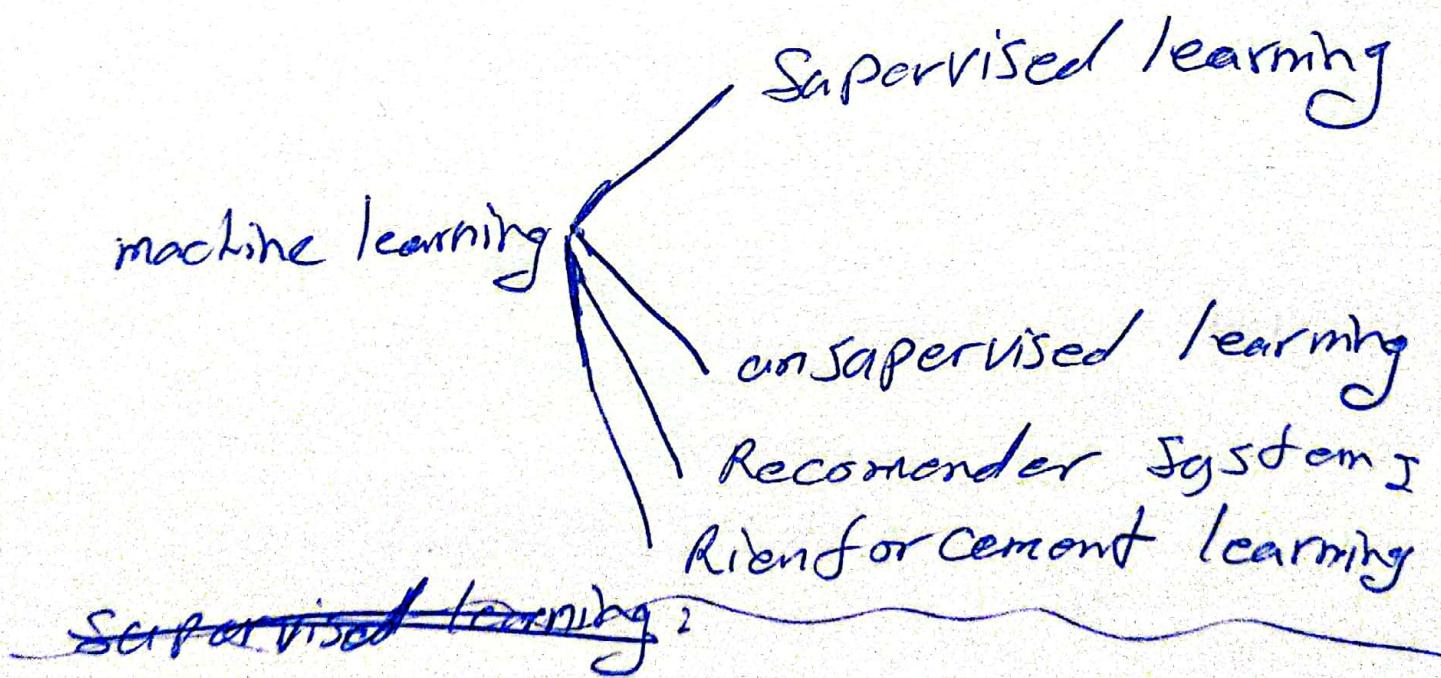


"Supervised machine learning" ①

(regression and classification)

machine learning:

"field of study that gives computers the ability to learn without being explicitly programmed" Arthur Samuel (1959)



Supervised learning:

$$X \longrightarrow Y$$

Input \longrightarrow output label

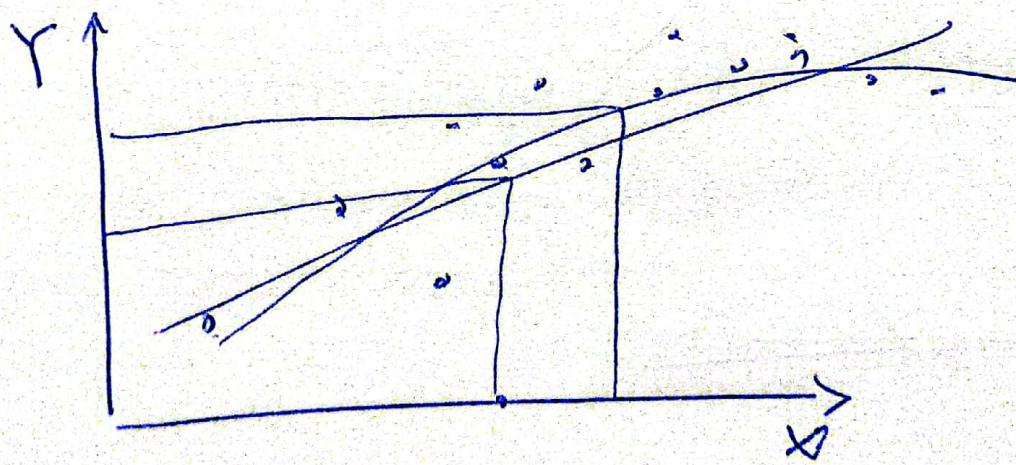
Learn from being given "right answers"

②

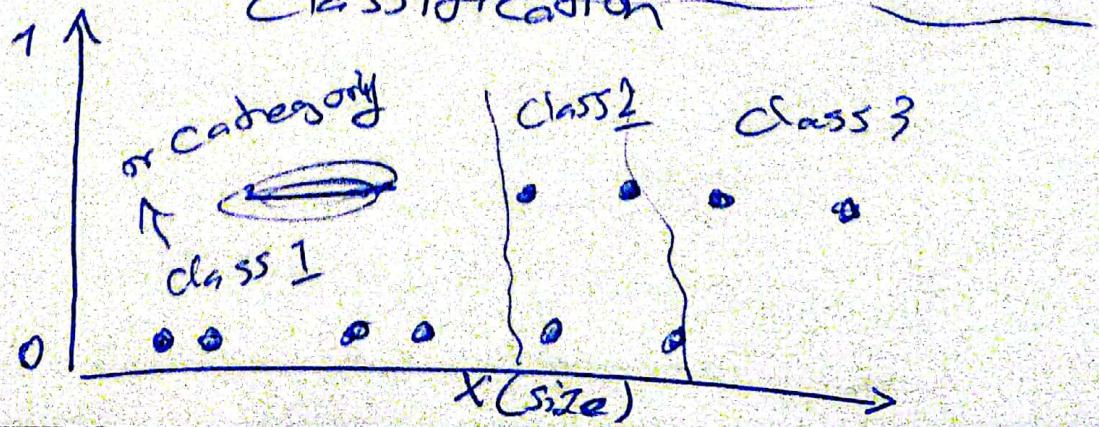
Supervised learning

| in part (X) | out part (Y) | application |
|-------------------|--------------------------------------|---------------------|
| email | \rightarrow spam (0/1) | spam filtering |
| audio | \rightarrow text transcript | speech recognition |
| English | \rightarrow Spanish | machine translation |
| ad user info | \rightarrow click? (0/1) | online advertising |
| image, radar info | \rightarrow position of other cars | self-driving car |
| image of phone | \rightarrow defect? (0/1) | visual inspection |

regression

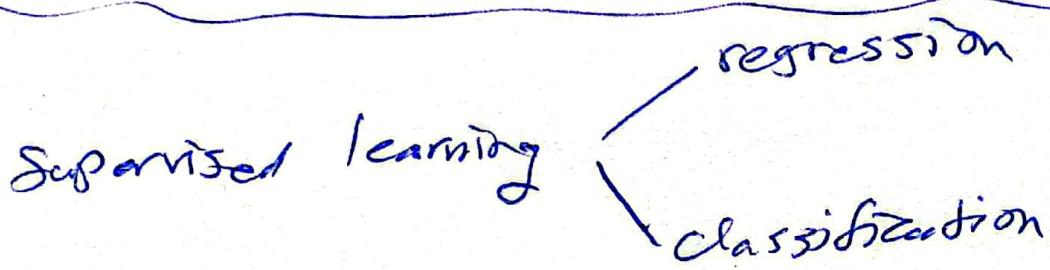
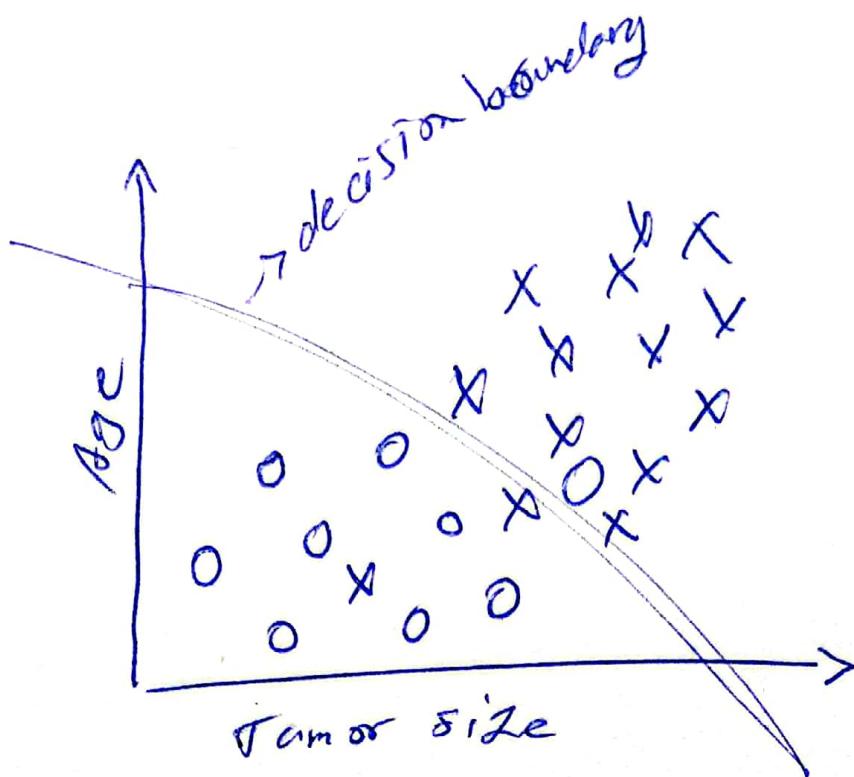


