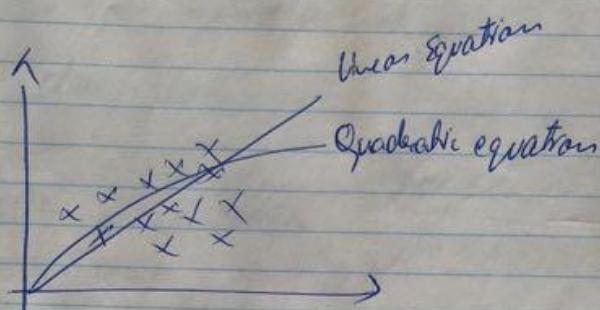


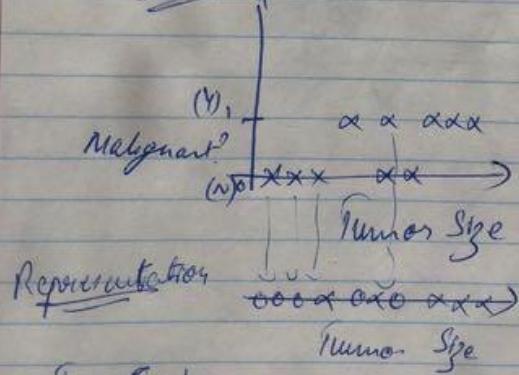
ML

Regression predict continuous valued output

Supervised



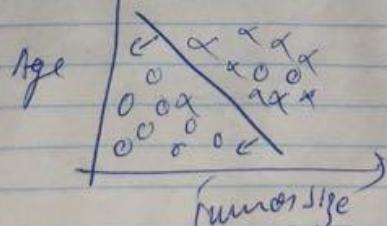
Ex One feature



Classification

Discrete values of
(0 or 1)
0, 1, 2, 3
Type 1
Type 2
Type 3
Type 4

Two Features

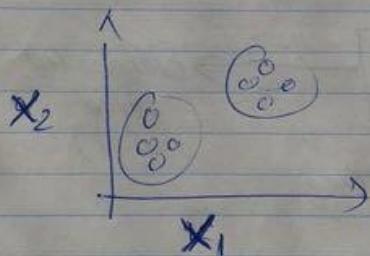


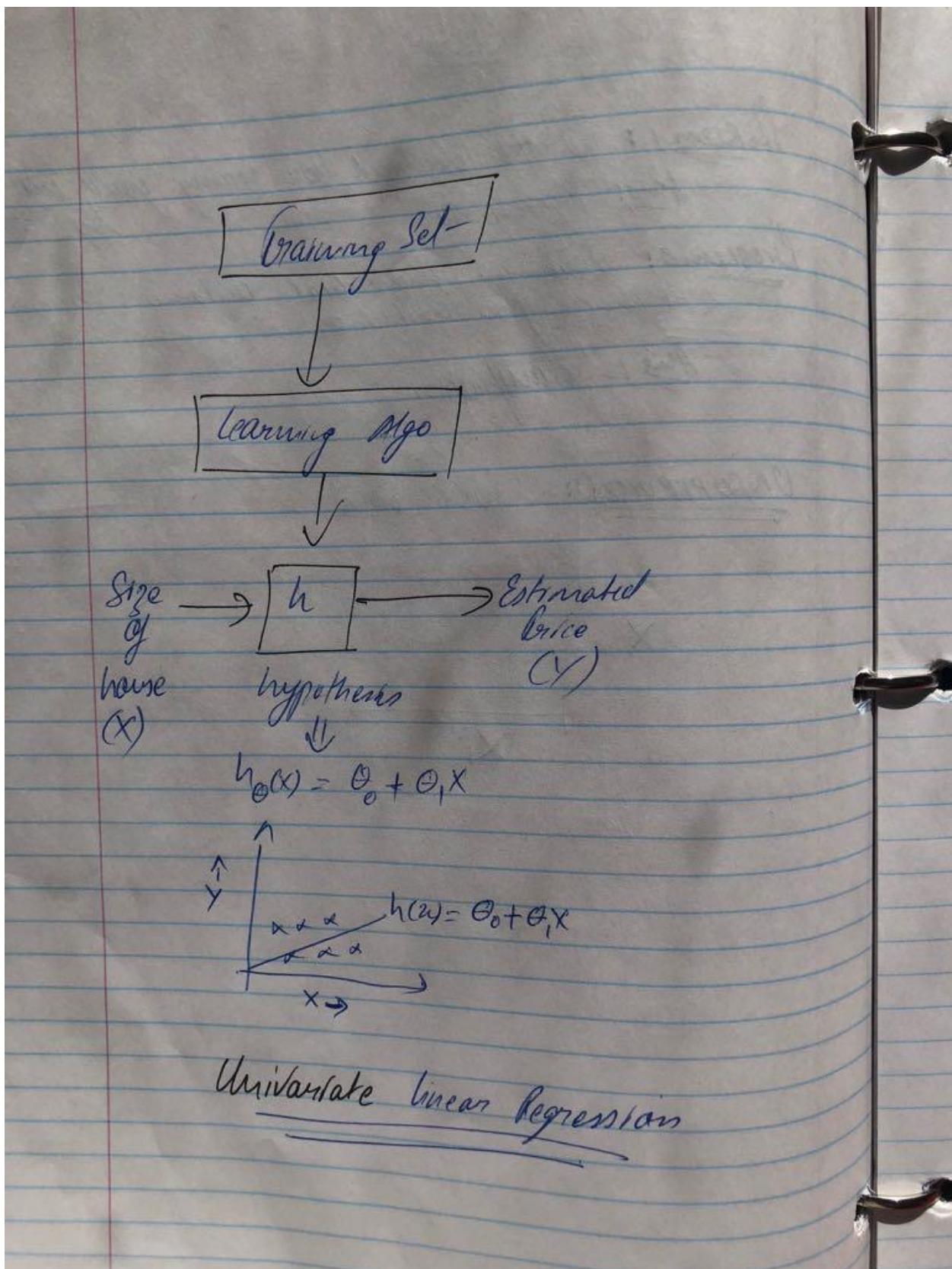
Problem 1: Item identical, how many you'll sell over ~~12~~ months
Ans: Regression

Problem 2: Manage individual customer accounts,
& decide if its hacked or not?

Ans: Classification

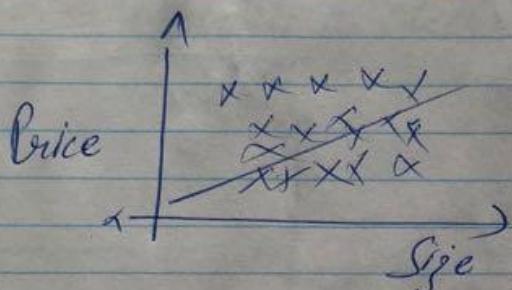
UNSUPERVISED: JUST DATA





LINEAR REGRESSION

Housing Prices [Regression Problem]



| <u>Size (m²) (X)</u> | <u>Price (\$) (Y)</u> | |
|---------------------------------|------------------------|----------|
| 2104 | 460 | |
| 1234 | 52 | |
| { } | { } | } m = 47 |

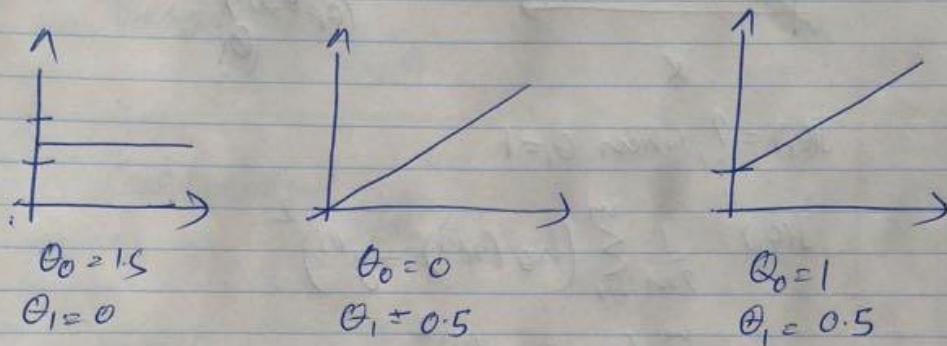
(x, y) - one example

$(x^{(i)}, y^{(i)})$ - i^{th} training example | $(x^{(1)}, y^{(1)}) = (2104, 460)$

