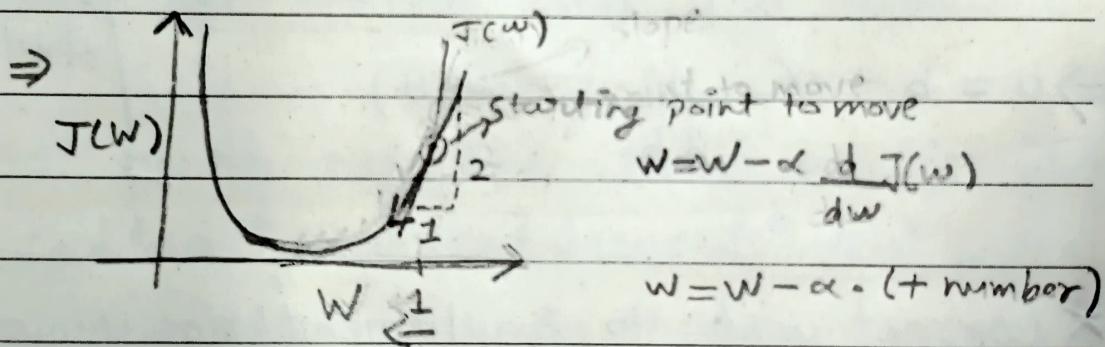
marking index

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

} simultaneously update

Suppose we have cost J^h or
 $J(w)$,

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



Learning Rate (α)

- Alpha also known as the learning rate, is a hyperparameter that determines the step size at each iteration of the gradient descent algorithm.
- It controls how quickly or slowly the algorithm converges to the optimal solution. The value of alpha is crucial as it can greatly impact the optimization process.

→ repeat until convergence:

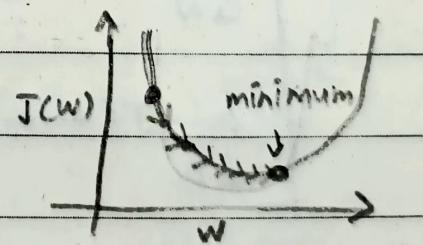
$$w = w - \alpha \frac{d}{dw} J(w)$$

alpha ← $\frac{d}{dw}$

} simultaneously update

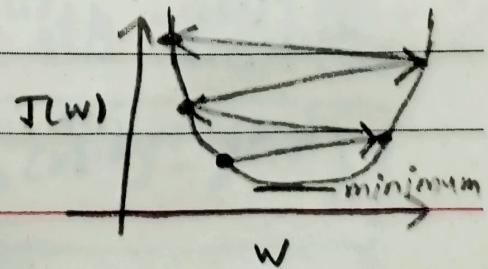
- if α is too small...

⇒ gradient descent may be slow



- if α is too large...

⇒ gradient descent may never reach minimum.



Gradient Descent for Linear Regression

Let's Recap:

- Linear regression model

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

- Cost Function

$$J(w, b) = \frac{1}{2m} \sum_{i=0}^{m-1} (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

- Gradient descent algorithm

repeat until convergence {

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

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

} simultaneously update

⇒ so, now derivatives of :

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

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

so, how we get these terms?

Let's find out : (optional)

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

Here we can see that we have to do partial derivatives with respect to 'w'.

$$\begin{aligned}\Rightarrow \frac{d}{dw} J(w, b) &= \frac{d}{dw} \left[\frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right] \\ &= \frac{d}{dw} \left[\frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2 \right] \\ &= \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)}) \times 2x^{(i)}\end{aligned}$$

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

similarly, with respect to 'b',

$$\begin{aligned}\Rightarrow \frac{d}{db} J(w, b) &= \frac{d}{db} \left[\frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right] \\ &= \frac{d}{db} \left[\frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2 \right] \\ &= \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)}) 2\end{aligned}$$

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

• So, now final gradient descent algorithm for linear Regression is:

⇒ repeat until convergence {

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

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

} simultaneously update

Week 2

Multiple Linear Regression

Recap: In previous week we have working with single variable feature:

size in feet ² (x)	price (\$) in 1000s (y)
2104	400
1416	232
1534	315
852	178
...	...

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

- But if we have multiple features like:

size in feet ²	Number of bedrooms	Number of floors	Age of home in years	Price
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

$x_1 \quad x_2 \quad x_3 \quad x_4$

$$\Rightarrow x_j = j^{\text{th}} \text{ feature} \quad // \quad x_1, x_2, x_3, x_4$$

$$n = \text{number of features} \quad // \quad n=4$$

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

$$\text{if } \vec{x}^{(2)} = [1534, 3, 2, 40]$$

$$x_j^{(i)} = \text{value of feature } j \text{ in } i^{\text{th}} \text{ training example} \quad // \quad x_3^{(2)} = 2$$

Model:

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

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

Example:

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

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

#size #bedrooms #floors #years base
price

- So, from above we can say that,
if we have n features, then model
will look like:

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

$$\Rightarrow \vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n] \text{ vectors}$$

$$\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$$

$b = b$ is a number

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

\uparrow
dot product

* A Linear Regression with multiple Feature is known as "multiple linear regression".

Parameters & features:

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

b is a number

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

Numpy

$$\Rightarrow \left[\begin{array}{l} \vec{w} = \text{np.array}([1.0, 2.5, -3.3]) \\ b = 4 \end{array} \right]$$

$$\vec{x} = \text{np.array}([10, 20, 30])$$

① without vectorization

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

$$\left[\begin{array}{l} \Rightarrow f = w[0] * x[0] + w[1] * x[1] + \\ w[2] * x[2] + b \end{array} \right]$$



② without vectorization (with loop)

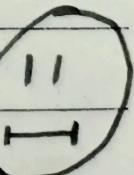
$$f_{\vec{w}, b}(\vec{x}) = \sum_{j=1}^n w_j x_j + b$$

$$\Rightarrow f = 0$$

for j in range(0, n):

$$f = f + w[j] * x[j]$$

$$f = f + b$$



With Vectorization

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

$$\Rightarrow f = np.dot(w, x) + b$$



Let's analyse code:

without vectorization

$$\begin{bmatrix} \Rightarrow \text{for } j \text{ in range}(0, 16): \\ f = f + w[j] * x[j] \end{bmatrix}$$

$$t_0 f + w[0] * x[0]$$

$$t_1 f + w[1] * x[1]$$

...