Classification

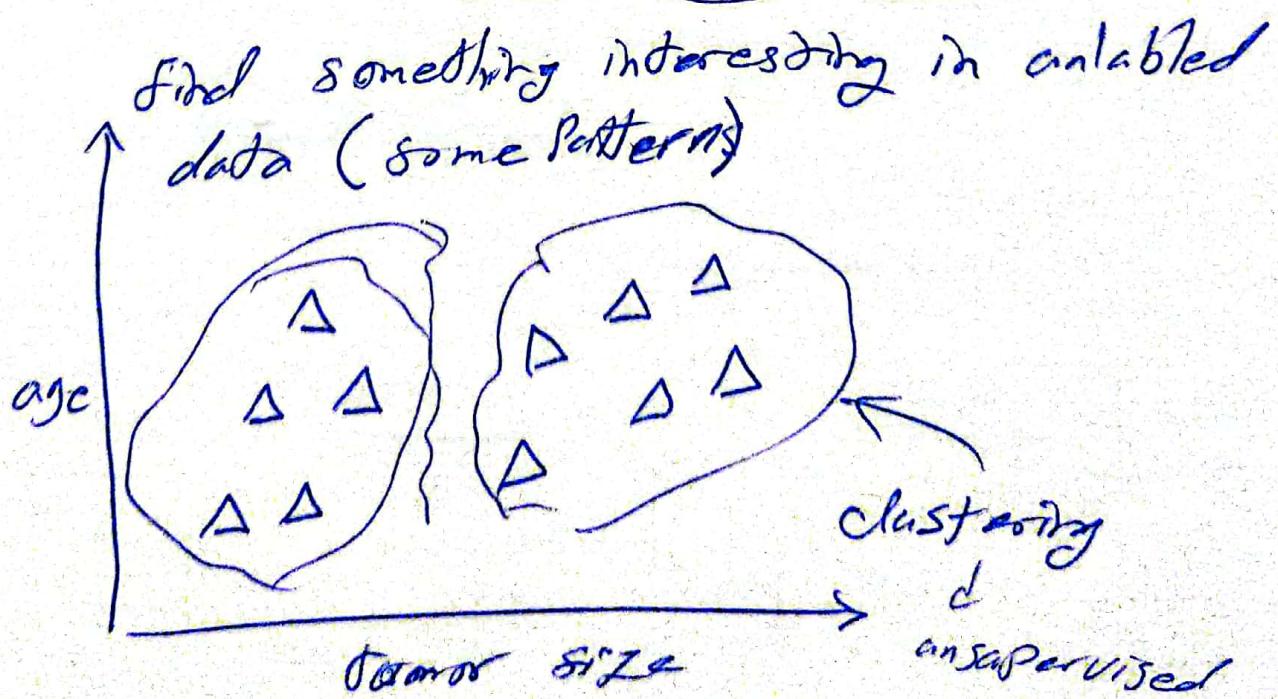


exp2: classification (2X5)

③



unsupervised learning

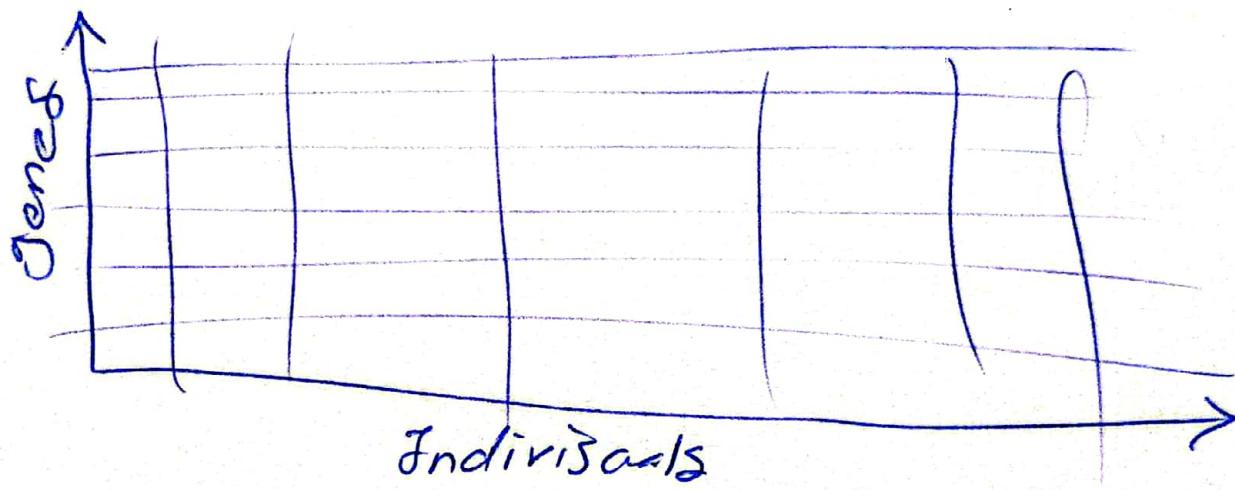


Exps of unsupervised learning

(4)

(suggestion) Google news (clustering)

- clustering: DNA microarray



- clustering: Grouping customers

unsupervised learning

Data only comes with inputs X , but not output labels y . algorithm has to find structure in the data.

Clustering

Group similar data points together

Anomaly detection

Find unusual data points.

Dimensionality reduction

Compress data using fewer numbers

Linear regression model

(5)

Terminology

- Training set: Data ~~set~~ used to train the model.
- x (lower case): input variable or feature
- y (lower case): output variable or target variable

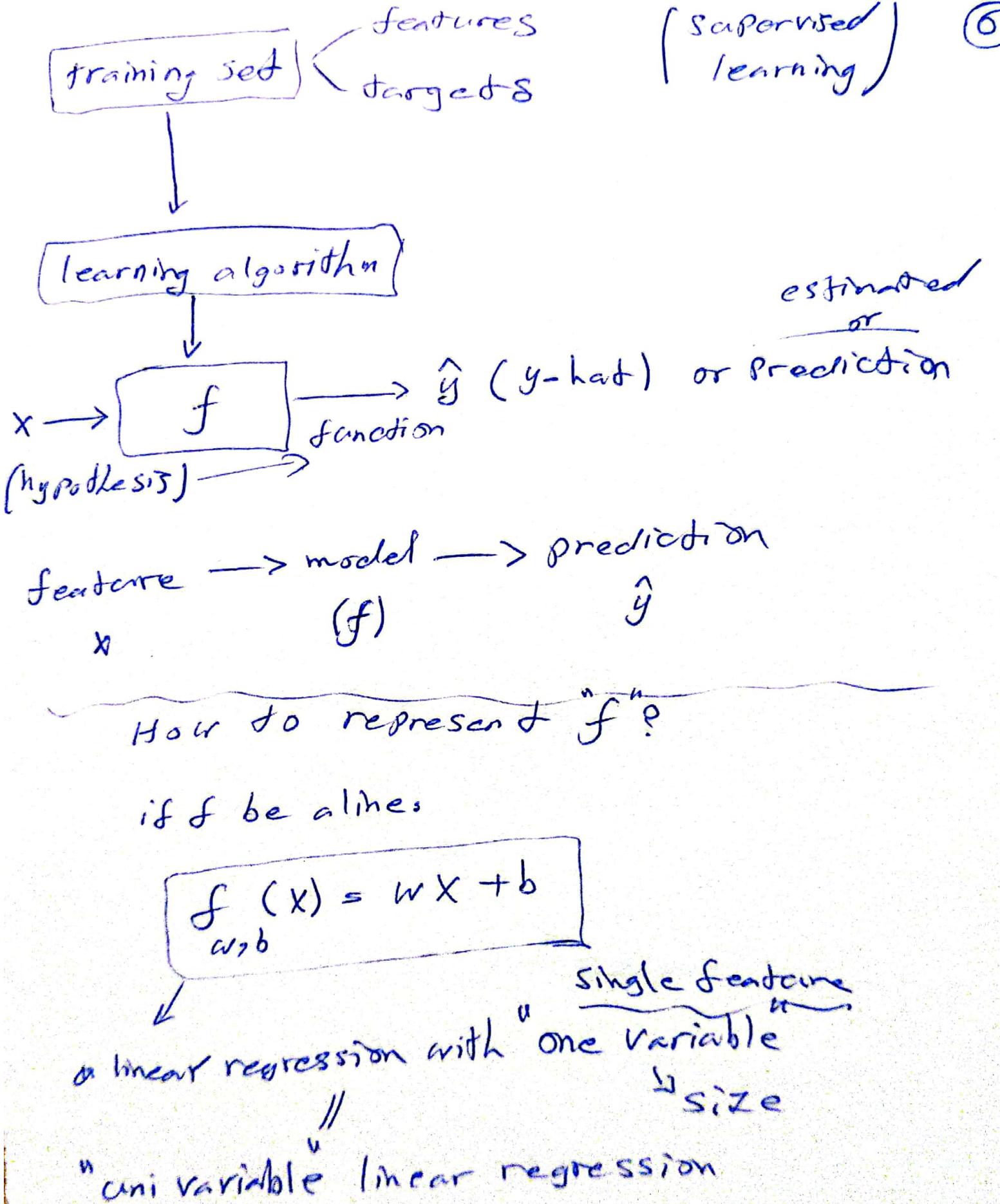
| (X) size of a foot | (Y) price in \$ |
|-------------------------|----------------------|
| 2014 | 400 |
| 1415 | 232 |
| : | : |
| : | : |

- m (lower case): number of training examples.
 (x, y) : single training example.

$(x^{(i)}, y^{(i)})$ = i th training example.

\downarrow
refers to a specific row in the table

exp: $(x^{(1)}, y^{(1)}) = (2014, 400)$



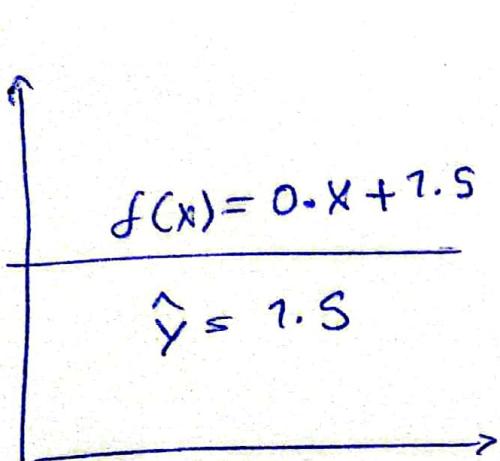
7

Cost function

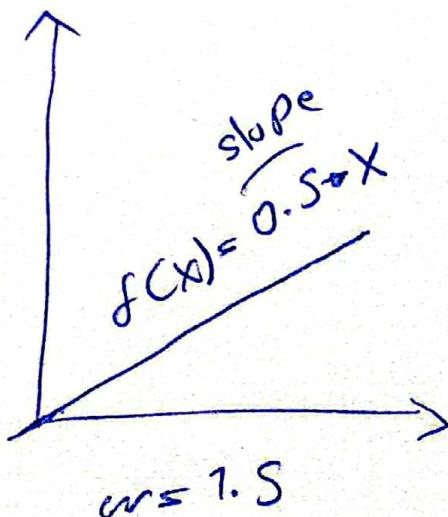
model: $f_{w,b}(x) = wx + b$

w and b are parameters.
coefficients

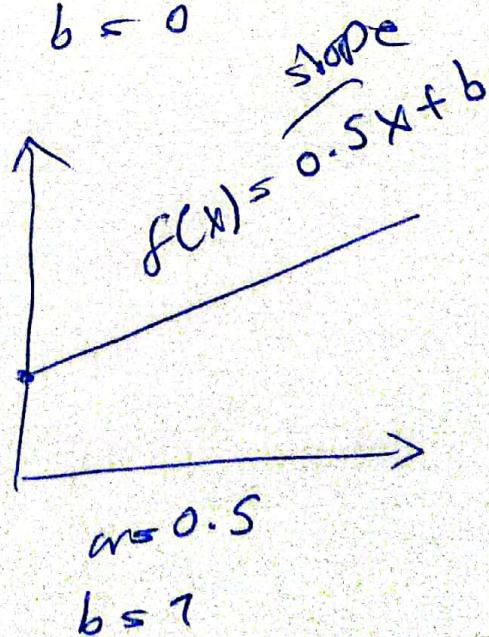
or
 weights



$w = 0$
 $b = 1.5$
 intercept



$w = 1.5$
 $b = 0$



$w = 0.5$
 $b = 1$

(8)