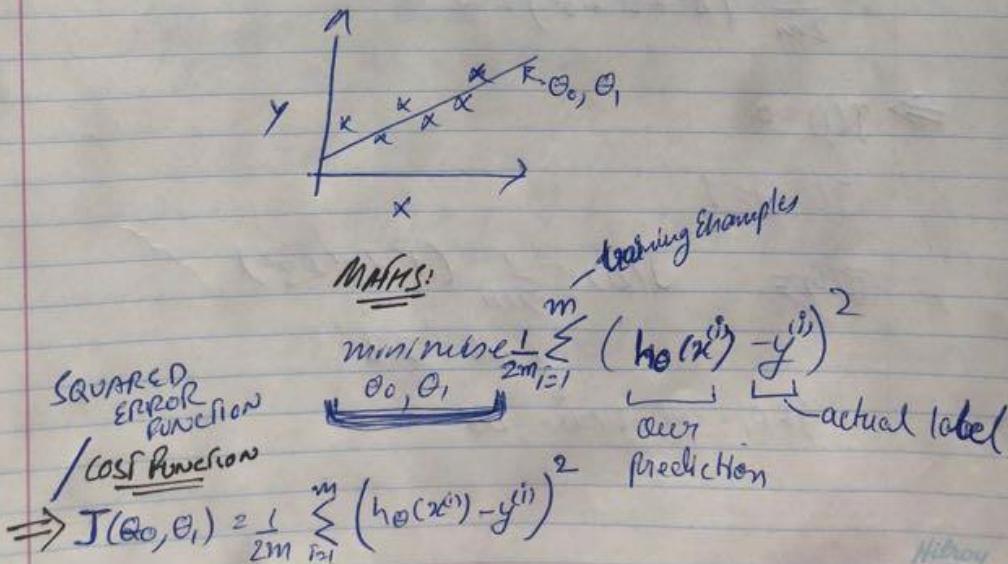
hypothesis:

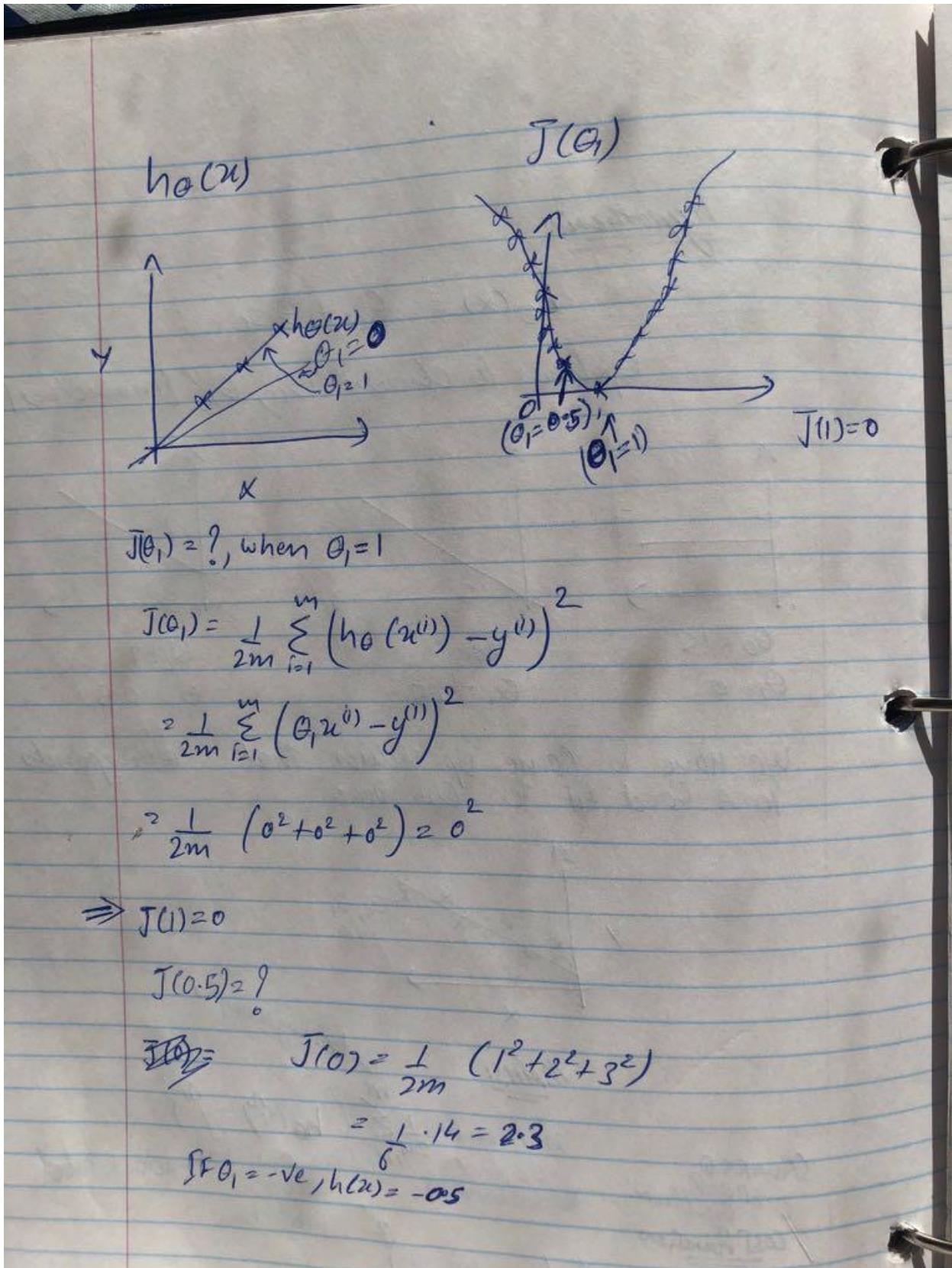
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Q: how to choose θ_0, θ_1 [parameters]



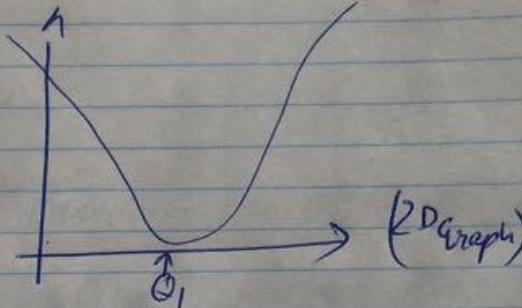
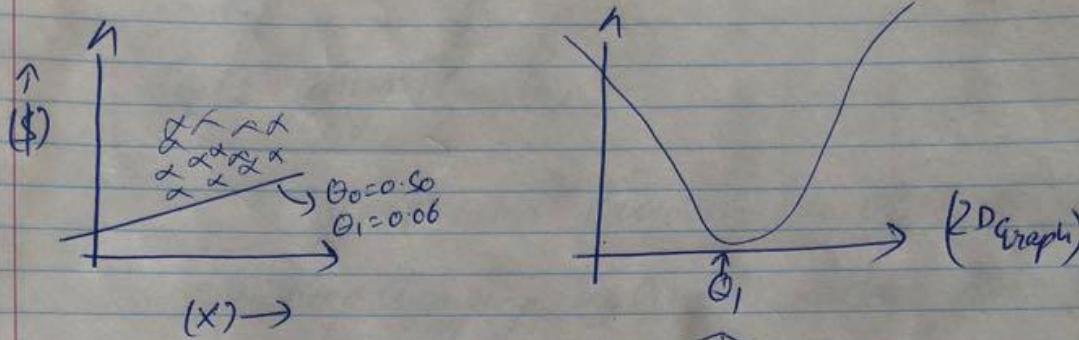
We have to come up values that corresponds to a good fit to our data.