$$t_{15} f + w[15] * x[15]$$

vectorization

$$\begin{bmatrix} \Rightarrow np.dot(w, x) \end{bmatrix}$$

$$\begin{bmatrix} w[0] & w[1] & \dots & w[15] \end{bmatrix}$$

$$\begin{bmatrix} * & * & \dots & * \end{bmatrix}$$

$$\begin{bmatrix} x[0] & x[1] & \dots & x[15] \end{bmatrix}$$

t_0

$$\begin{bmatrix} w[0]*x[0] + w[1]*x[1] + \dots + w[15]*x[15] \end{bmatrix}$$

$$\dots + \boxed{w[15]*x[15]}$$

Gradient Descent for Multiple Regression

vector notation

Parameters:

$$\vec{w} = [w_1 \dots w_n]$$

 b is a number

Model :

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

Cost Function: $J(\vec{w}, b)$

Gradient Descent: repeat {

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

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

}; simultaneously update

Gradient Descent

one feature

repeat {

$$w = w - \alpha \frac{1}{m} \sum_{i=0}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

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

}; simultaneously update

multiple features

repeat {

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

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

} simultaneously update

Gradient Descent in Practice

Feature scaling

• Feature scaling is the process of normalizing the range of feature in a dataset. Real-world dataset are more often various, in range, degree, etc. Therefore, in order for machine learning models to interpret these features on the same scale, we need to perform feature scaling.

Let,

$$\widehat{\text{price}} = w_1 x_1 + w_2 x_2 + b$$

$\downarrow \quad \quad \quad \downarrow$
 size #bedrooms

x_1 : size(feet²)
 x_2 : # bedrooms

$\Rightarrow x_1$: [range = 300 - 2,000] large

x_2 : [range = 0 - 5] small

House: $x_1 = 2000$, $x_2 = 5$, $\text{price} = \$500k$

\Rightarrow size of the parameters w_1, w_2 ?

\Rightarrow suppose,

$$w_1 = 50, w_2 = 0.1, b = 50$$

$$\widehat{\text{price}} = 50 * 2000 + 0.1 * 5 + 50$$

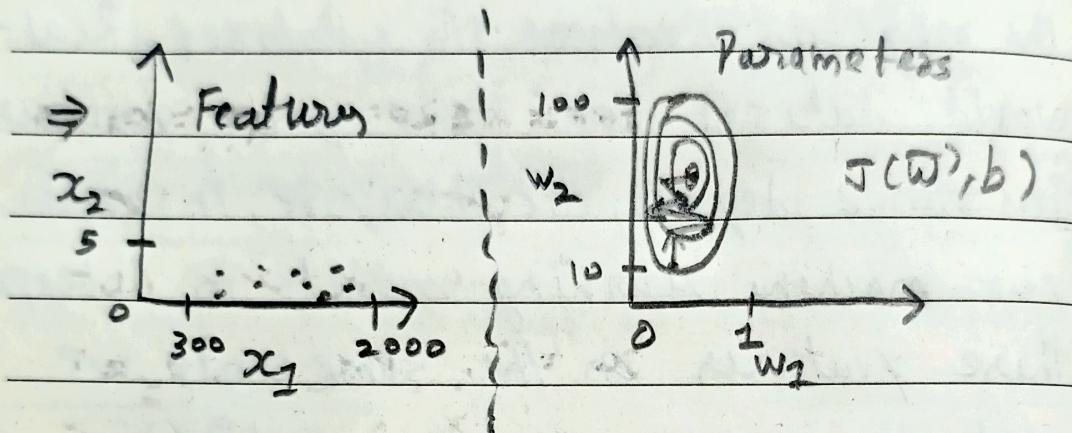
$$= \$100,050.5k > \$500k$$

\Rightarrow sub func,

$$w_1 = 0.1, w_2 = 50, b = 50$$

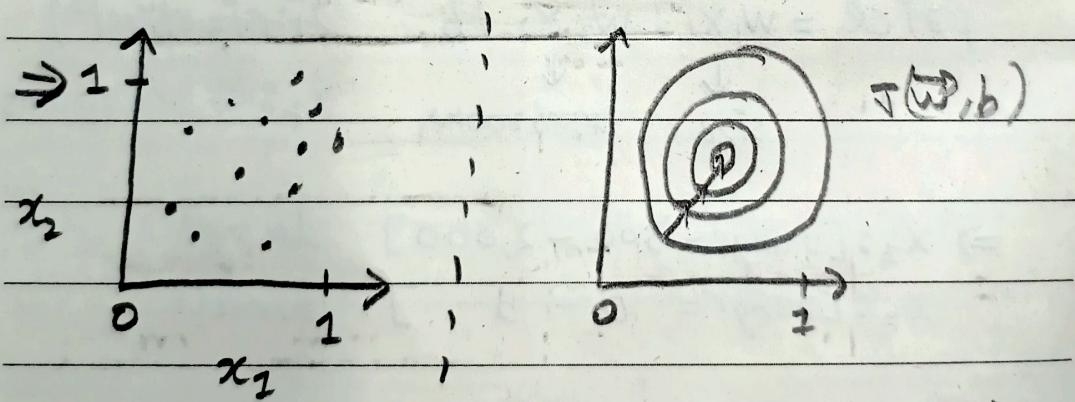
$$\widehat{\text{price}} = 0.1 \times 2000 + 50 \times 5 + 50$$

$$= \$500K \text{ more reasonable}$$



After feature scaling

[doing some transformation]