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

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

cost function: (squared error)

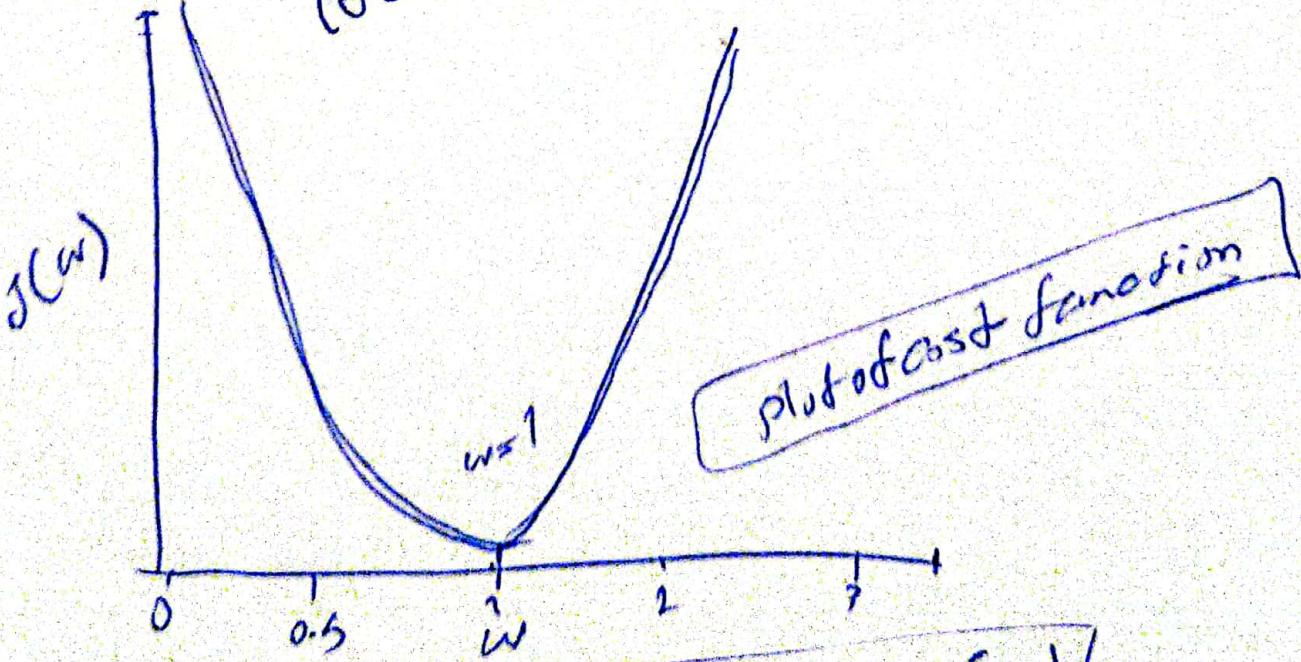
$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m \underbrace{(\hat{y}^{(i)} - y^{(i)})^2}_{\text{error}}$$

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

goal: minimize $J(w, b)$

w, b

(function of w)



choose w to minimize $J(w)$

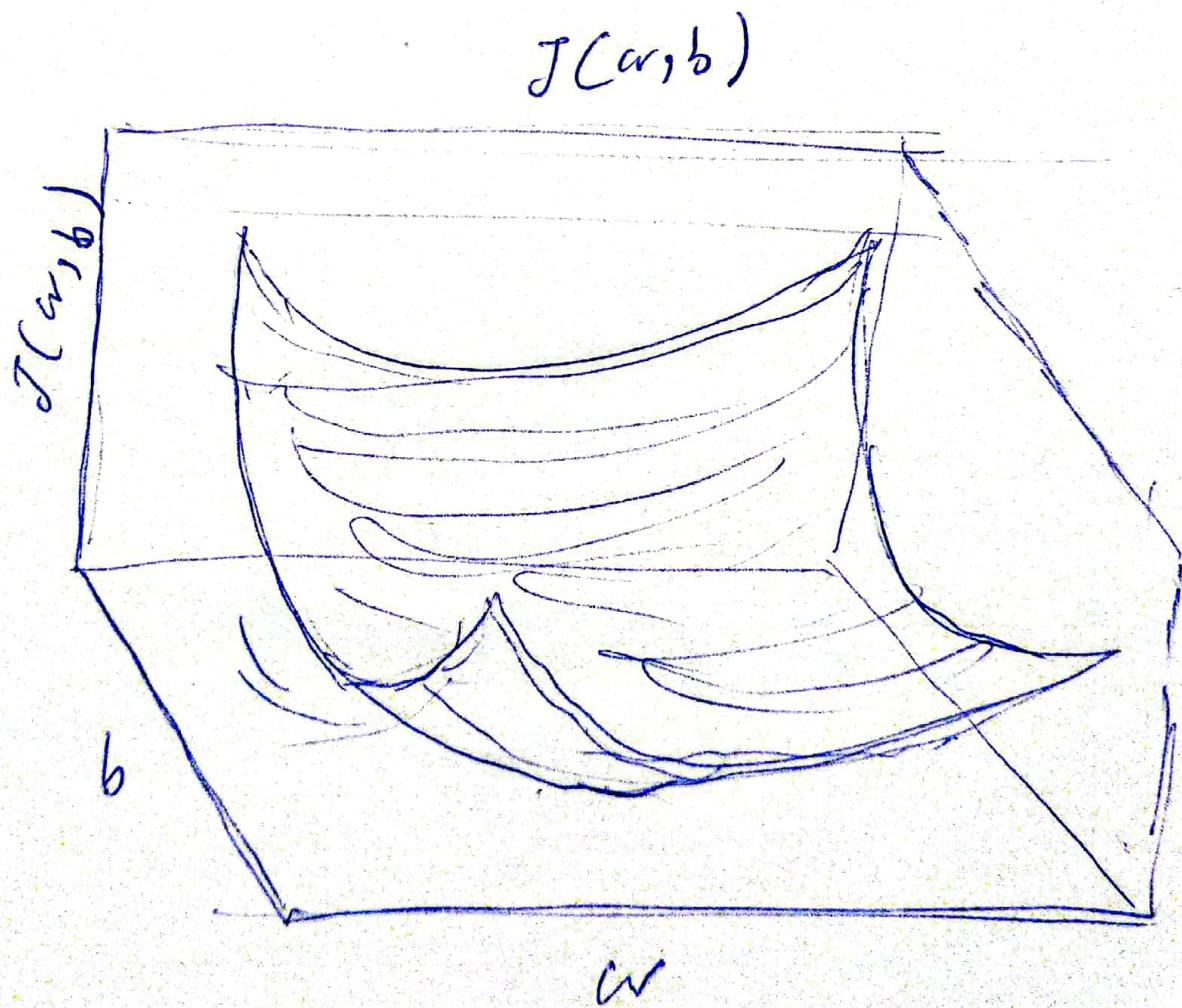
Visualizing the cost function ⑨

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

Parameters: w, b

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

objective: minimize $J(w,b)$
 w, b



Gradient Descent

cost function: $J(w, b)$

want to: $\min_{w, b} J(w, b)$

In general: $J(w_1, w_2, \dots, w_n, b)$

goal: $\min_{w_1, w_2, \dots, w_n, b} J(w_1, w_2, \dots, w_n, b)$

out line:

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

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

until we settle at ~~the~~ or near a minimum.

may have \rightarrow 1 minimum

Gradient Descent algorithm

(11)

Assignment sign

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

α = learning rate (a small positive number)

$\frac{d}{dw} J(w, b)$ = derivative of cost function with respect to weight.

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

simultaneously update w and b

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

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

↓

$$w = \text{tmp_}w$$

$$b = \text{tmp_}b$$

Gradient Descent is an algorithm for finding values of parameters w and b that minimize the cost function ($J(w, b)$)