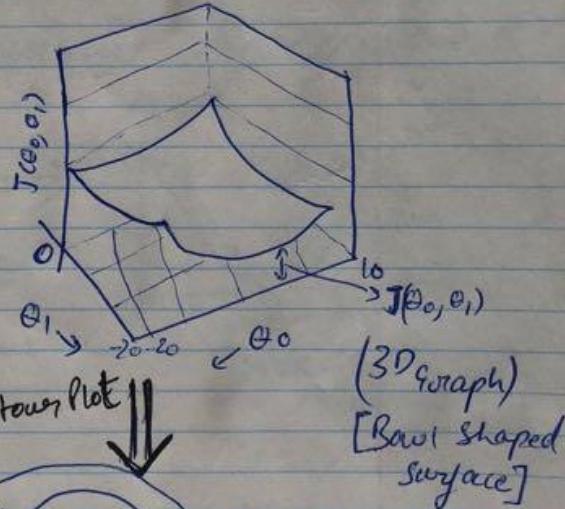


Two Parameters : $J(\theta_0, \theta_1)$

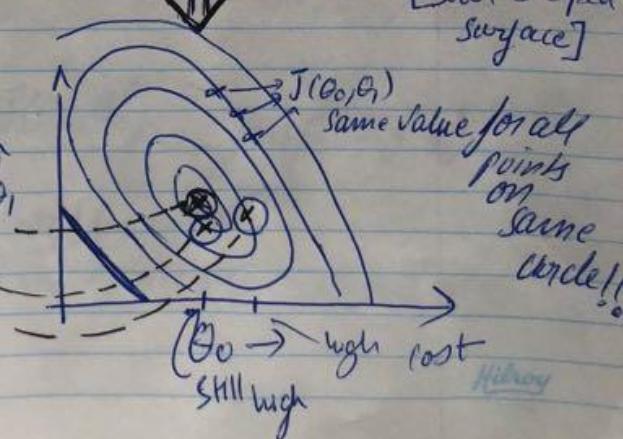
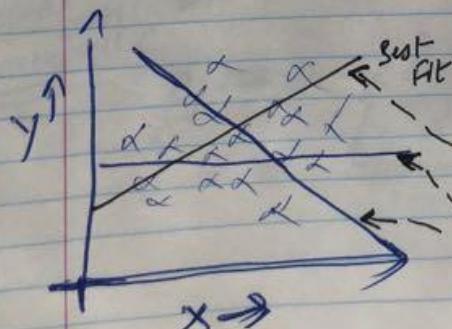
$h_{\theta}(x)$



$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



(3D graph)
[Bowl shaped surface]



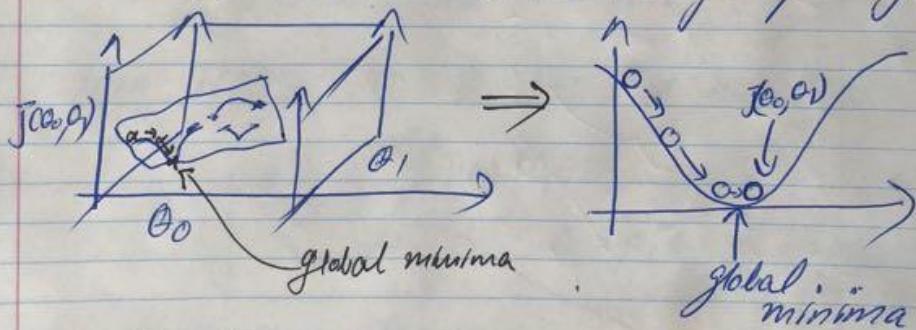
GRADIENT DESCENT

A way to minimize $J(\theta)$ / cost function

$$J(\theta_0, \theta_1)$$

$$\text{Goal: } \min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

- ① Start with Random (θ_0, θ_1)
- ② Keep changing (θ_0, θ_1) to reduce $J(\theta_0, \theta_1)$
- ③ Find minimum by repeating ②.



Algorithm:

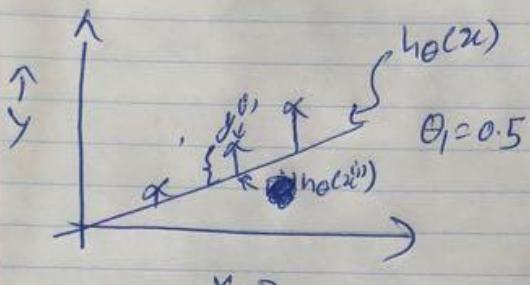
~~Repeat until convergence~~ {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{for } (j=0 \& j=1)$$

} learning rate [Any number]
[how big of a step to take??]

Hilary

$h_{\theta}(x)$
(for fixed θ_1 , this function of x)



$$J(0.5) = ?$$

$$= \frac{1}{2m} \left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right]$$

$$= \frac{1}{2 \times 3} (3.5) = \frac{3.5}{6} = 0.58$$

$$\Rightarrow J(0.5) = 0.58$$

So, our goal minimize $J(\theta_1)$

We want to fit the line (straight line) to our data

And, as we saw $(\theta_1=1)$ gave the line that best fit the data!

Cost function

Hypothesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Parameters:

$$\theta_0, \theta_1$$

Cost Function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Goal:

$$\underset{\theta_0, \theta_1}{\text{minimise}} J(\theta_0, \theta_1)$$

Explanation [One parameter]
 $\Rightarrow h_{\theta}(x) = \theta_1 x$

$$\Rightarrow \theta_1$$

$$\Rightarrow J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\Rightarrow \underset{\theta_1}{\text{minimise}} J(\theta_1)$$

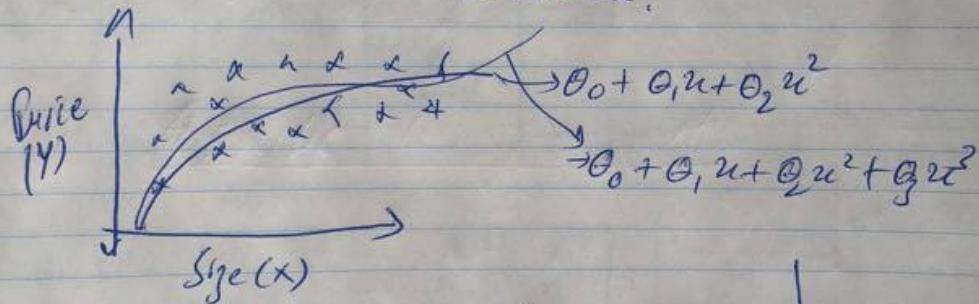
POLYNOMIAL REGRESSION

House Price Prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{u_1} + \theta_2 \times \underbrace{\text{depth}}_{u_2}$$

New feature \Rightarrow Area
 $x = \text{frontage} * \text{depth}$

$$h_{\theta} = \theta_0 + \theta_1 x \quad \text{land area.}$$



$$h_{\theta}(u) = \theta_0 + \theta_1 u_1 + \theta_2 u_2 + \theta_3 u_3$$

$$= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

$$u_1 = \text{size}$$

$$u_2 = (\text{size})^2$$

$$u_3 = (\text{size})^3$$

Size : 1-1000

Size² : 1-100,000

Size³ : 1-10⁹

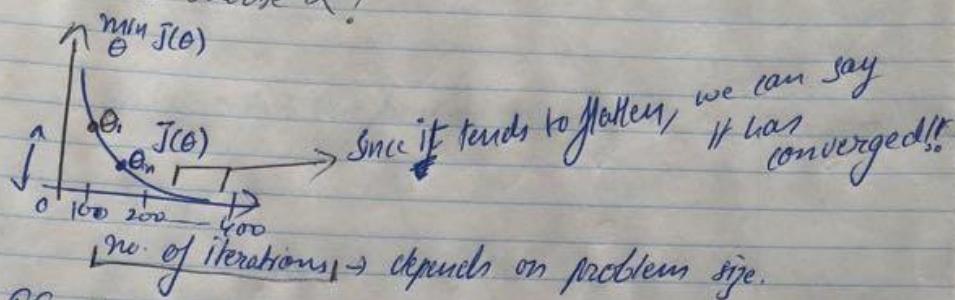
or we choose something like:

$$h_{\theta}(u) = \theta_0 + \theta_1 (\text{size}) + \theta_2 \sqrt{\text{size}}$$

Hilroy

LEARNING RATE $[\alpha]$

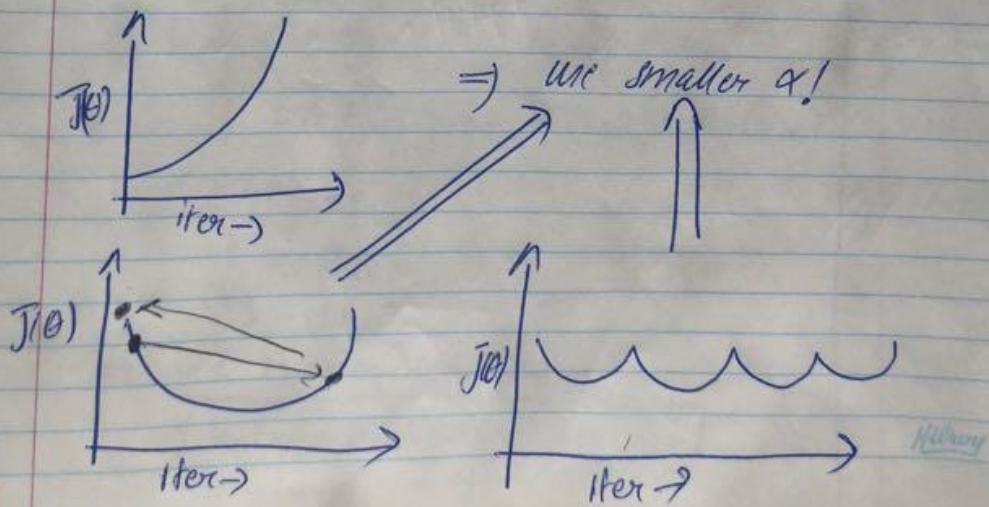
How to choose α ?