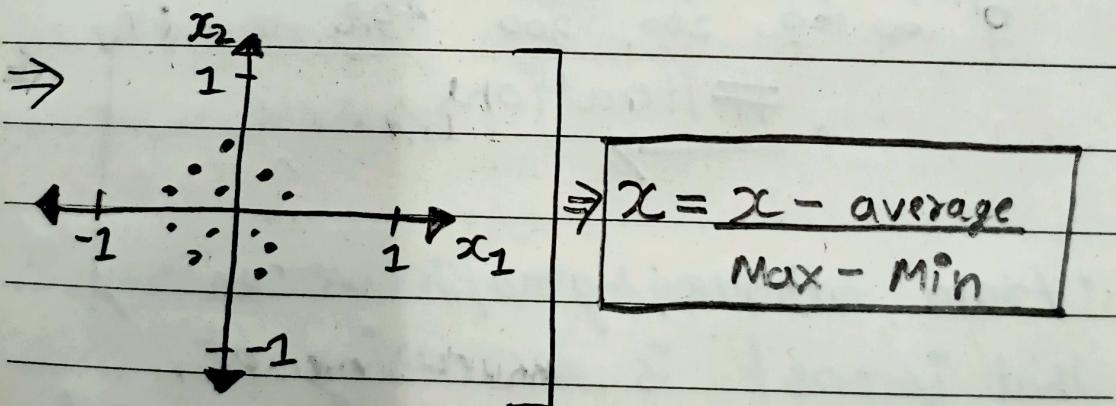
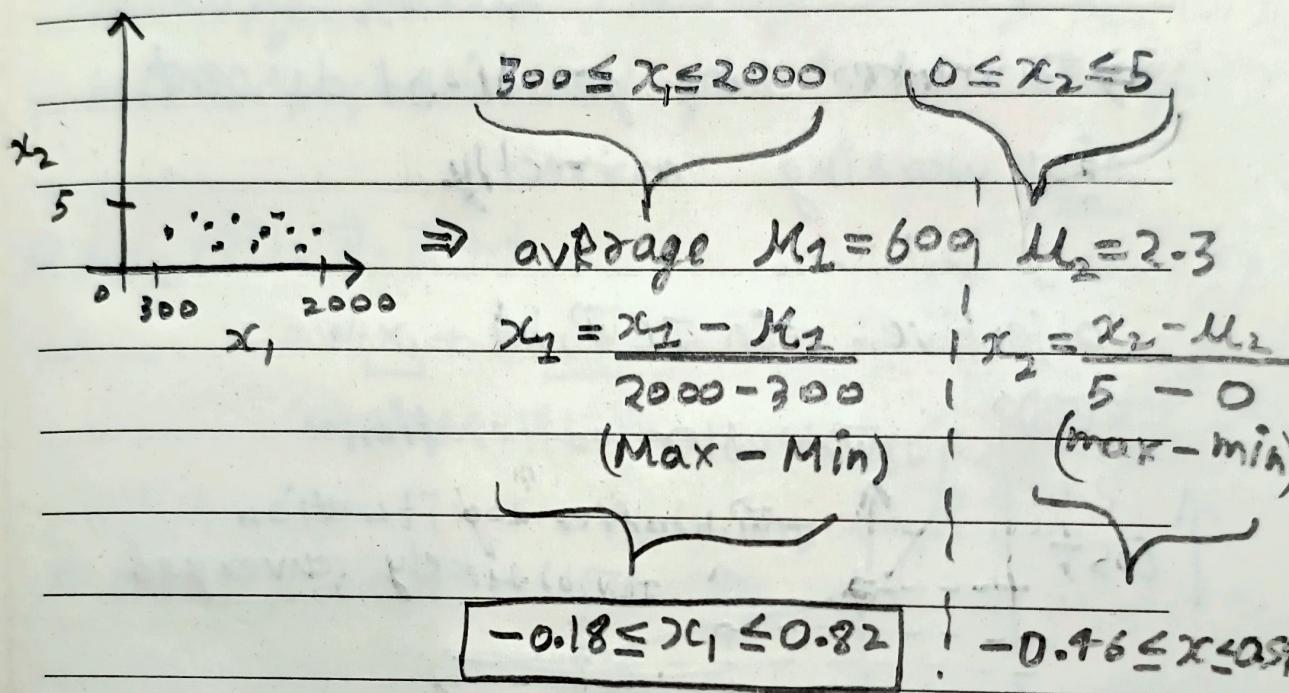


much more faster

accurate than
previous.

How do we do feature scaling?

⇒ One way to do this is using
"Mean normalization"

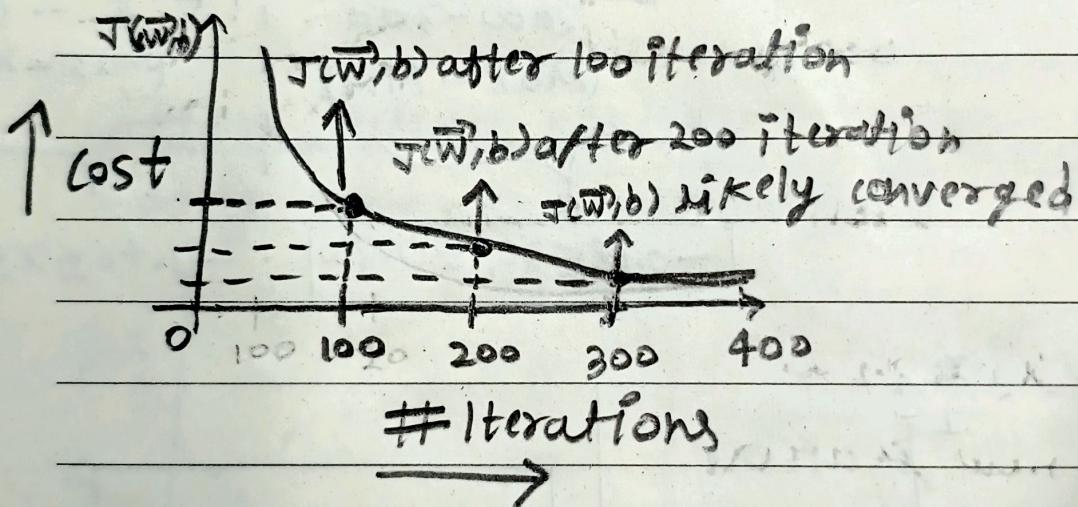


"Mean Normalization"

- Now after learning so many things, how should we know our algorithm (gradient descent) is converging?

\Rightarrow To make sure gradient descent is working correctly

objective: $\min J(\bar{w}, b)$



From above graph we can say that if cost is decreasing at every iteration then our gradient descent working properly, else there is some issue in our gradient descent algorithm.

Feature engineering

- Feature engineering refers to manipulation - addition, deletion, combination, mutation - of your data set to improve machine learning model training, leading to better performance & greater accuracy.