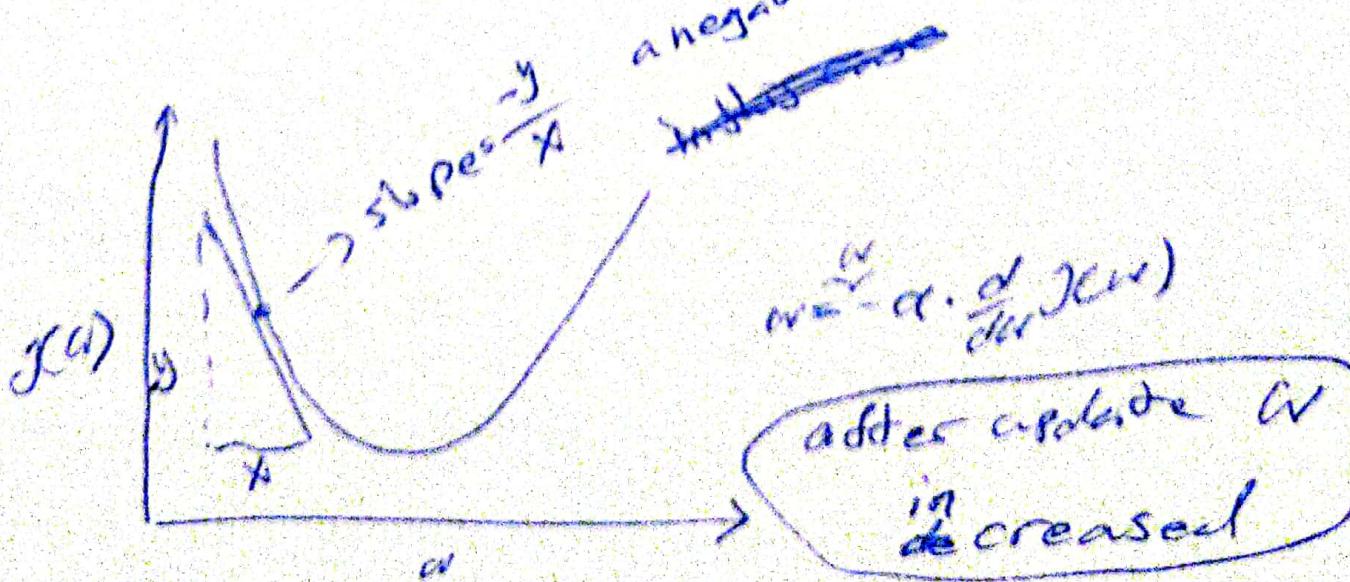
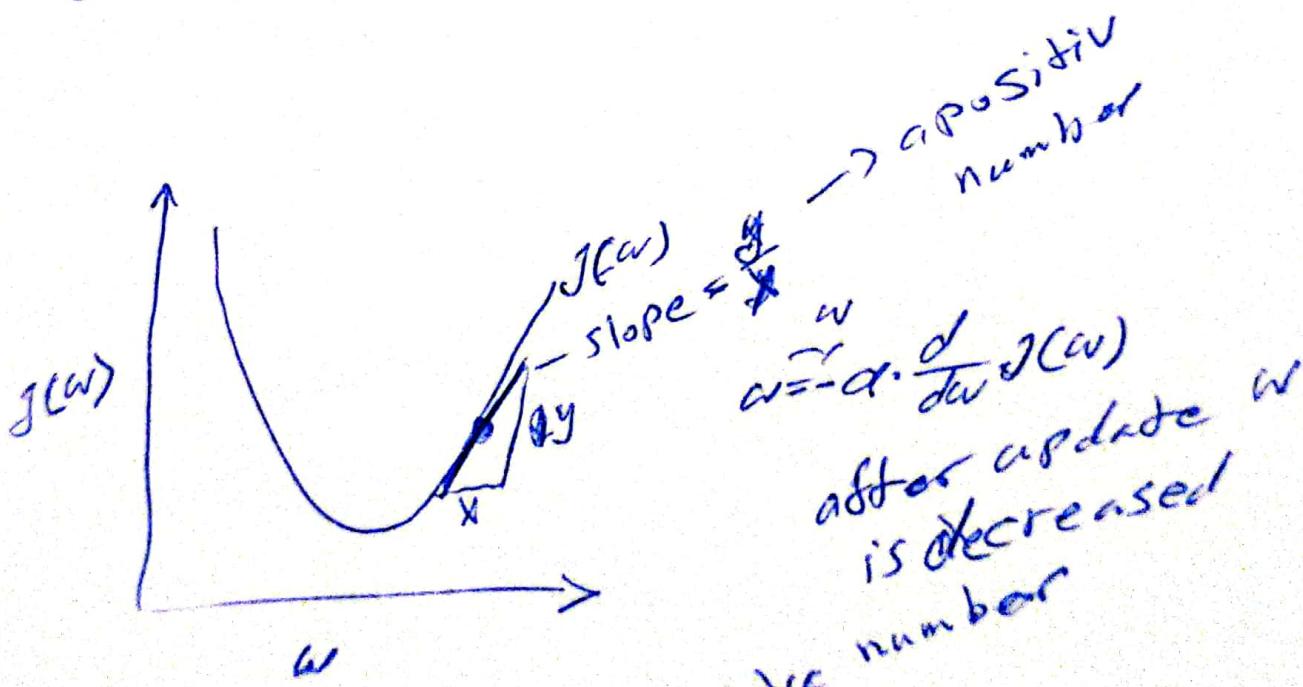
"Gradient Descent intuition"

repeat until convergence {

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

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

}

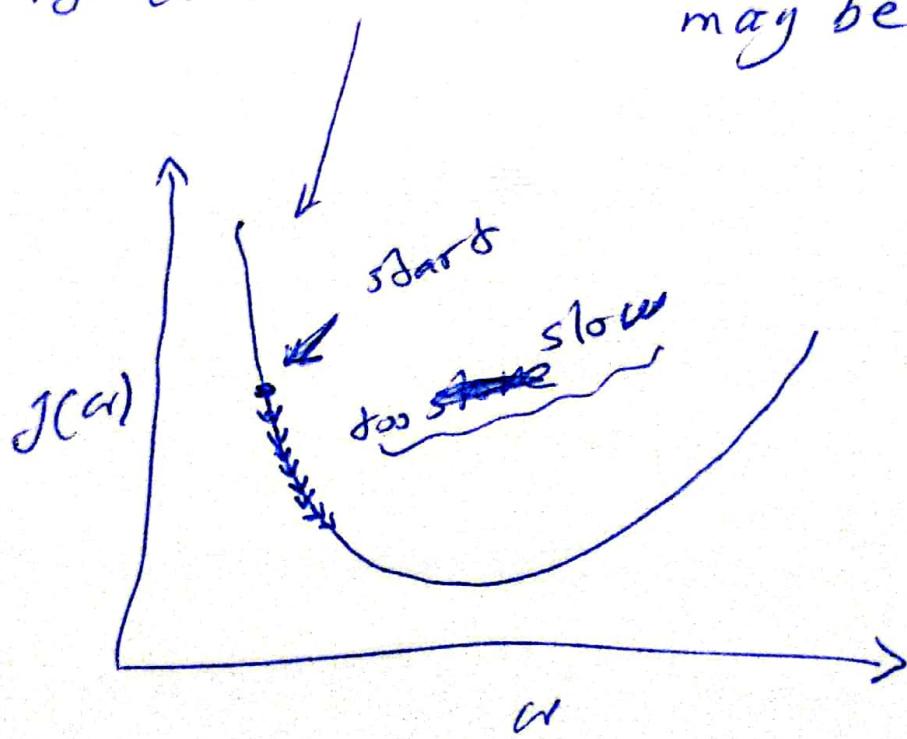


Learning Rate (α)

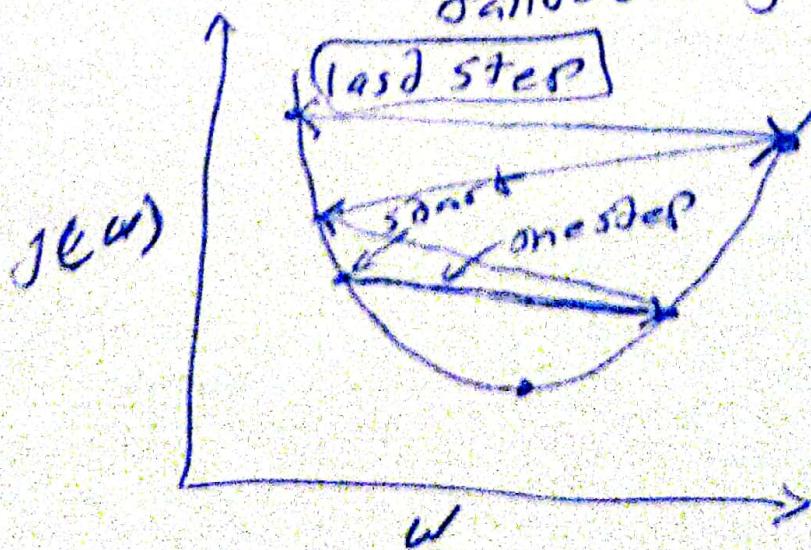
(13)

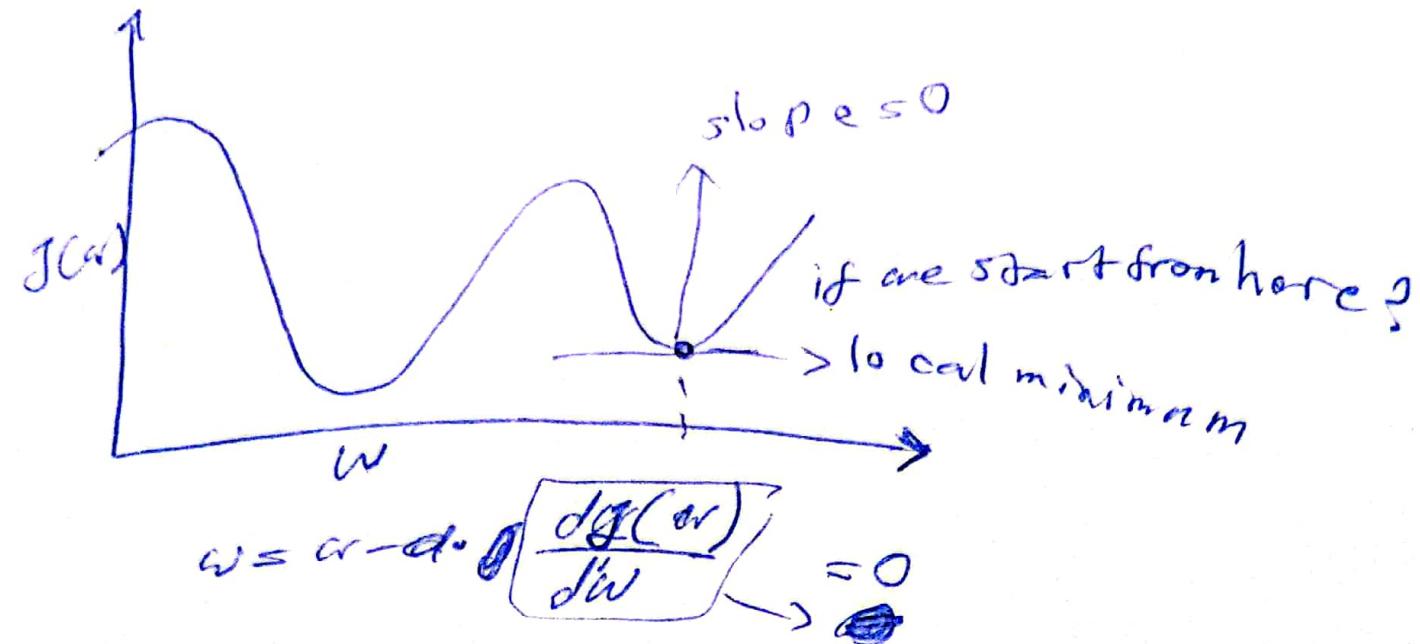
$$w = w - \alpha \cdot \frac{d}{dw} J(w)$$

if α is too small: Gradient descent may be slow.



if α is too large: Gradient descent may overshoot, never reach minimum. fail to converge or may diverge





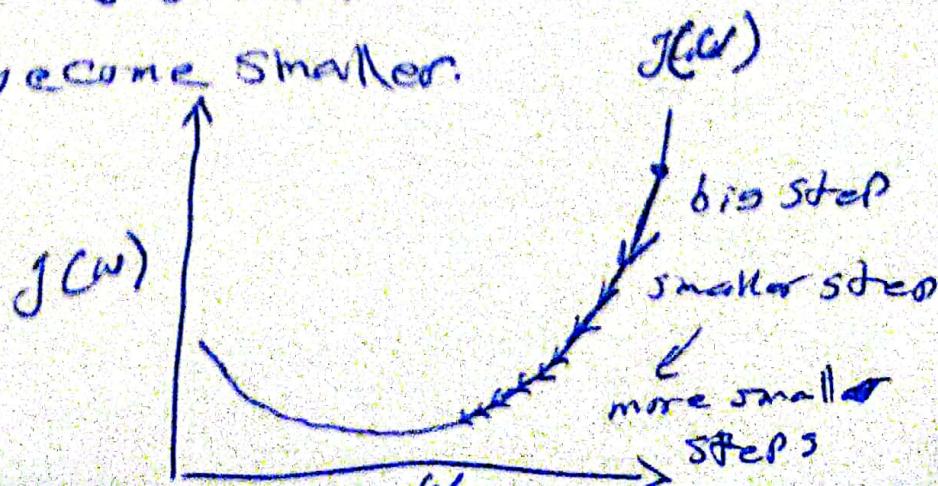
so

$$w = w - \alpha \cdot 0$$

$$\Delta w = w - w$$

it means if we start from a local minimum
gradient descent algorithm changes nothing

Even if α is a fixed number
updated steps change over iteration
 α near a local minimum
derivative becomes smaller
so update steps become smaller.