If our algo works properly then our $J(\theta)$ will always keep on decreasing.

Automatic Convergence Test

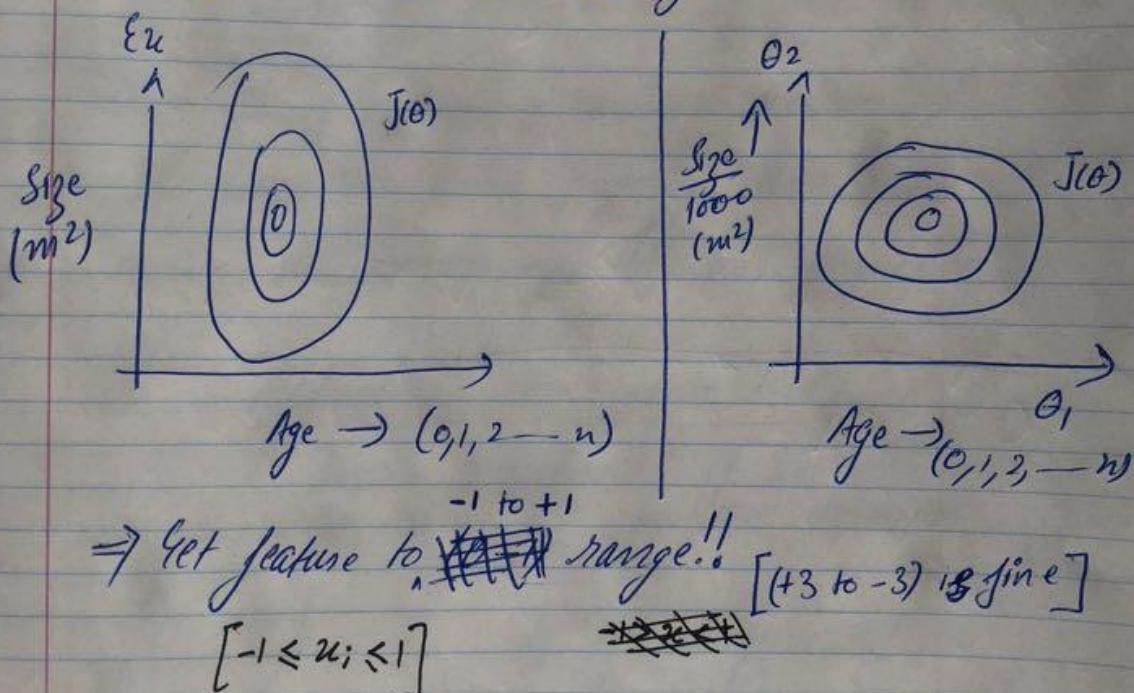
Declare convergence if $|J(\theta)|$ decreases by less than 10^{-3} in one iteration.



Feature Scaling

Make sure feature are on same scale so that gradient descent converges quickly.

If the scale is too big then the contours graph will be very elongated causing it harder to converge.



One way := normalization!! [Don't apply $\theta_0 = 1$ with it]

$$x_j \leftarrow \frac{x_j - \mu_j}{\sigma_j} \quad \text{mean of } x_j$$

$$\text{range}_{\min}^{(x_j)} - \text{range}_{\max}^{(x_j)}$$

Hilary

Gradient descent:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

} [Simultaneous update for every $j=0 \rightarrow n$]

Algorithm

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(Simultaneous update θ_j for
 $j=0, 1, \dots, n$)

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 x_0^{(i)}$$

$$\theta_1 = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 x_1^{(i)}$$

$$\vdots$$

$$\theta_n = \theta_n - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 x_n^{(i)}$$

Hilary

$$\begin{aligned}
 & \left[\theta_0, \theta_1, \theta_2, \dots, \theta_n \right] \quad \left[\begin{array}{c} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right] \\
 & \text{---} \\
 & (\text{n+1}) \times 1 \quad x \\
 & \theta^T \\
 & \Rightarrow h_{\theta}(x) = \theta^T x
 \end{aligned}$$

Multivariate Linear Regression

Hypothesis: $h_{\theta}(x) = \theta^T x \rightarrow x_0 = 1$ [Assumption]

$$h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Parameters: $\theta = [\theta_0, \theta_1, \theta_2, \dots, \theta_n]$ n+1-dimensional vector

Cost Function:

$$\begin{aligned}
 J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) &= \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\
 &\downarrow \\
 J(\theta)
 \end{aligned}$$

| <u>'n' Parameters</u> | | | | | |
|-----------------------|-------|-------|---------|-------|-----|
| x_1 | x_2 | x_3 | \dots | x_n | y |
| 1 | 5 | 28 | | | |
| 2 | 6 | 98 | | | |
| 3 | 7 | 76 | | | |
| 4 | | | | | |

MATRIX X

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}^T$$

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$x^{(i)} = \text{input features of } i^{\text{th}} \text{ example}$

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

$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Eg: $\theta_0 + 0.1x_1 + 0.01x_2 + \dots + 2x_4 - 2x_5$

let $x_0 = 1 \Rightarrow [x_0^{(i)} = 1]$

$\Rightarrow x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$

$\Rightarrow h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = \theta^T x$

{ Repeat Until Convergence

$$\theta_0 := \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

OR

Repeat Until Convergence

$$\left. \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \end{array} \right] \text{update } \theta_0, \theta_1$$

$$\left. \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \end{array} \right] \text{simultaneously}$$

{ Repeat Until Convergence

$$\left. \begin{array}{l} \text{temp 0} = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \text{temp 1} = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \end{array} \right]$$

$$\theta_0 := \text{temp 0}$$

$$\theta_1 := \text{temp 1}$$

{ The cost func. for linear regression will always be bowl shaped because it's a convex function.

Linear Regression Model

$h_{\theta}(u) = \theta_0 + \theta_1 u$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(u^{(i)}) - y^{(i)}]^2$$

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

Gradient Descent Algo

repeat until convergence

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} (J(\theta_0, \theta_1))$$

$$\text{for } j=1 \text{ to } j=0 \text{ do}$$

OR

repeat until convergence

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (J(\theta_0, \theta_1))$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} (J(\theta_0, \theta_1))$$