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

$$= w_1 x_1 + w_2 x_2 + b$$

frontage depth



$$\Rightarrow \text{area} = \text{frontage} \times \text{depth}$$

$$x_3 = x_1 x_2$$

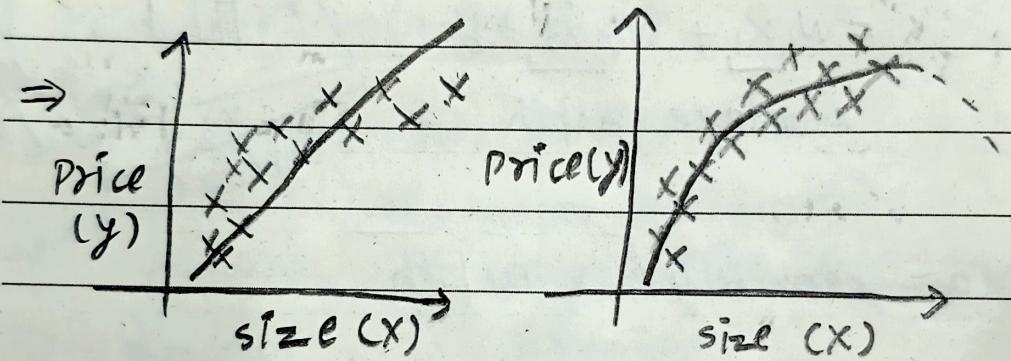
new feature

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

- Feature engineering: Using intuition to design new features, by transforming or combining original features.

Polynomial regression

- Polynomial regression is a form of linear regression where due to the non-linear relationship between dependent & independent variable; we add some polynomial terms to linear regression to convert it into polynomial regression.



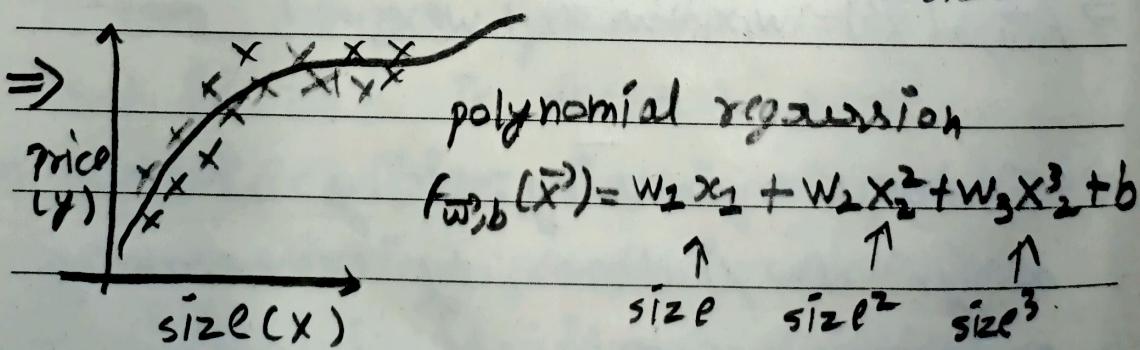
Linear Regression

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

Polynomial Regression

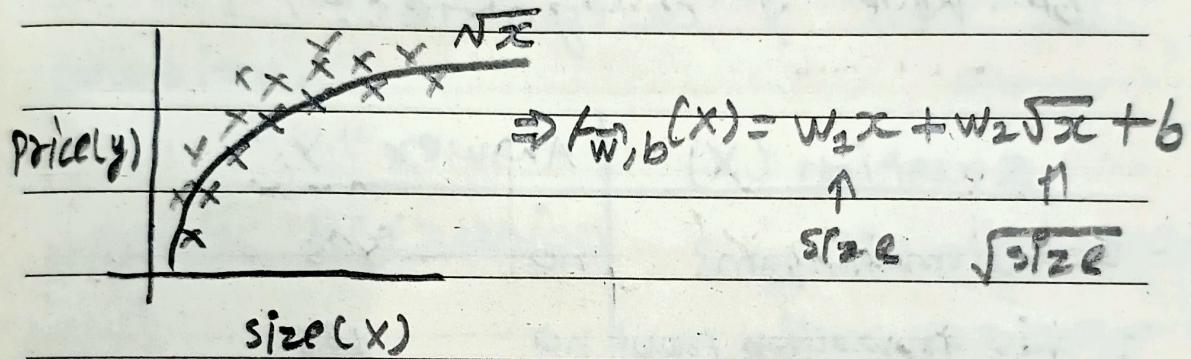
$$f_{w,b}(\vec{x}) = w_1 \cdot x_1 + w_2 \cdot x_2^2 + b$$

size size²



* when we are using polynomial regression then it is very necessary to use feature scaling.

• we can also use another eqn for polynomial suggestion:



$\Rightarrow X_{\text{poly}} = \text{np.column_stack}$ (

[x^{**i} for i in range(degree)]

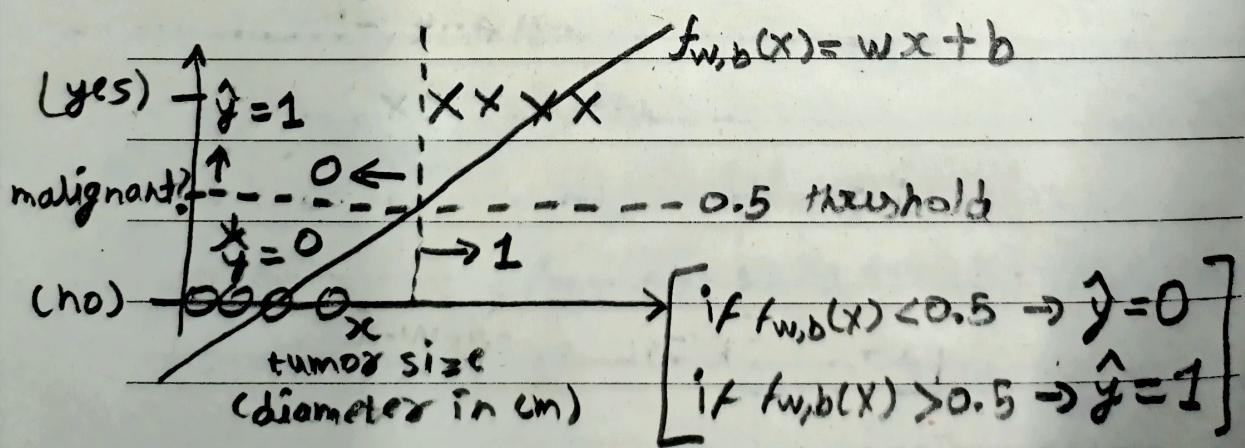
)

Week - 3# Classification with logistic regression
classification

- Classification problem is another type of problem in supervised learning where feature input X we have y categories.