Gradient Descent for Linear

(15)

Regression.

Linear regression: $f_{w,b}(x) = wX + b$

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

Gradient descent:

repeat until convergence

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

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

?

Partial derivatives

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

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

"Partial derivatives calculating" (10)

$$\begin{aligned}
 \frac{\partial}{\partial w} J(w, b) &= \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - g^{(i)})^2 \\
 &= \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (\omega x^{(i)} + b - g^{(i)})^2 \\
 &= \frac{1}{2m} \sum_{i=1}^m (\omega x^{(i)} + b - g^{(i)}) \cancel{2} x^{(i)} \\
 &\boxed{= \frac{1}{m} \sum_{i=1}^m (\omega x^{(i)} + b - g^{(i)}) x^{(i)}}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial b} J(w, b) &= \frac{\partial}{\partial b} \cdot \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - g^{(i)})^2 \\
 &= \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (\omega x^{(i)} + b - g^{(i)})^2 \\
 &= \frac{1}{2m} \sum_{i=1}^m (\omega x^{(i)} + b - g^{(i)}) \cancel{2} \\
 &\boxed{= \frac{1}{m} \sum_{i=1}^m (\omega x^{(i)} + b - g^{(i)})}
 \end{aligned}$$

Gradient Descent algorithm

(17)

repeat until convergence {

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

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

}

if you use ~~the cost function~~
with linear regression, (squared error cost)

it does not have multiple local minimum,
it has in this case just a "single"
"global minimum"

Batch gradient descent (normal gradient descent)
refers to ~~use~~ use all training set to
update the parameters in each iteration.

multiple features (variable)

(18)

$x_j = j^{\text{th}}$ feature

$n = \text{number of features}$

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

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

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

$w = [w_1, w_2, \dots, w_n]$ a row vector of weights.

b is a number

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

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

dot product

multiple linear regression

Gradient descent for
multiple linear regression

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

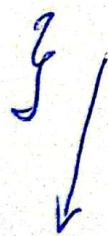
$$J(\vec{w}, b)$$

\leftarrow cost function

repeat {

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

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



repeat {

$$w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - g^{(i)}) x_1^{(i)}$$

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

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

simultaneously update

w_j (for $j=1, \dots, n$) and b

→ normal equation

only for linear regression

solve for w, b without iteration.

disadvantages:

- does not generalize to other learning algorithms.
- slow when number of features is large ($> 10,000$)

Feature scaling

mean normalization

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

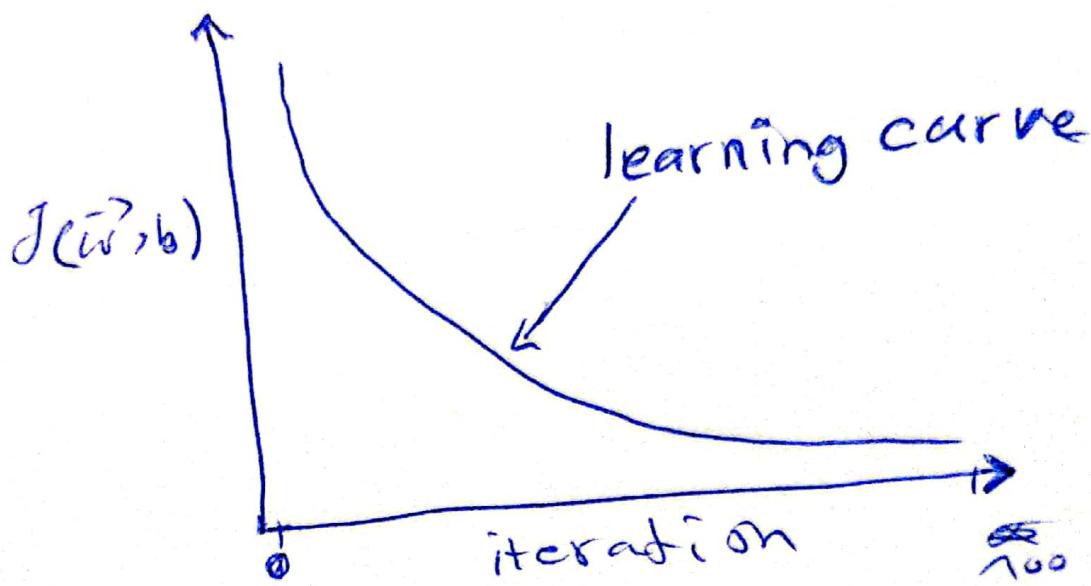
I-score normalization

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

How to check if GD converges?

(27)

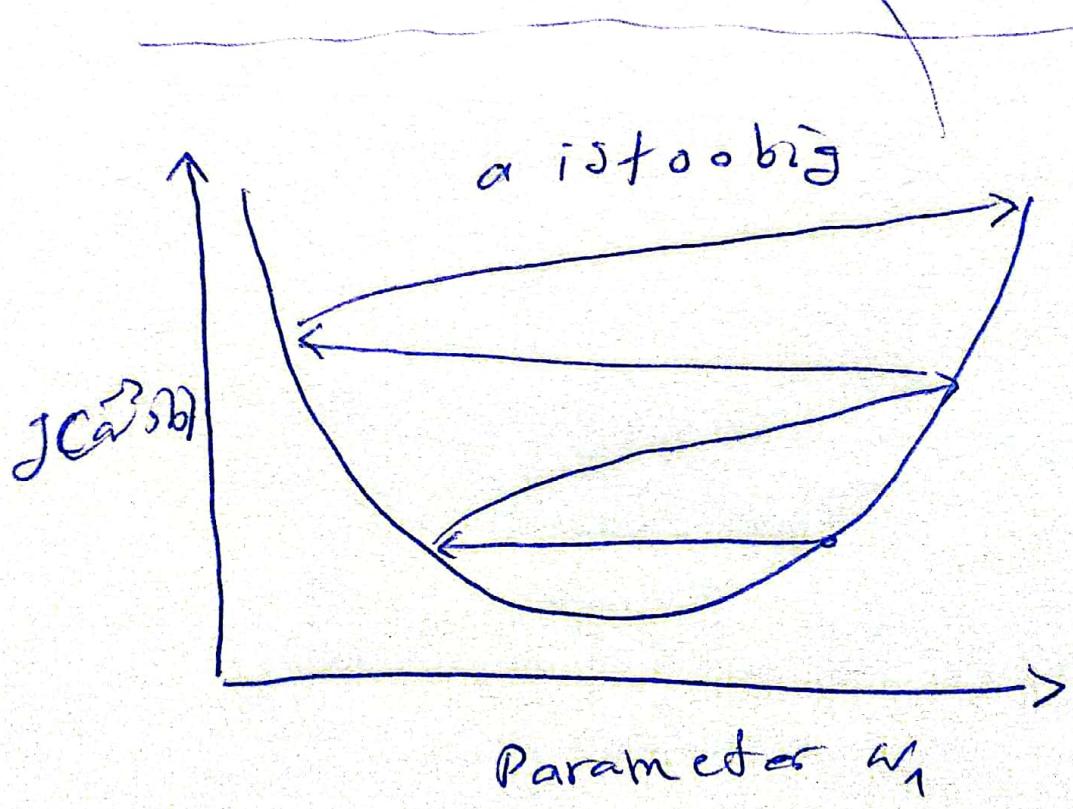
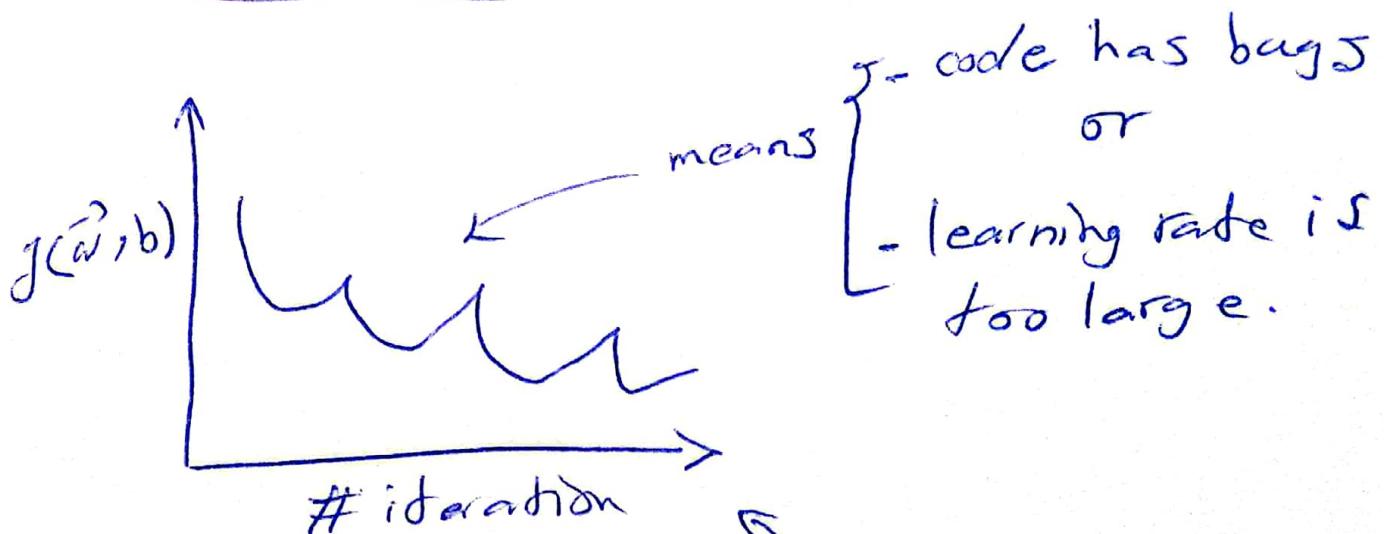
plot the cost function:



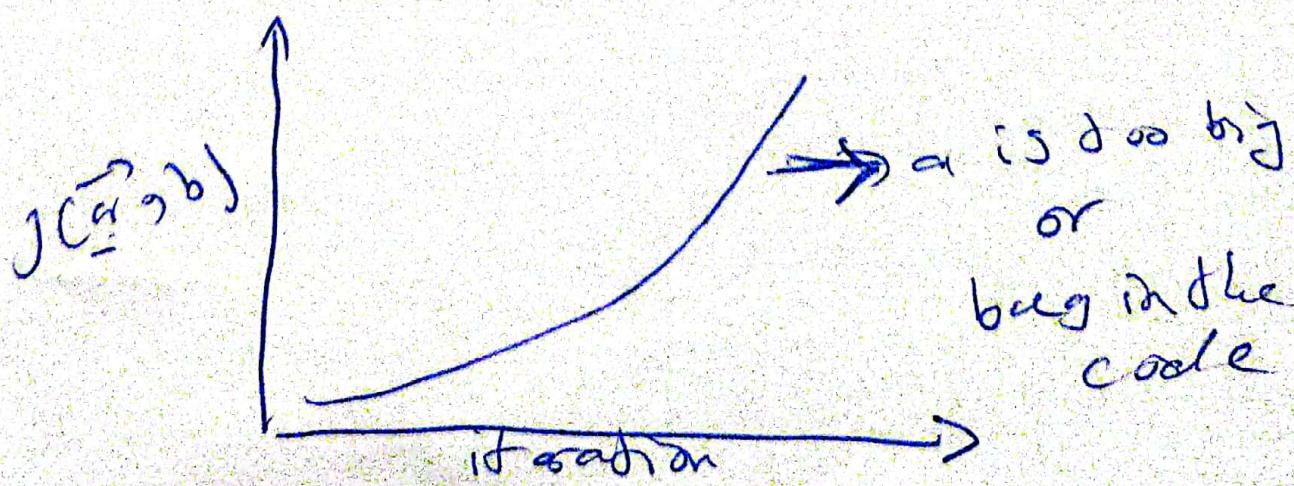
Automatic convergence test

if $J(\bar{w}, b)$ decreases by $\leq \epsilon$
in one iteration, declare
convergence.

Identify problem with gradient descent (22)



use a smaller learning rate.



- with a small enough α , $J(\vec{w}, b)$

should decrease on every iteration.

feature engineering

- create a new feature

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

$$x_3 = x_1 \times x_2 \quad \leftarrow \text{new feature}$$

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

exp: price of a house

x_1 = frontage

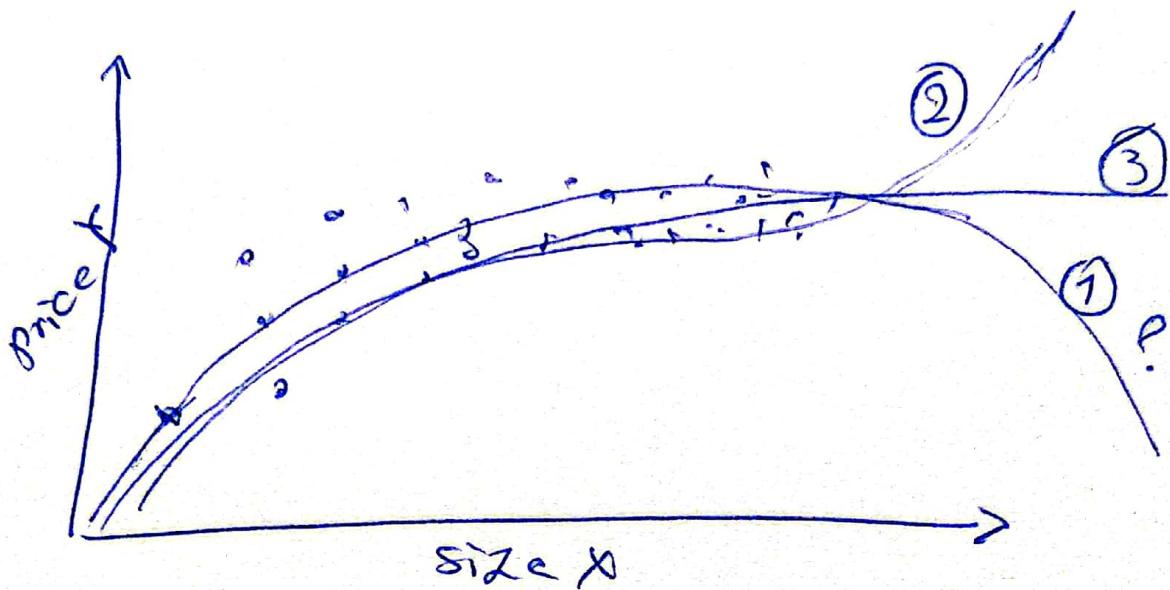
x_2 = depth

x_3 = area (frontage \times depth)

feature engineering, using intuition to design
new features, by transforming or combining -

original features.

Polynomial regression



① quadratic function: $f(\vec{w}, b) = w_1 x + w_2 x^2 + b$
 not very appropriate

② cubic function: $f(\vec{w}, b) = w_1 x + w_2 x^2 + w_3 x^3 + b$
 a better model

feature scaling really important because of
"powers"

③ square root of $f(\vec{w}, b)$: $f(x) = w_1 x + \cancel{w_2} + w_2 \sqrt{x} + b$

Classification

Binary classification: when thy target (y) can be only categorized in two classes.

Exps:

| | |
|-----------------|------------------|
| <u>1</u> | <u>0</u> |
| True | False |
| Yes | No. |

- is this email spam?

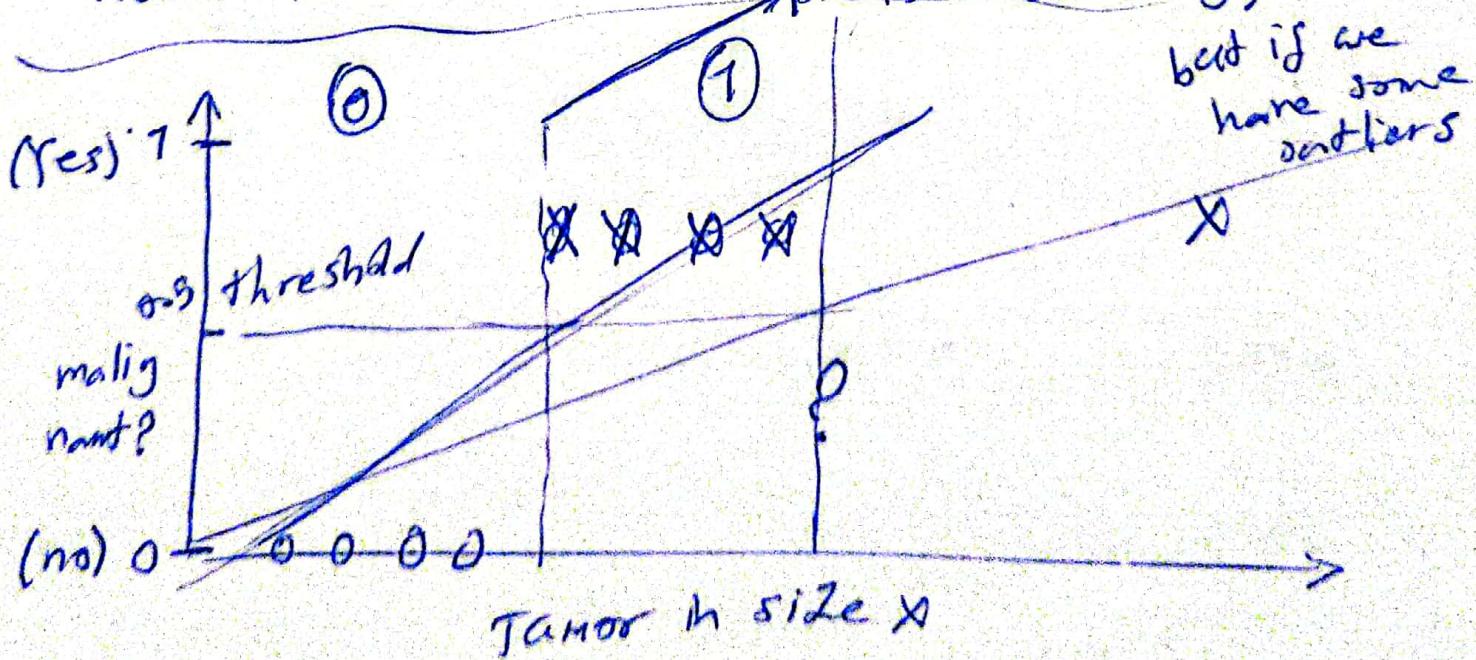
- is the transaction fraudulent? Yes No.

- is the tumor malignant? Yes No.

false/~~0~~/no \rightarrow negative class.

True/1/Yes \rightarrow positive class.

(decision boundary)