$$\theta_0 := \text{temp}_0$$

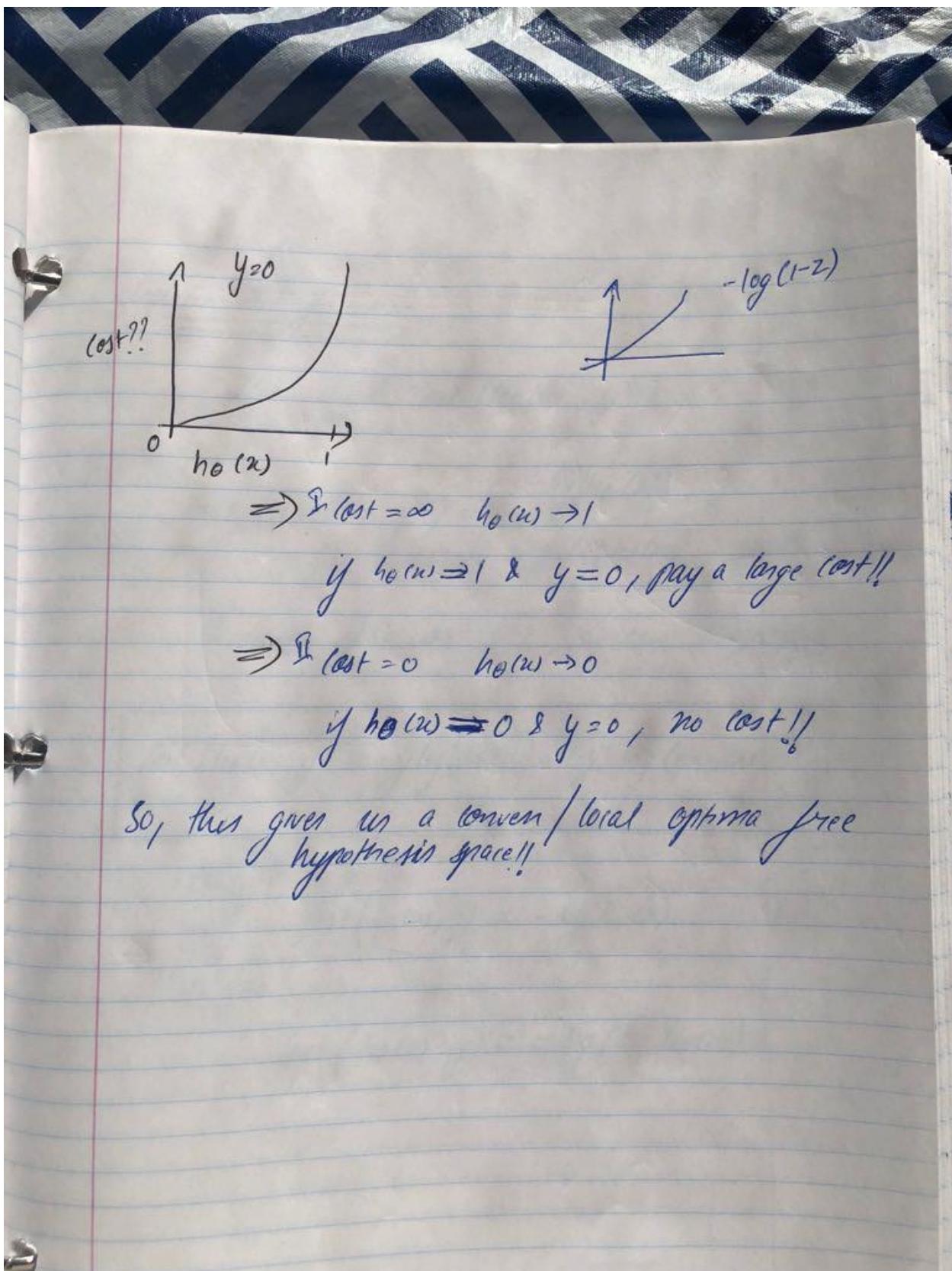
$$\theta_1 := \text{temp}_1$$

FROM

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(u^{(i)}) - y^{(i)}]^2$$

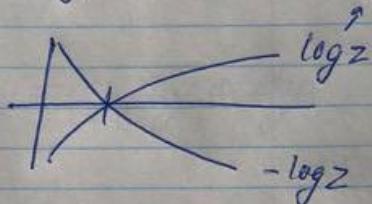
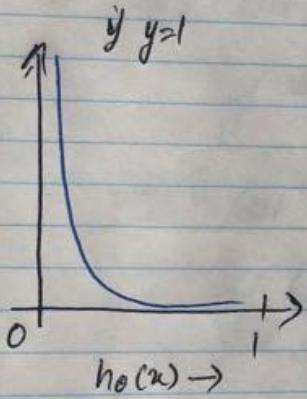
$$\Rightarrow \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

Hilroy



Logistic Regression Loss Function

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log h_{\theta}(x) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$



Cost = 0 if $y=1, h_{\theta}(x)=1$

But as

$$h_{\theta}(x) \rightarrow 0 \\ \text{Cost} \rightarrow \infty$$

⇒ captures intuition that if $h_{\theta}(x)=0$

$$\text{predict } P(y=1/x; \theta) = 0$$

but $y=1$ for that x , then

we'll penalize learning algorithm by a
very large cost.

Hibway

Example

Training Set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

m examples: $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad x_0 = 1, y \in \{0, 1\}$

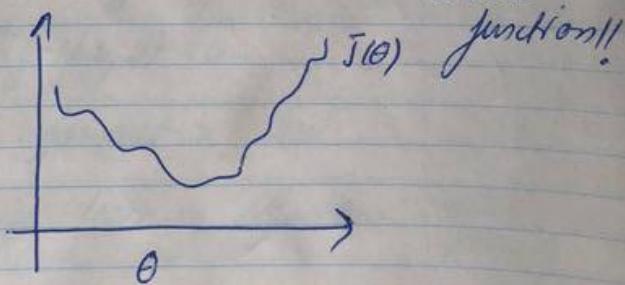
Hypothesis

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose θ ??

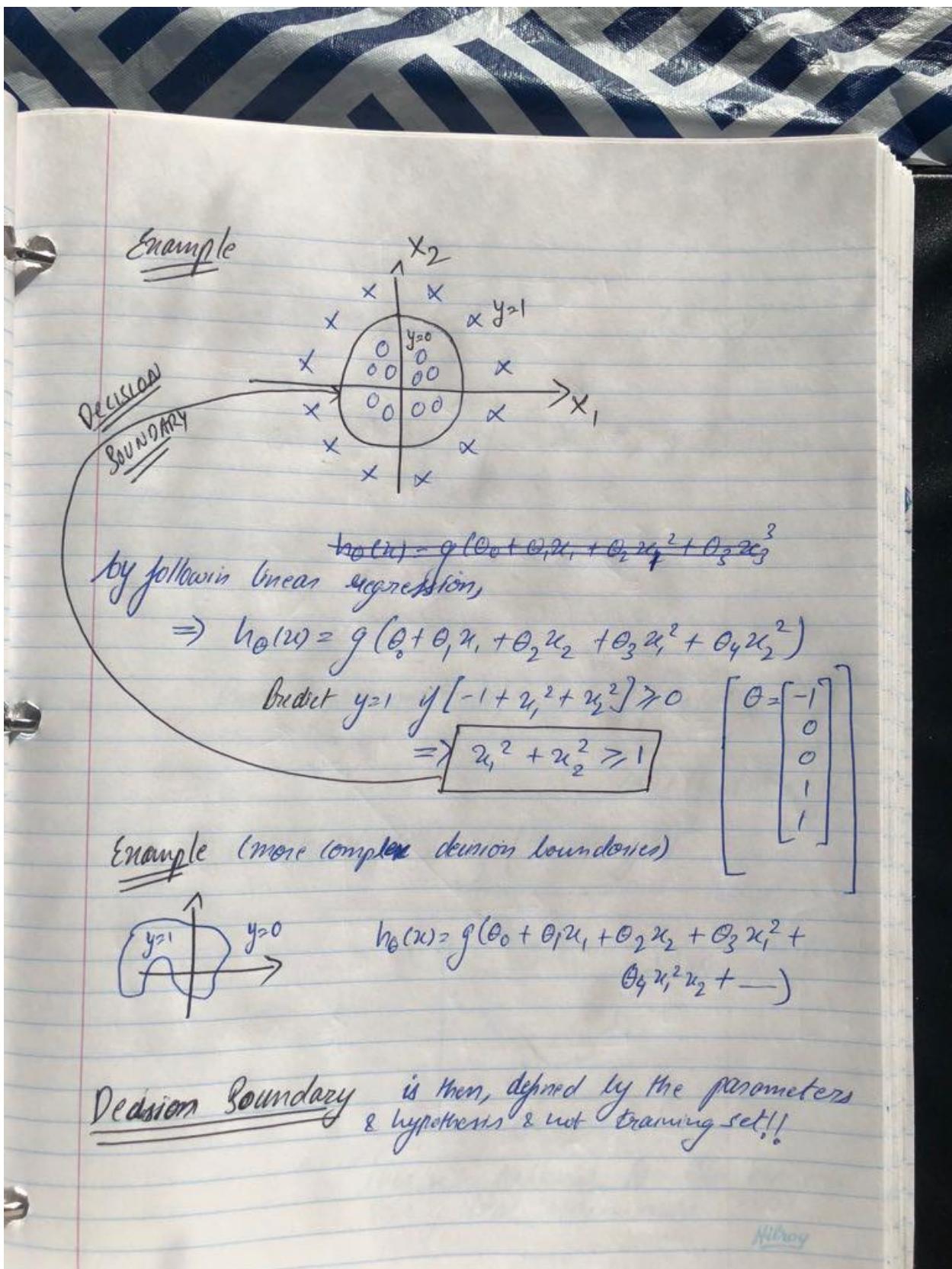
$$\text{Linear Regression} := J(\theta) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2}_{\text{cost}} (\text{cost})(h_{\theta}(x^{(i)}), y)$$