<u>Question (X)</u>	<u>Answer "y"</u>	
Is this email spam?	no	yes
Is this transaction fraud?	no	yes
Is the tumor malignant?	no	yes

- "y" can only be one of two values
- "binary classification" (0 or 1)



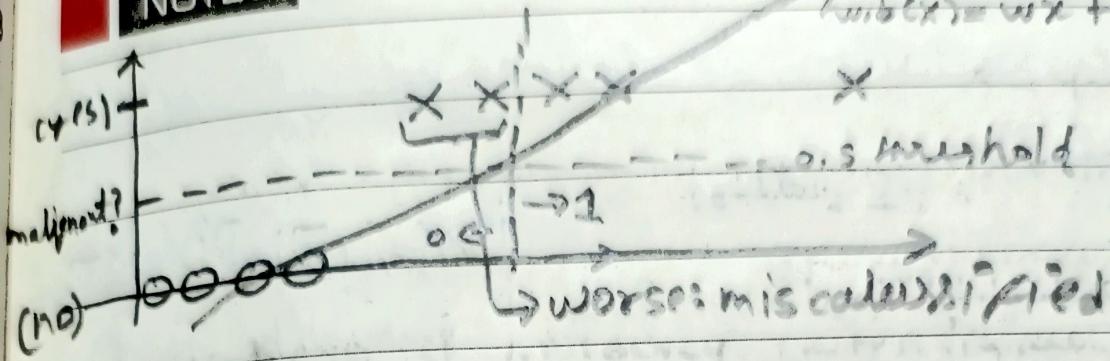
- In above data set our "Linear Regression" is working good.

NOTES

(decision boundary)

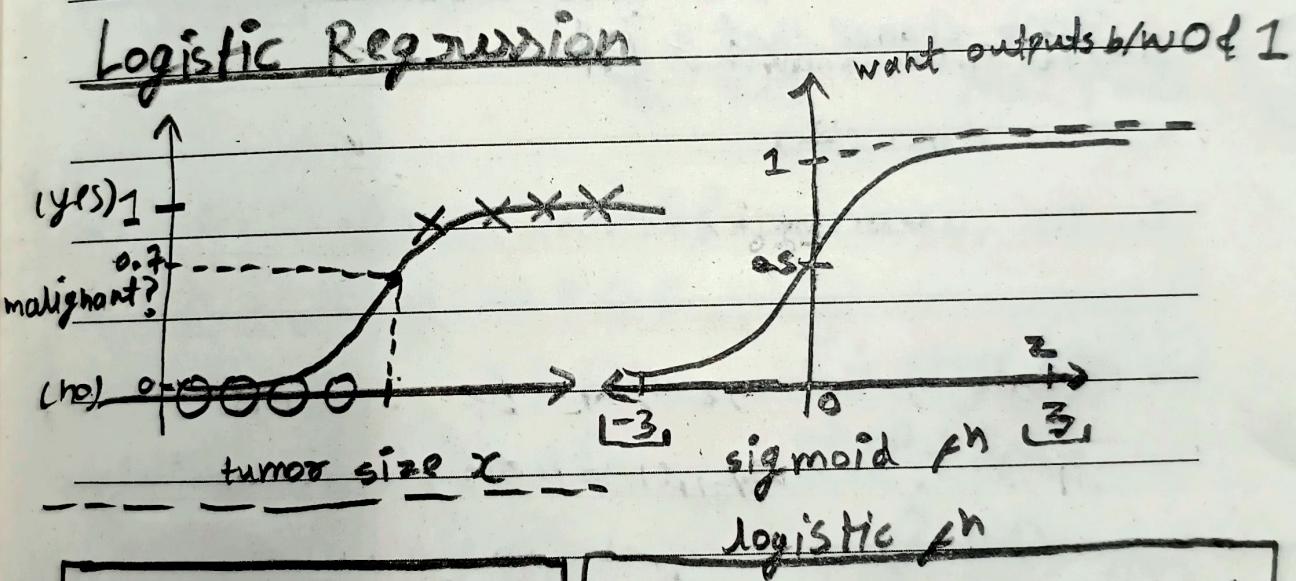
Date :

$$(\text{malignant}) = w_0 + w_1 x_1 + b$$



- Now "Logistic Regression" comes in picture. Logistic Regression is a technique to solve classification problem.

Logistic Regression



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

$$z = \vec{w} \cdot \vec{x} + b$$

$$\downarrow z$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\Rightarrow g(z) = \frac{1}{1 + e^{-z}}, 0 < g(z) < 1$$

outputs b/w 0 & 1 only.

$$\Rightarrow f_{\vec{w}, b}(\vec{x}) = g(\underbrace{\vec{w} \cdot \vec{x} + b}_z)$$

$$\text{"Logistic Regression"} = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

Interpretation of logistic regression output

$$\bullet f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{(\vec{w} \cdot \vec{x} + b)}}$$

"probability" that class is 1

• Example:

x is "tumor size"

y is 0 (not malignant)
or 1 (malignant)

$$\Rightarrow f_{\vec{w}, b}(\vec{x}) = 0.7$$

70% chance that y is 1

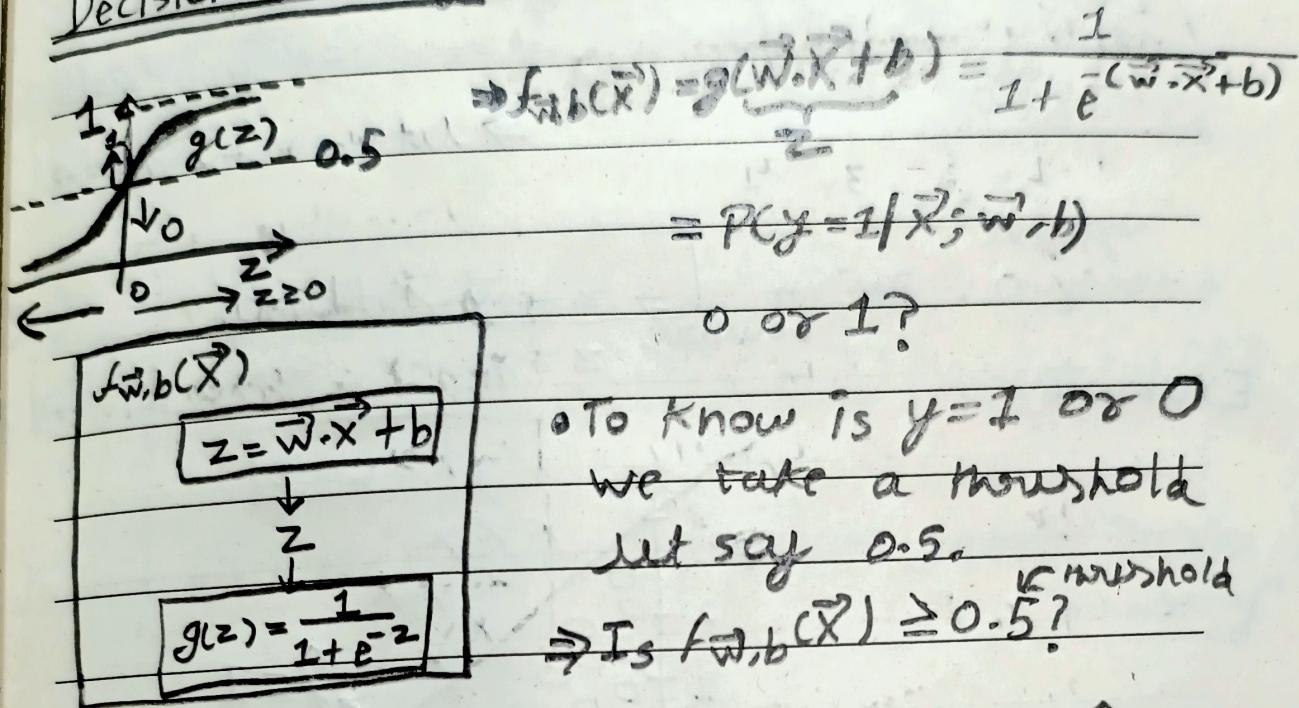
$$\Rightarrow f_{\vec{w}, b}(\vec{x}) = P(y=1|\vec{x}; \vec{w}, b)$$