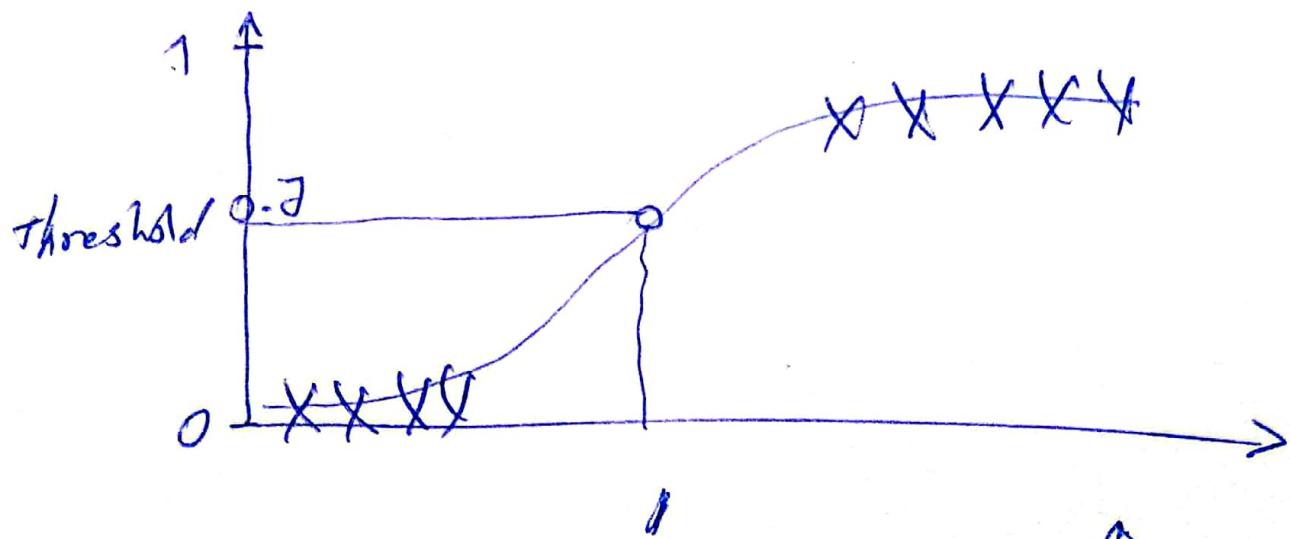


if $f_{w,b}(x) < 0.5 \rightarrow \hat{y} = 0$

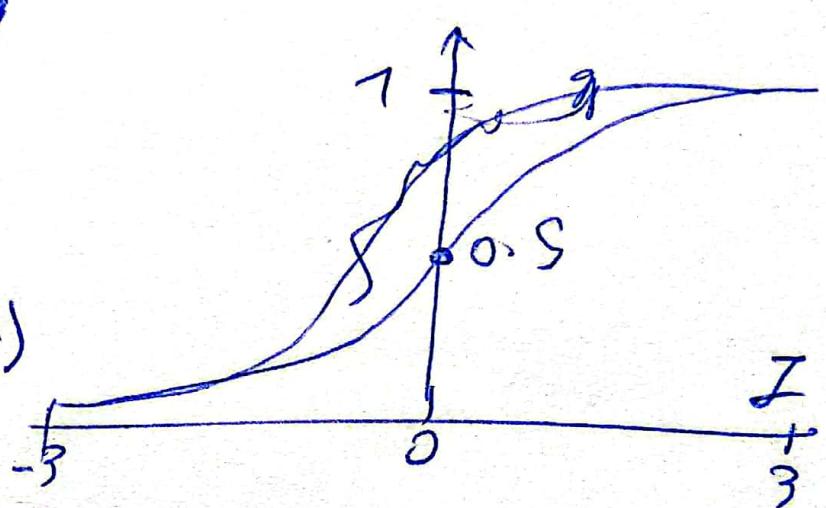
if $f_{w,b}(x) \geq 0.5 \rightarrow \hat{y} = 1$

logistic regression

(26)



sigmoid function
(logistic function)



output between 0 and 1

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

$$\text{when } z \gg \phi \quad \boxed{g(z) \approx 1}$$

$$\text{when } z = 0 \rightarrow g(z) = 0.5$$

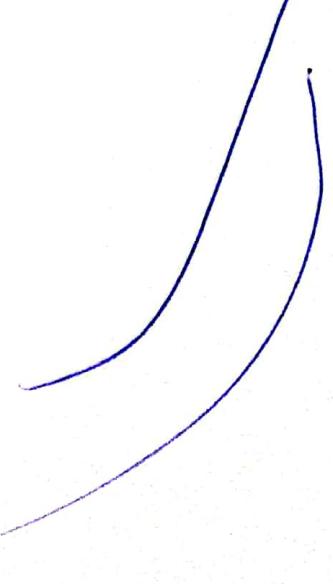
$$f(\vec{w}, b) \circ (\vec{x}) = \vec{w} \cdot \vec{x} + b$$

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

\downarrow

Z
 \downarrow

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



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

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

The last part of logistic regression is
the "probability" of that class is 1
 number between 0 and 1

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

probability that y is 1,
given input \vec{x} , Parameters \vec{w}, b

Decision boundary

logistic regression

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

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

$$\begin{matrix} \downarrow \\ z \\ \downarrow \\ 1 \end{matrix}$$

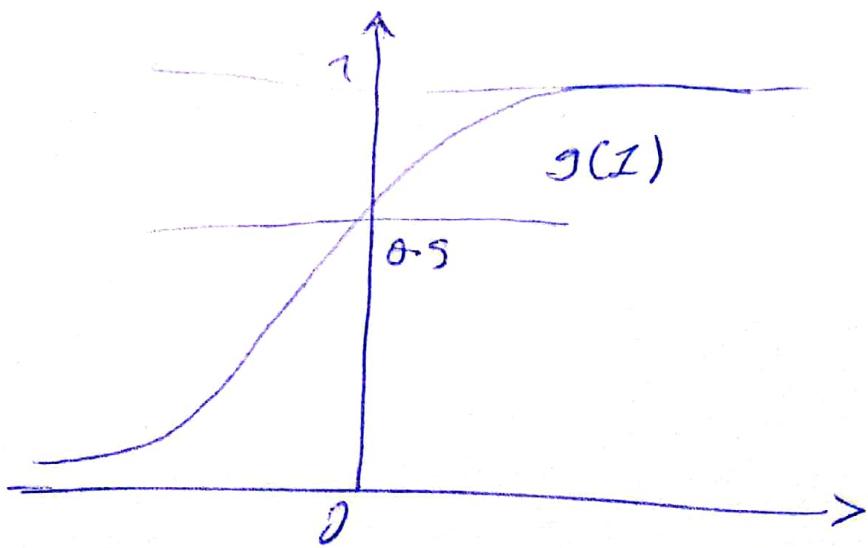
$$g(z) = \frac{1}{1+e^{-z}}$$

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

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

$$f_{\vec{w}, b}(\vec{x}) > 0.5$$

yes: $\hat{y} = 1$ no: $\hat{y} = 0$



when is $g(z)$ greater or equal 0.5?

$g(z)$ is greater than or equal 0.5 whenever z is greater or equal ~~less~~ to zero. it means when z is 0 or Positive.

when is $\vec{w} \cdot \vec{x} + b \geq 0.5$?

$$g(z) \geq 0.5$$

$$z \geq 0$$

$$\vec{w} \cdot \vec{x} + b \geq 0$$

$$\vec{w} \cdot \vec{x} + b \geq 0$$

$$\hat{y} = 1$$

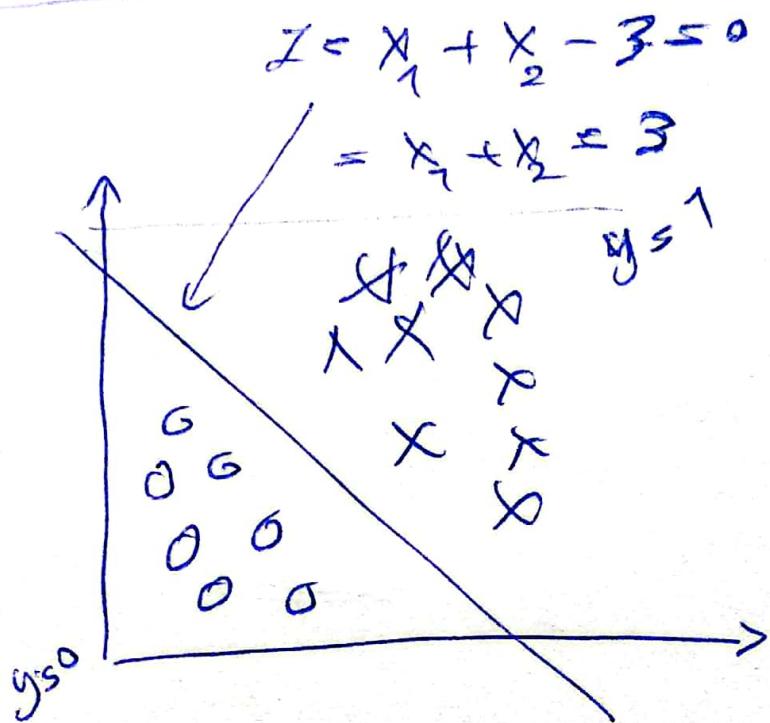
$$\vec{w} \cdot \vec{x} + b < 0$$

$$\hat{y} = 0$$

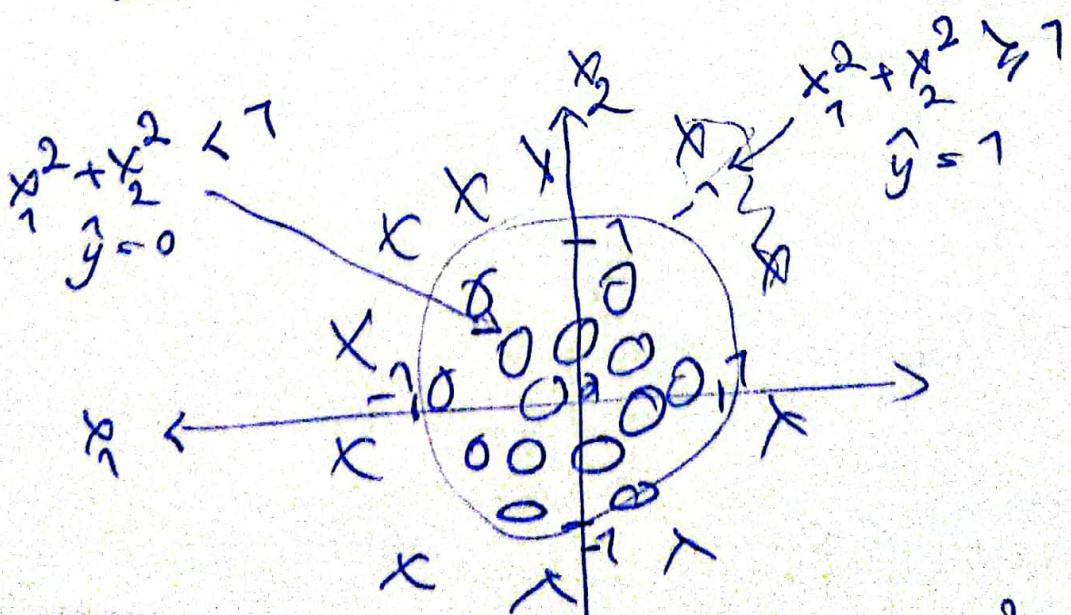
$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + b)$$

(30)

decision boundary: $z = \vec{w} \cdot \vec{x} + b = 0$



non linear decision boundary



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

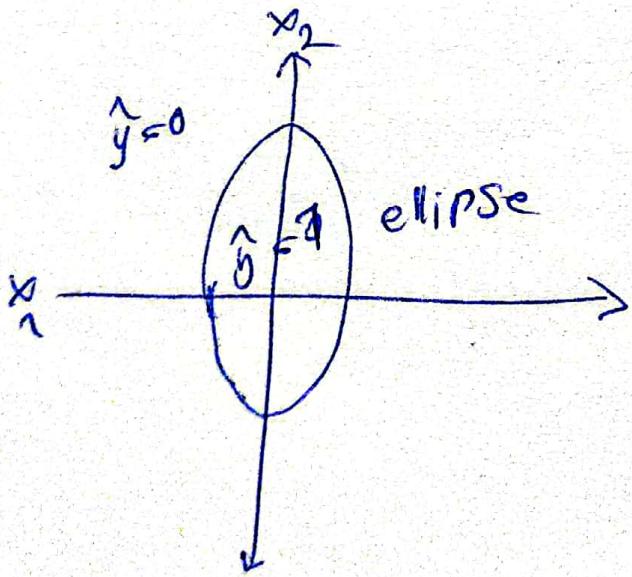
$$I = x_1^2 + x_2^2 - 1$$

decision boundary $\leq I = 0$

$$\leq x_1^2 + x_2^2 - 1 \leq 0$$

$$\leq x_1^2 + x_2^2 \leq 1$$

$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2 + \dots + b)$$



detected

- let's say we are creating a tumor algorithm. our algorithm will be used to flag potential tumors for future inspection by a specialist. what ~~should~~ value should we use for a threshold?