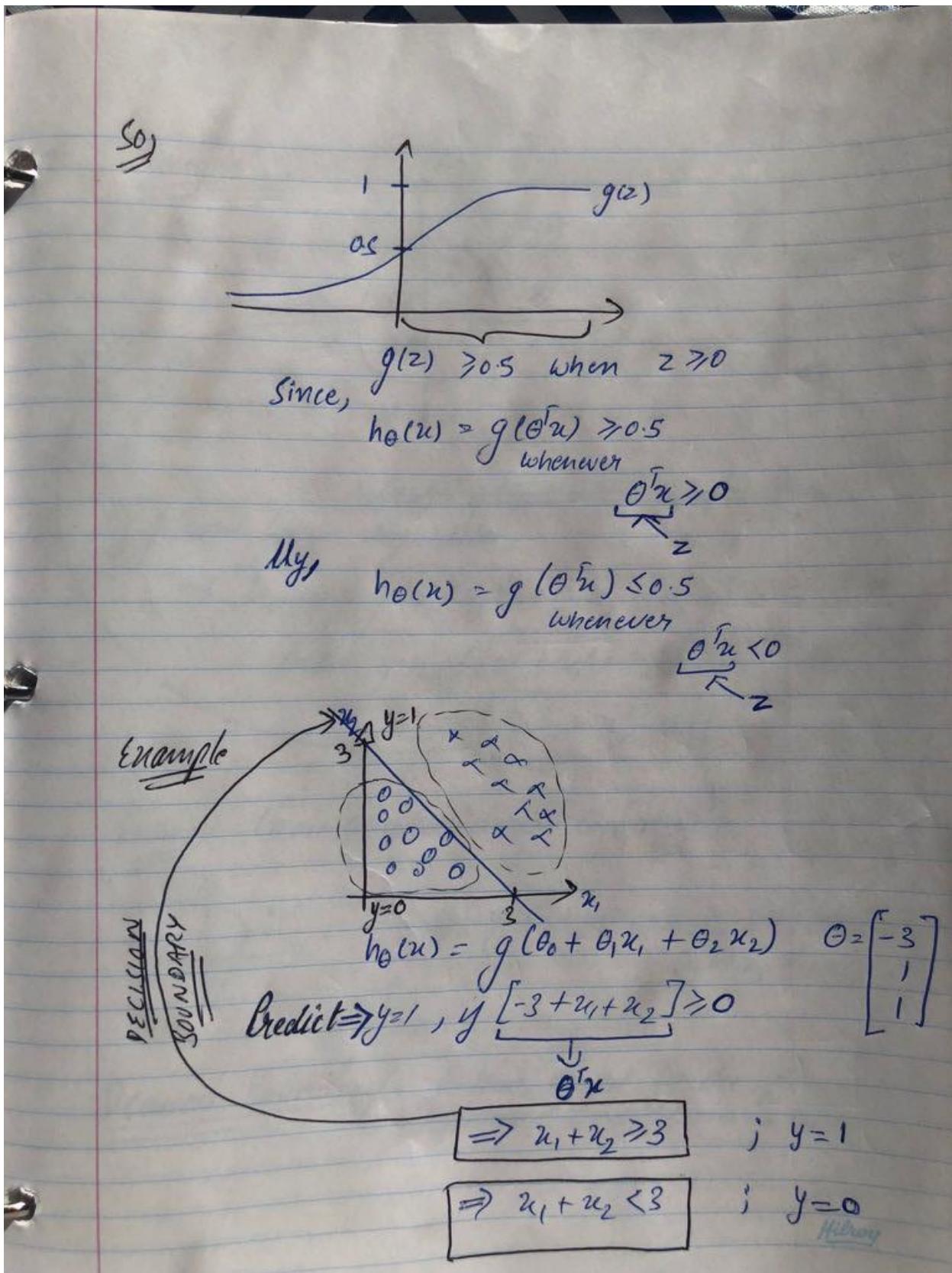
$$\Rightarrow \text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



This is a problem! As there are too many local minimas.

History





Interpretation

$h_{\theta}(u)$ = estimated probability that $y=1$ on input x .

Example

$$u = \begin{bmatrix} x_0 \\ u_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{turnsize} \end{bmatrix}$$

we get,
 $h_{\theta}(u) = 0.7 ; y=1$

That means ' u ' has 0.7 probability
of being class 1!

OR,

$$h_{\theta}(u) = P(y=1 | u; \theta)$$

"Probability that $y=1$, given u
parameterized by ' θ '."

So,

$$\Rightarrow P(y=0 | u; \theta) + P(y=1 | u; \theta) = 1$$

$$\Rightarrow \boxed{P(y=0 | u; \theta) = 1 - P(y=1 | u; \theta)}$$

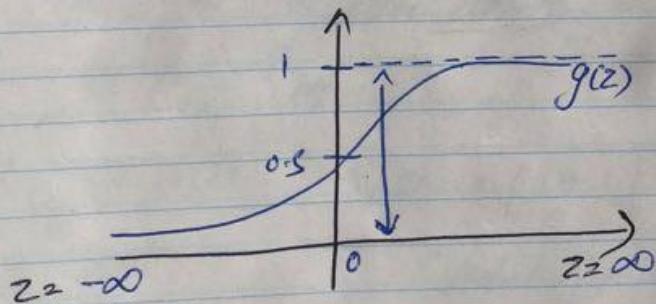
Want $0 \leq h_{\theta}(u) \leq 1$

We can make

$$h_{\theta}(u) = \Theta^T u$$

Hypothesis $\quad \quad \quad$ to

Sigmoid function $\quad \quad \quad$
$$\left. \begin{array}{l} h_{\theta}(u) = g(\Theta^T u) \\ g(z) = \frac{1}{1+e^{-z}} \end{array} \right\} h_{\theta}(u) = \frac{1}{1+e^{-\Theta^T u}}$$

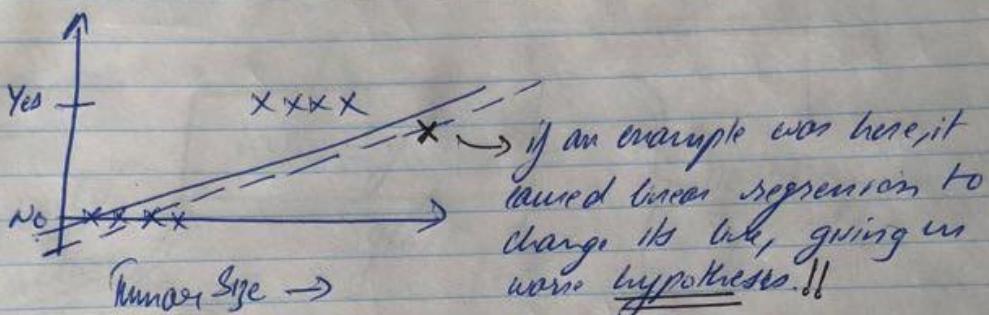


So, when $|z|$ reaches $+\infty/-\infty$, $g(z)$ reaches $1/0$

Hence $h_{\theta}(u) = g(z)$ will always be less than 1 & more than 0.

LOGISTIC REGRESSION (SIGMOID functions)

Example1



Threshold classifier output $h_\theta(x)$ at 0.5:

if, $h_\theta(x) \geq 0.5$ then $y=1$

else, $h_\theta(x) < 0.5$ then $y=0$

Example2 $y=0$ or $y=1$

$h_\theta(x)$ can be >1 or <0 !!

Logistic Regression:

$$0 \leq h_\theta(x) \leq 1$$

| <u>m Training Examples, n features</u> | |
|---|--|
| <u>Gradient Descent</u> | <u>Normal Equation</u> |
| Need to choose α | No Need to choose α |
| Needs many iterations | Don't Need to iterate |
| Works well even when $'n'$ is large $\uparrow n = 10^6 \rightarrow$ | Need to compute $(X^T X)^{-1}$ Slow if ' n ' is very large. n_1 few n_2 one n_3 few |
| <u>A.K.A</u> \Rightarrow Gradient Descent when $n > 10^6$ \Rightarrow Normal Equation when $n < 10^6$ | |

NORMAL EQUATION

$$\Theta = (X^T X)^{-1} X^T y$$

Set $A = X^T X$

$$(X^T X)^{-1} = A^{-1}$$

m examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$; n features

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

(design
Matrix)