probability that

y is 1, given
input \vec{x} , parameters

\vec{w}, b .

$$\Rightarrow P(y=0) + P(y=1) = 1$$

Decision boundary

\Rightarrow When is $f_{\vec{w}, b}(\vec{x}) \geq 0.5$?

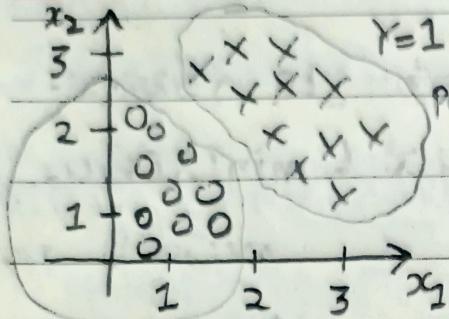
$$\begin{cases} g(z) \geq 0.5 \\ z \geq 0 \\ \vec{w} \cdot \vec{x} + b \geq 0 \end{cases} \quad [\because \text{above graph}] \quad \vec{w} \cdot \vec{x} + b < 0$$

$\hat{Y}=1 \quad \hat{Y}=0$

- In Logistic Regression, "Decision Boundary" is a linear line, which separates the positive & negative points/classes.

NOTES

Date :



$$f_{w,b}(\vec{x}) = g(z)$$

$$= g(w_1x_1 + w_2x_2 + b)$$

$$\Rightarrow \text{let } w_1, w_2 = 1 \text{ and } b = -3$$

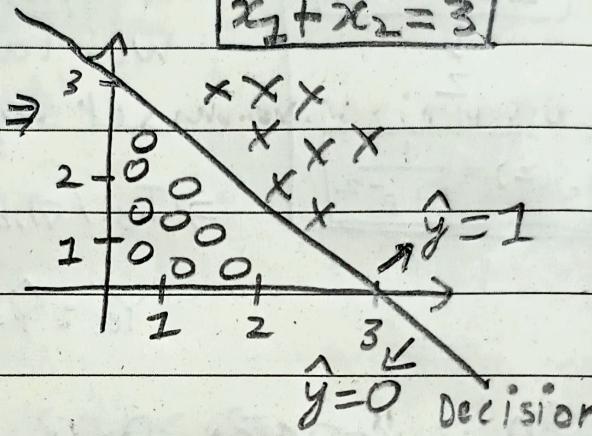
negative

$$Y=0 \quad \text{Decision boundary} \quad z = \vec{w} \cdot \vec{x} + b = 0$$

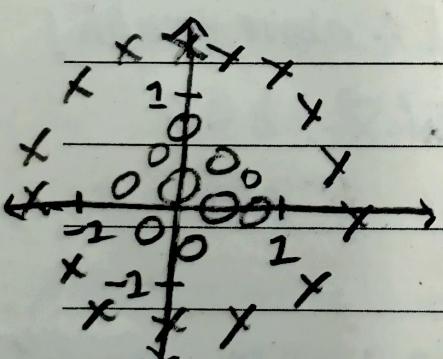
$$z = x_1 + x_2 - 3 = 0$$

Example

$$x_1 + x_2 = 3$$



Non-linear Decision boundaries



$$z = w_1 x_1^2 + w_2 x_2^2 + b$$

$$\Rightarrow f_{w,b}(\vec{x}) = g(z) = g(w_1 x_1^2 + w_2 x_2^2 + b)$$

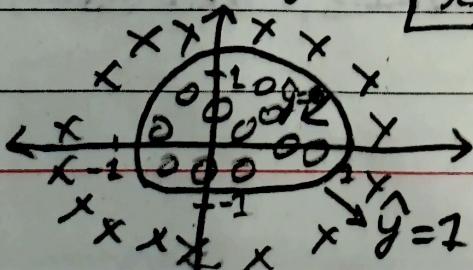
$$\text{let } w_1, w_2 = 1 \text{ and } b = -1$$

$$\Rightarrow z = x_1^2 + x_2^2 - 1$$

$$\text{Decision Boundary} \quad z = x_1^2 + x_2^2 - 1 = 0$$

Boundary

$$x_1^2 + x_2^2 = 1$$



$$\Rightarrow x_1^2 + x_2^2 \geq 1 \quad | \quad x_1^2 + x_2^2 < 1$$

$$\hat{y} = 1 \quad ; \quad \hat{y} = 0$$

outside

inside

code for creating decision boundaries:

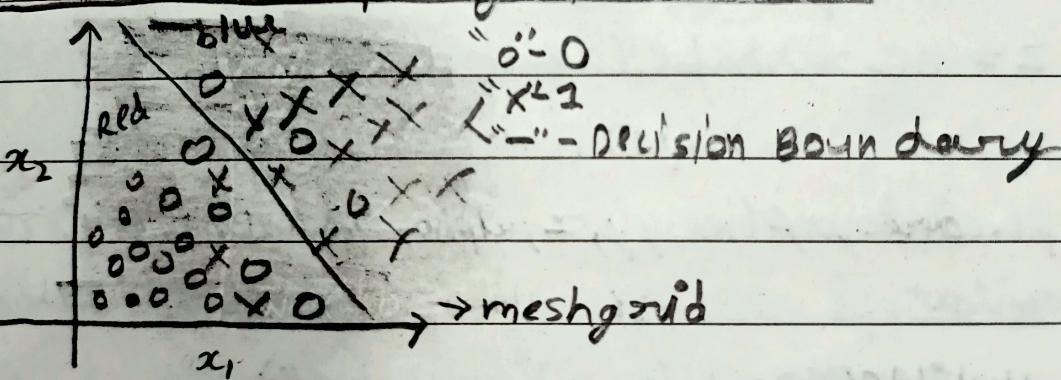
x-values = np.linspace(min(XE), 0],
max(XE), 100)