high esp 0.9

low esp 0.2

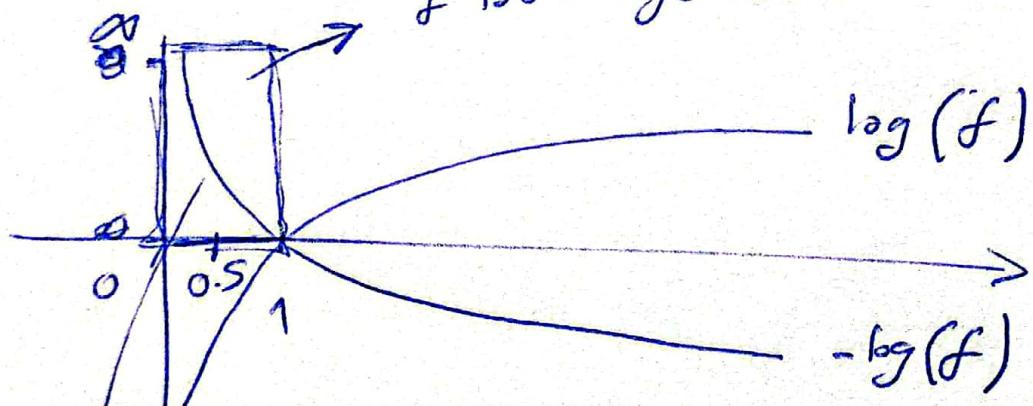
cost function for logistic regression

32

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

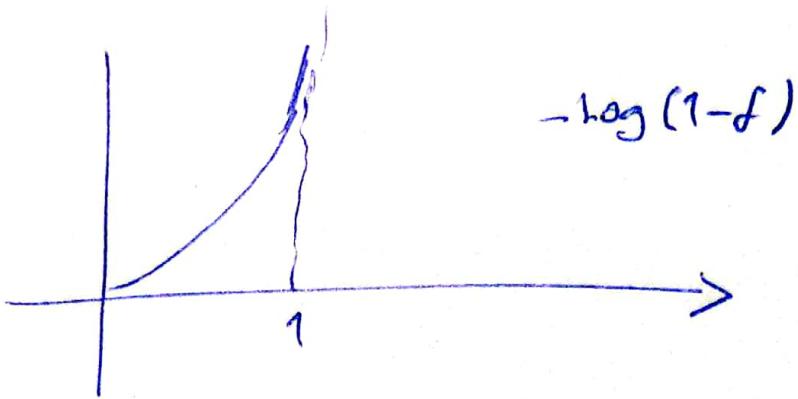
$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

f is always between 0 and 1. So



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

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



~~if $y = 0$~~

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

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

simplified loss function

~~$L(f_{\vec{w}, b}(\vec{x}^{(i)}))$~~

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

if $y^{(i)} = 1$

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

if $y^{(i)} = 0$

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

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

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

Training logistic regression

find \vec{w}, b (optimal)

given new \vec{x} , output

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

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

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

repeat {

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

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

}

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

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

"gradient descent for logistic regression"

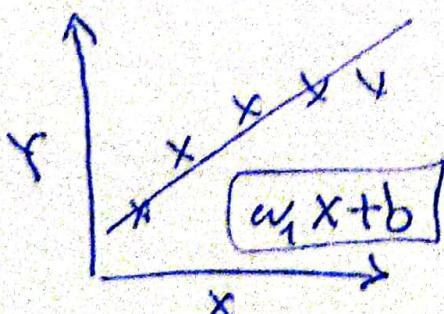
repeat {

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

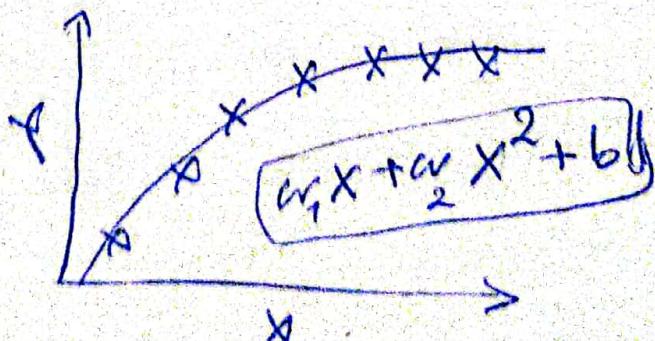
$$b = b - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \right]$$

} simultaneous updates

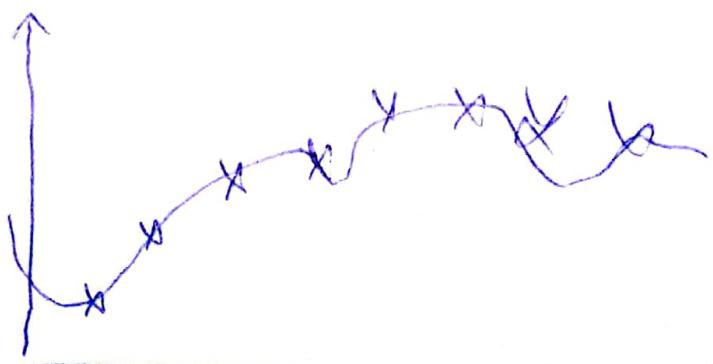
over fitting



(underfit
or
high bias)



generalization



$$y = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \alpha_3 x^3 + \dots + \alpha_9 x^9 + b$$

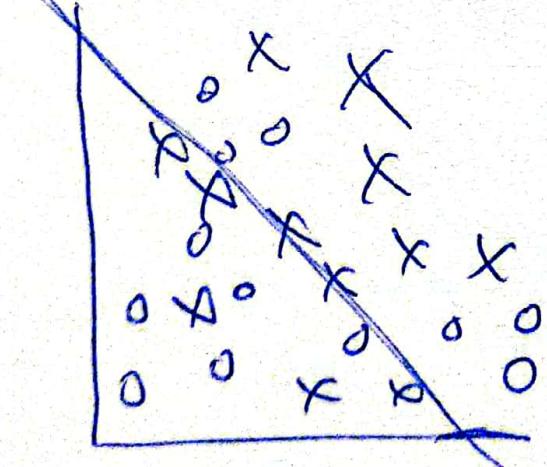
The model is overfitted
 $\text{cost} \approx 0$ but

The model has (overfitting problem)
 or

(high variance)
 in classification

decision boundary

decision boundary



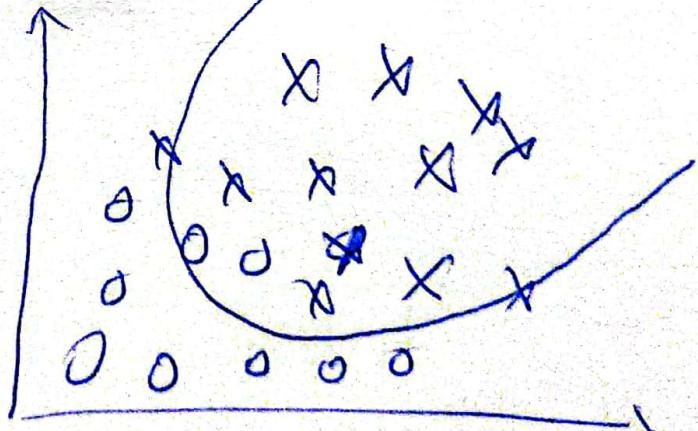
$$Z = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + b$$

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

underfit

or

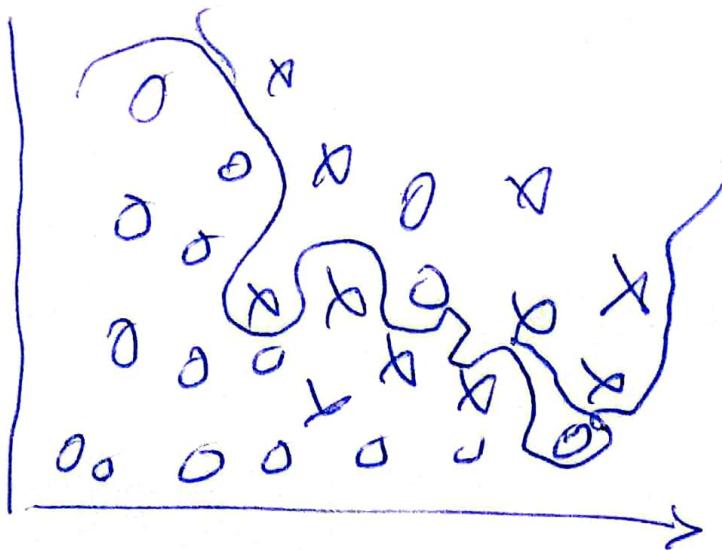
high bias



$$Z = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \alpha_4 x_4 + b$$

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

fit



$$Z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 + \dots$$

$\vdots \quad + b$

over fit

if the model ~~if~~ has overfitting

what we should do?

1. collect more training examples
2. Select features to include/exclude (feature selection)
3. regularization (reduce the impact of $f(x)$)

~~some~~ / some features

$$f(x) = 13x - 0.23x^2 + 0.000019x^3 - 0.000007x^6 + b$$

$\swarrow \searrow$

Small values for w_j or w_j^2

~~skip~~

models $w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

if we want to regularize w_3 and w_4 :

~~skip~~
~~skip~~

~~cost function~~ $J_{\vec{w}, b}(\vec{x}) = \frac{1}{2m} \sum_{i=1}^m \left(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)} \right)^2 + \underbrace{\frac{1}{1000} w_3^2 + \frac{1}{10000} w_4^2}_{\text{large number}}$

~~skip~~

if we have a model with a high number of features (e.g., 100 features) we can

regularize all of them.

So:

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

sometimes:

$$+ \frac{1}{2m} b^2$$

n = number of features.

m = number of data points.

λ (lambda) = regularization parameter λ small number

if lambda is 0 model is underfit.

if lambda is too large model is overfit.

regularized linear regression

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

Gradient descent \mathcal{L}

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

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

\downarrow

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

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

if we regularize b

then $+ \frac{1}{m} w_b$

Gradient descent with regularization

Q10

$$\text{GD} \left\{ \begin{array}{l} w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m [f_{\bar{w}, b}(\vec{x}^{(i)}) - y^{(i)}] x_j^{(i)} \right] + \frac{1}{m} w_j \\ b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\bar{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \end{array} \right.$$

simultan easily update.

Regularized logistic regression

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

GD
repeat

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

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

}

$$\frac{\partial}{\partial w_j} j(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \left[(\hat{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \right] + \frac{1}{m} \omega_j$$

$$\frac{\partial}{\partial b} j(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (\hat{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$