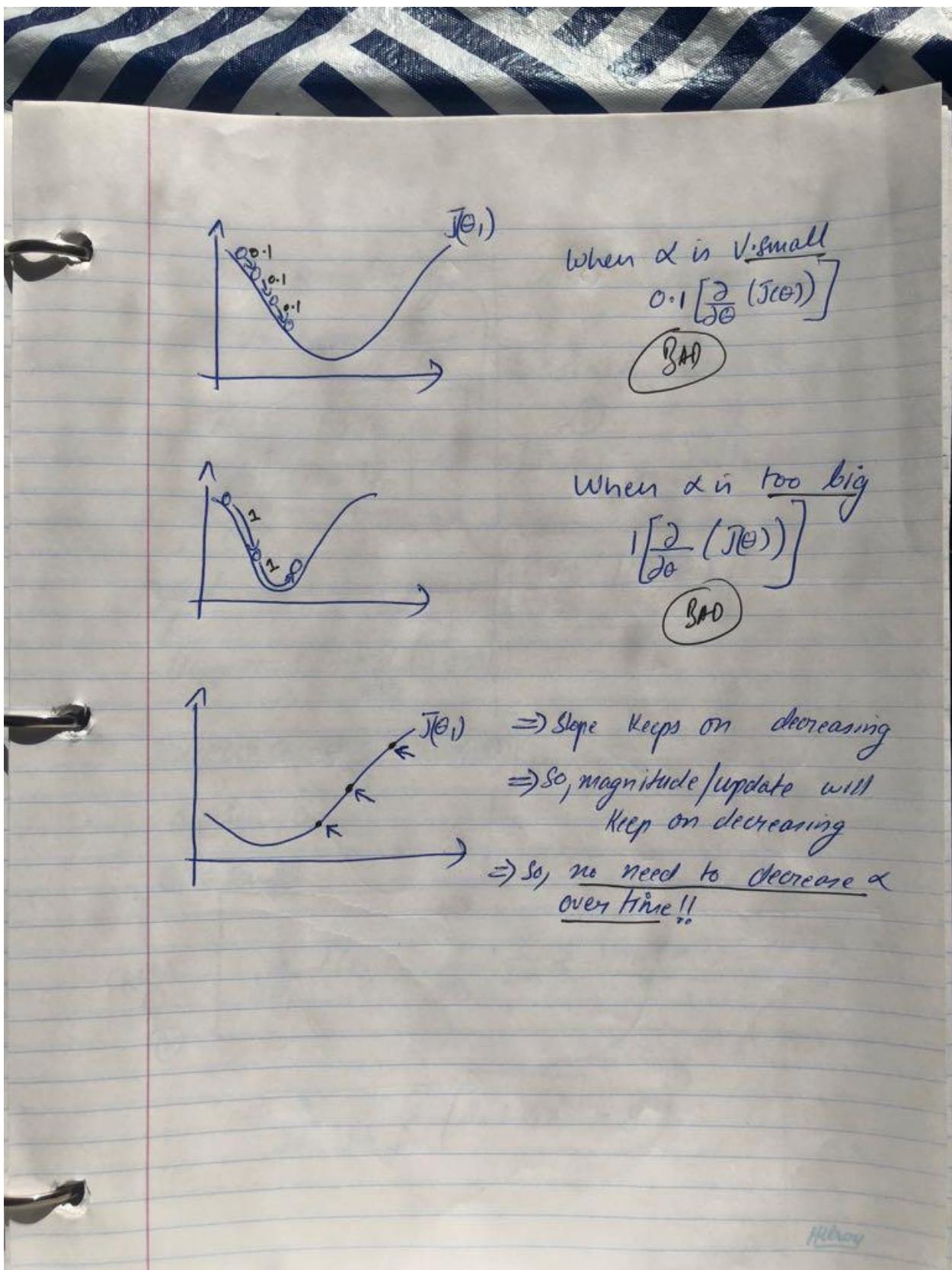
$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

Eg $y = x^{(i)} = \begin{bmatrix} 1 \\ x_i^{(i)} \end{bmatrix}$

$$x = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_2^{(1)} \\ \vdots & \vdots \\ 1 & x_m^{(1)} \end{bmatrix}_{m \times 2}$$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\Rightarrow \boxed{\Theta = (X^T X)^{-1} X^T y}$$



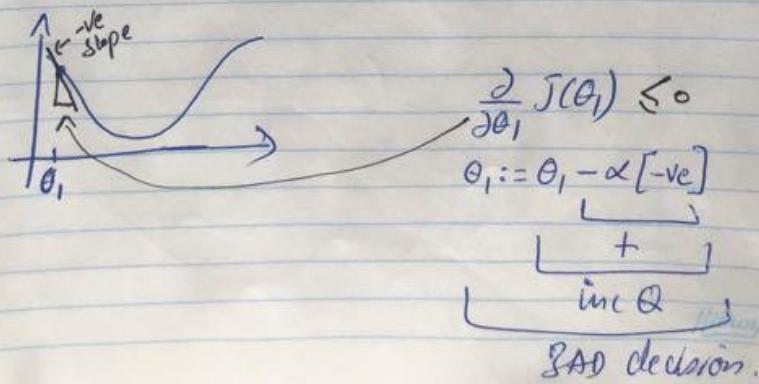
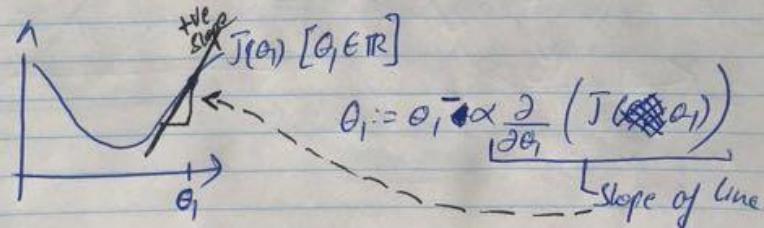
Correct Implementation

$$\left\{ \begin{array}{l} \text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \end{array} \right.$$

$\theta_0 := \text{temp}_0$] → update only after
 $\theta_1 := \text{temp}_1$ computing $J(\theta_0, \theta_1)$
 } Repeat Until convergence

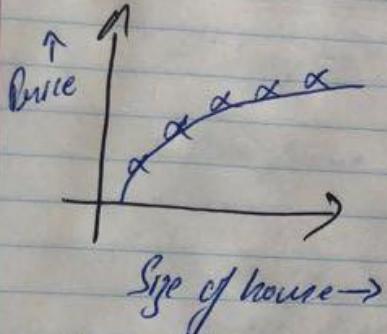
Simultaneous Update
or
Gradient Descent

$$\min_{\theta_1} J(\theta_1) \quad \theta_1 \in \mathbb{R} \quad [\text{all real numbers}]$$

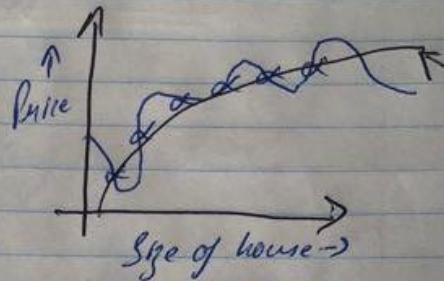


REGULARISATION

Intuition



$$\theta_0 + \theta_1 u + \theta_2 u^2$$



$$\theta_0 + \theta_1 u + \theta_2 u^2 + \theta_3 u^3 + \theta_4 u^4$$

↑ if they $\theta_3, \theta_4 \approx 0$

Suppose we penelise & make θ_3, θ_4 really small

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \left[\begin{array}{l} \text{SQUARED ERROR} \\ \text{COST FUNCTION} \end{array} \right]$$

$$+ 1000 \theta_3^2 + 1000 \theta_4^2$$

To make it small, θ_3, θ_4 has to be small!

minimising So, as we will end up, $\theta_3 \approx 0, \theta_4 \approx 0$

A quadratic function that is good!
By adding the 2 new terms!!