2-D plot

$y_{\text{values}} = -(w_{\text{-updated}}[0] * x_{\text{values}} + b_{\text{-updated}}) / w_{\text{-updated}}[1]$

```
plt.plot(x-values, y-values, label="Decision boundary")
```

If we want to plot graph like this:



`x-values = np.linspace(min(x[:,0]), max(x[:,0]), 100)`

$y\text{-value} = -(w[0] * x\text{-value} + b) / w\text{-updated}[1]$

```
plt.plot(x-values, y-values, label="DB", color="blue")
```

scatter plot of data points.

import seaborn as sns
 sns.scatterplot(x="x1", y="x2", hue="class", data=df)

Scatter library

meshgrid for more attractive visualization
 $x_{\min}, x_{\max} = \text{data}['x1'].min() - 1, \text{data}['x1'].max() + 1$
 $y_{\min}, y_{\max} = \text{data}['x2'].min() - 1, \text{data}['x2'].max() + 1$
 $xx, yy = \text{np.meshgrid}(\text{np.arange}(x_{\min}, x_{\max}, 0.1),$
 $\text{np.arange}(y_{\min}, y_{\max}, 0.1))$

$Z = \text{np.dot}(\text{np.c_}[xx, yy], w) + b$

$Z = \text{np.sign}(Z)$

$z = \text{np.reshape}(xx, shape)$

plt.contourf(xx, yy, z, alpha=0.3, color=['blue', 'red'])

plt.title("Decision boundary")

plt.xlabel("Feature 1 (x1)")

plt.ylabel("Feature 2 (x2)")

plt.legend()

plt.grid(True)

Cost Function for Logistic Regression.

Training set

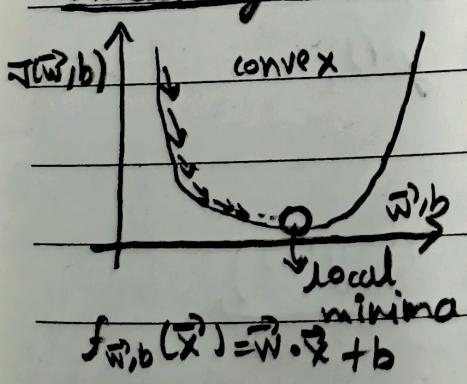
	tumor size (cm)	...	Patient's age malignant?	$i=1, \dots, m$	train example
$i=1$	10	...	52	1	$j=1, \dots, n$ feature
:	2	...	73	0	forget y is
	5	...	55	1	0 or 1
$i=m$	$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

⇒ How to choose $\vec{w} = [w_1, w_2, \dots, w_n]$ & b ?

Squared error cost J^h

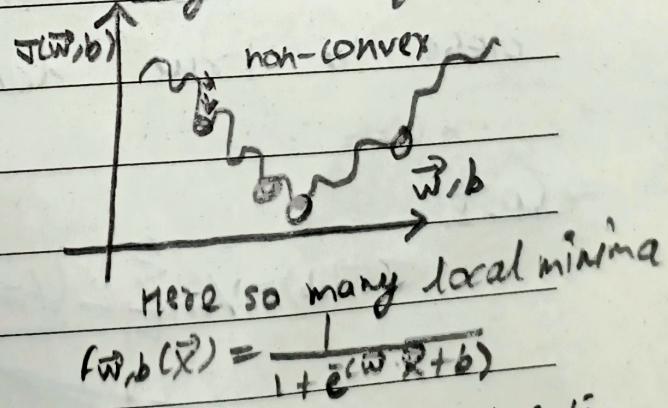
$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$$

Linear Regression



- If we are using squared error cost J^h for logistic regression, it will create so many local minima which will create the problem.

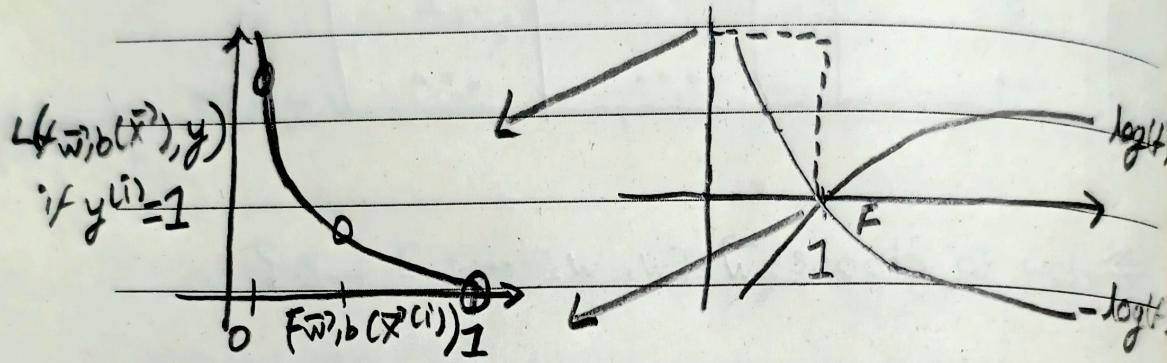
Logistic Regression



• So, for logistic regression we have loss function.

Logistic Loss function

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)}=1 \\ -\log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)}=0 \end{cases}$$



As $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 1$ then loss $\rightarrow 0$ //

As $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 0$ then loss $\rightarrow \infty$ //

\Rightarrow Loss is lowest when $f_{\vec{w}, b}(\vec{x}^{(i)})$ predicts close to true label $y^{(i)}$.

Cost

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}),}_{\text{Loss}}$$

$$= \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)}=1 \\ -\log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)}=0 \end{cases}$$

it can reach a global minimum (convex)

Simplified loss function.

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)}=1 \\ -\log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)}=0 \end{cases}$$

$$\Rightarrow L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)}))$$

let $y^{(i)}=1,$

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -\log(f_{\vec{w}, b}(\vec{x}^{(i)}))$$

proof $\underline{\underline{=}}$

let $y^{(i)}=0,$

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -\log(1-f_{\vec{w}, b}(\vec{x}^{(i)}))$$

Simplified cost function

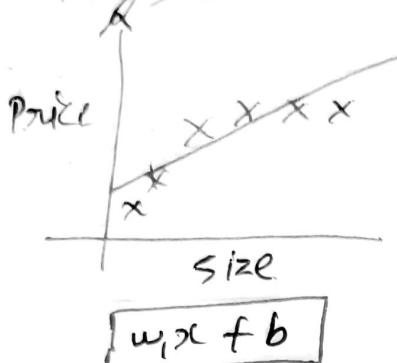
$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})]$$

$$\Rightarrow J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)}))]$$

↑
—
(Loss function)

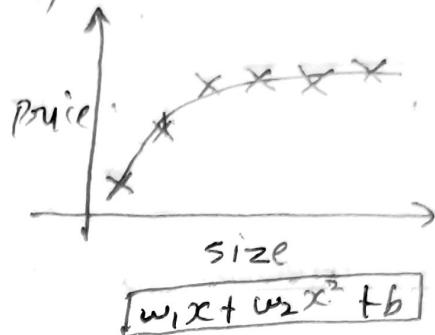
The problem of overfitting

Regression Example



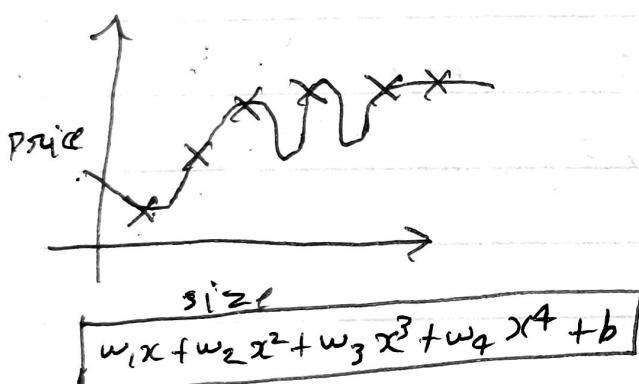
Under fit

- Does not fit the training set well
→ high bias



Just right

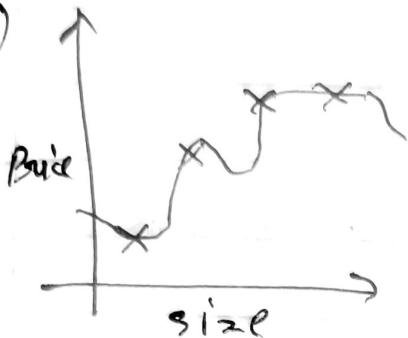
- Fits training set pretty well
→ Generalization



- Fits the training set extremely well
Over fit
→ high variance

Addressing overfitting

①



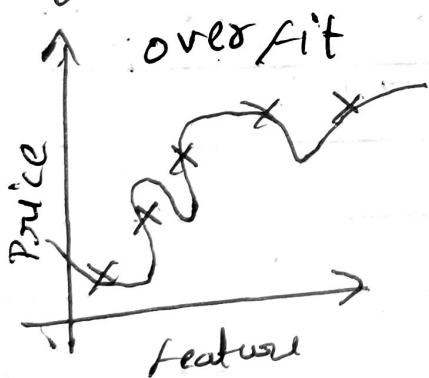
② select feature to include/exclude

All feature

↓
insufficient data
↓
overfit

selected feature
+
insufficient data
↓
just right.

③ * Regularization. (Reduce the size of parameters)



$$f(x) = 28x - 385x^2 + 39x^3 - 174x^4 + 100$$

large values for w_j



$$f(x) = 13x - 0.023x^2 + 0.00014x^3 - 0.001x^4 + 10$$

small feature for w_j value

summarise

options:

- 1). collect more data
- 2). select features
- 3). Reduce size of parameters

Cost function for Regularization

$w_1, w_2, w_3, \dots, w_{100}, b$ n features $n=100$
regularization term

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

"lambda" $\lambda > 0$

regularization parameter

mean squared error

regularization term

$$\Rightarrow \min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right]$$

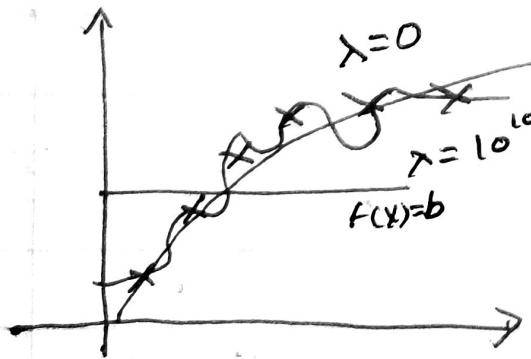
fit data , keep w_j small

$\lambda=0$ $\lambda=10^{-10}$ λ balance both goals

choose $\lambda=0$, overfit

choose $\lambda=10^{-10}$, underfit

choose λ in b/w , just right



Regularized linear regression

cost function:

$$\min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right]$$

gradient descent

repeat {

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

$$= \frac{1}{m} \sum_{i=0}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

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

$$= \frac{1}{m} \sum_{i=0}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

} simultaneous update

don't have to regularize b .

Implementation of gradient descent

repeat {

$$w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=0}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j \right]$$

$$b = b - \alpha \left[\frac{1}{m} \sum_{i=0}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \right]$$

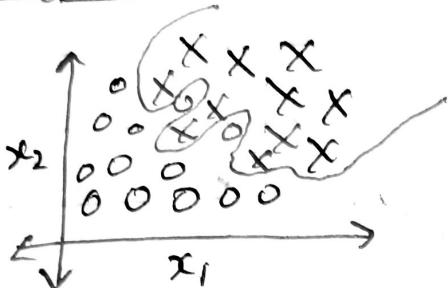
} simultaneous update

option {

$$\text{rewrite: } w_j = w_j \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_i^{(i)} \right]$$

just of the previous code.

Regularized Logistic Regression



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + \dots + b$$

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-z}}$$

cost J^n :

$$\begin{aligned} J(\vec{w}, b) &= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) \right] \\ &\quad + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{regularized term}} \end{aligned}$$

simple cost J^n

Gradient descent

repeat {

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

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

} simultaneous update

$$= \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

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

$$\Rightarrow f_{\vec{w}, b} = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}} = \frac{1}{1 + e^{-z}}$$