Addressing Overfitting

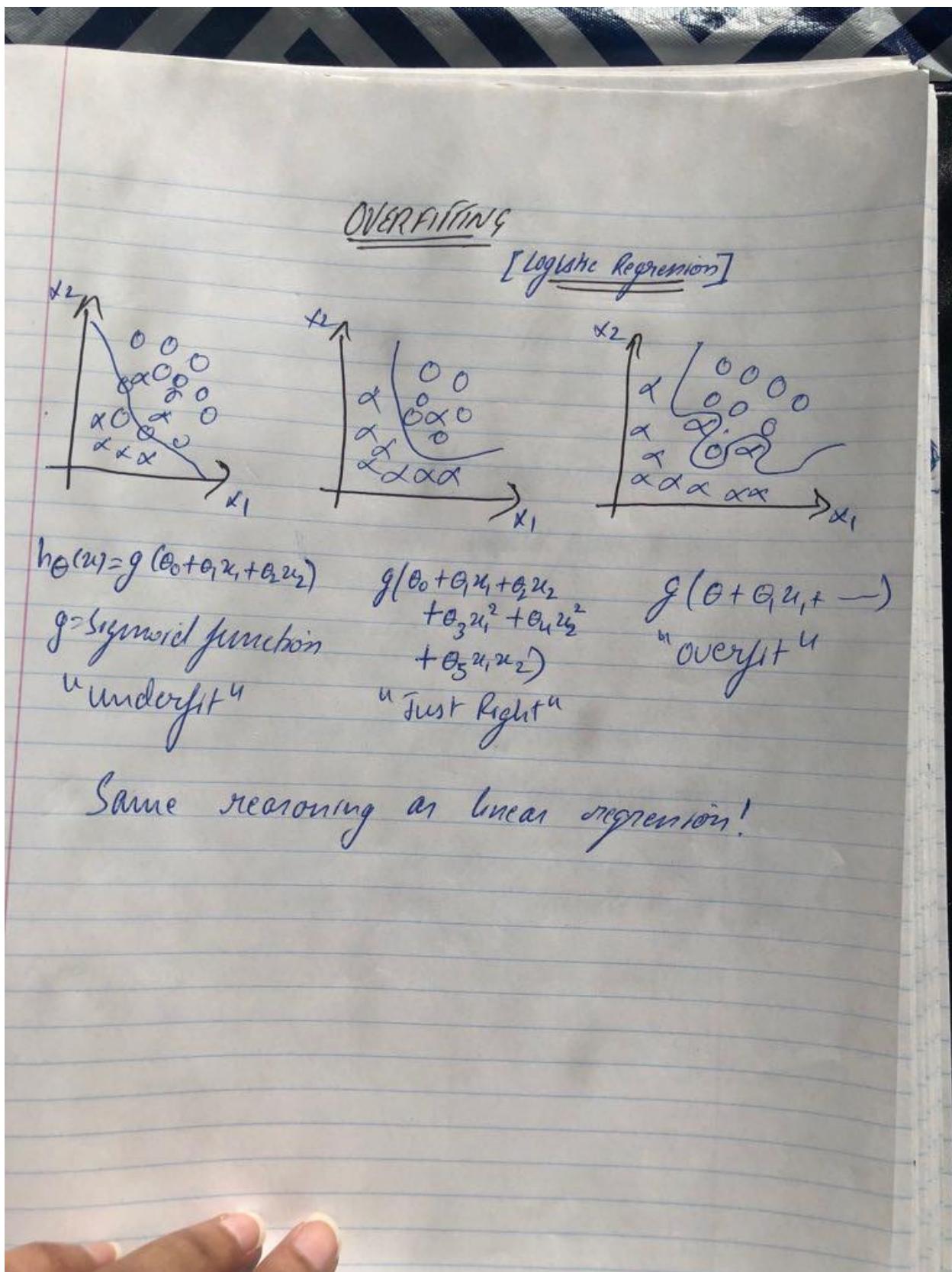
Reduce Number of Features

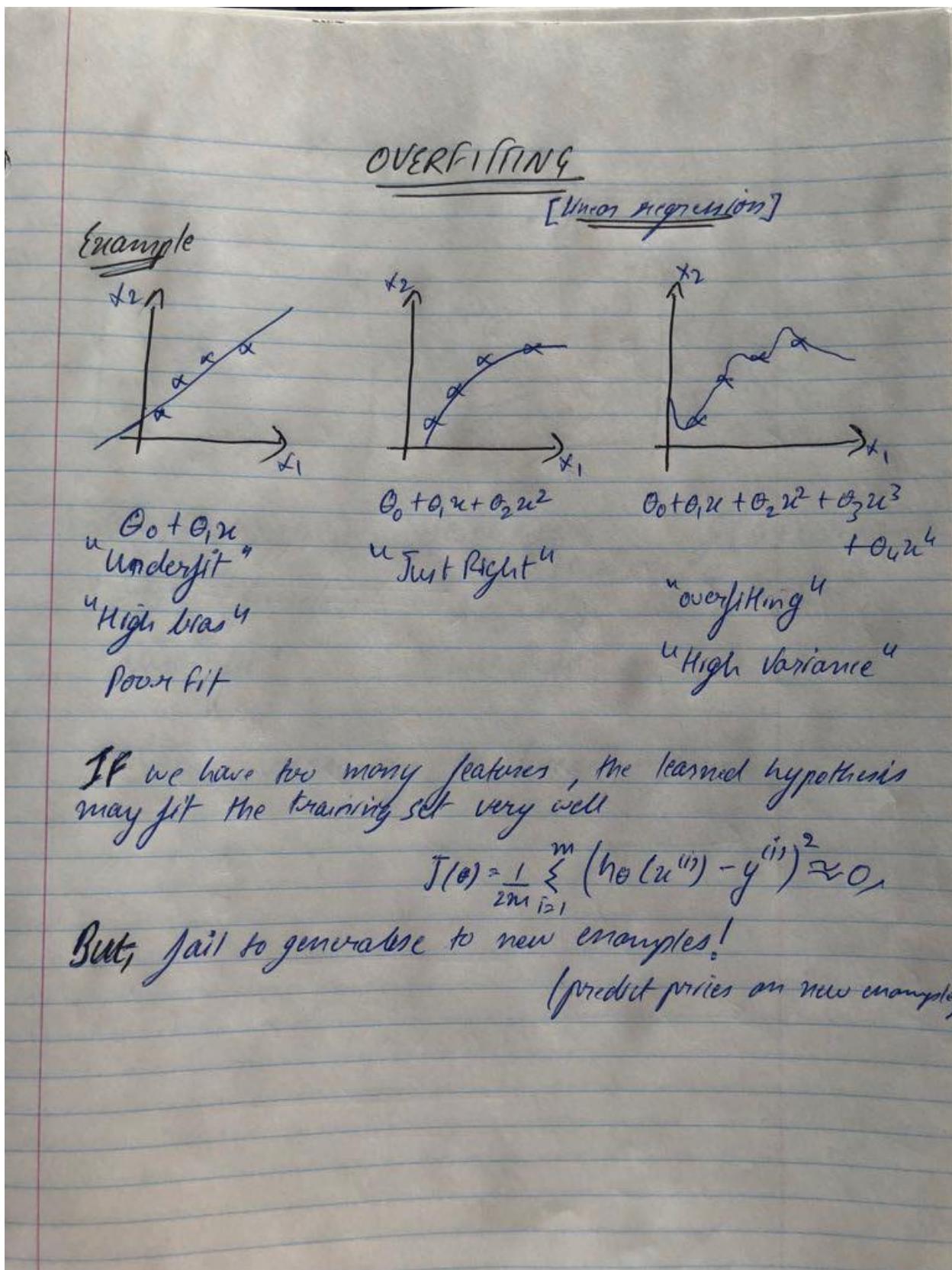
- manually
- model selection algorithms!

What if we throw useful feature away?

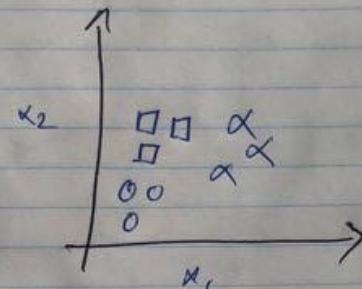
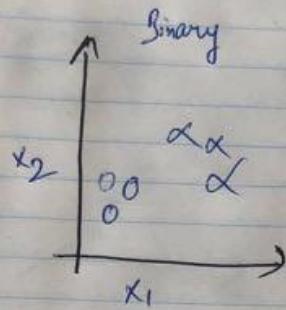
REGULARIZATION

- Keep all feature, but reduce magnitude / values of parameter θ_j
- Works well when we have a lot of features each of which contributes a bit to predicting y .





MULTICLASS CLASSIFICATION



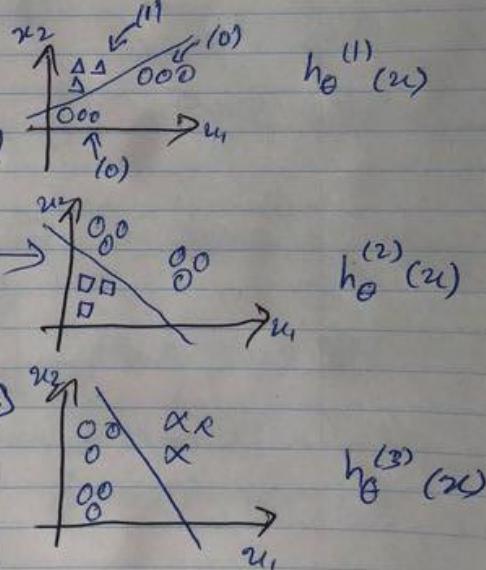
One vs All Classification

Class 1: Δ
 Class 2: \square
 Class 3: \times

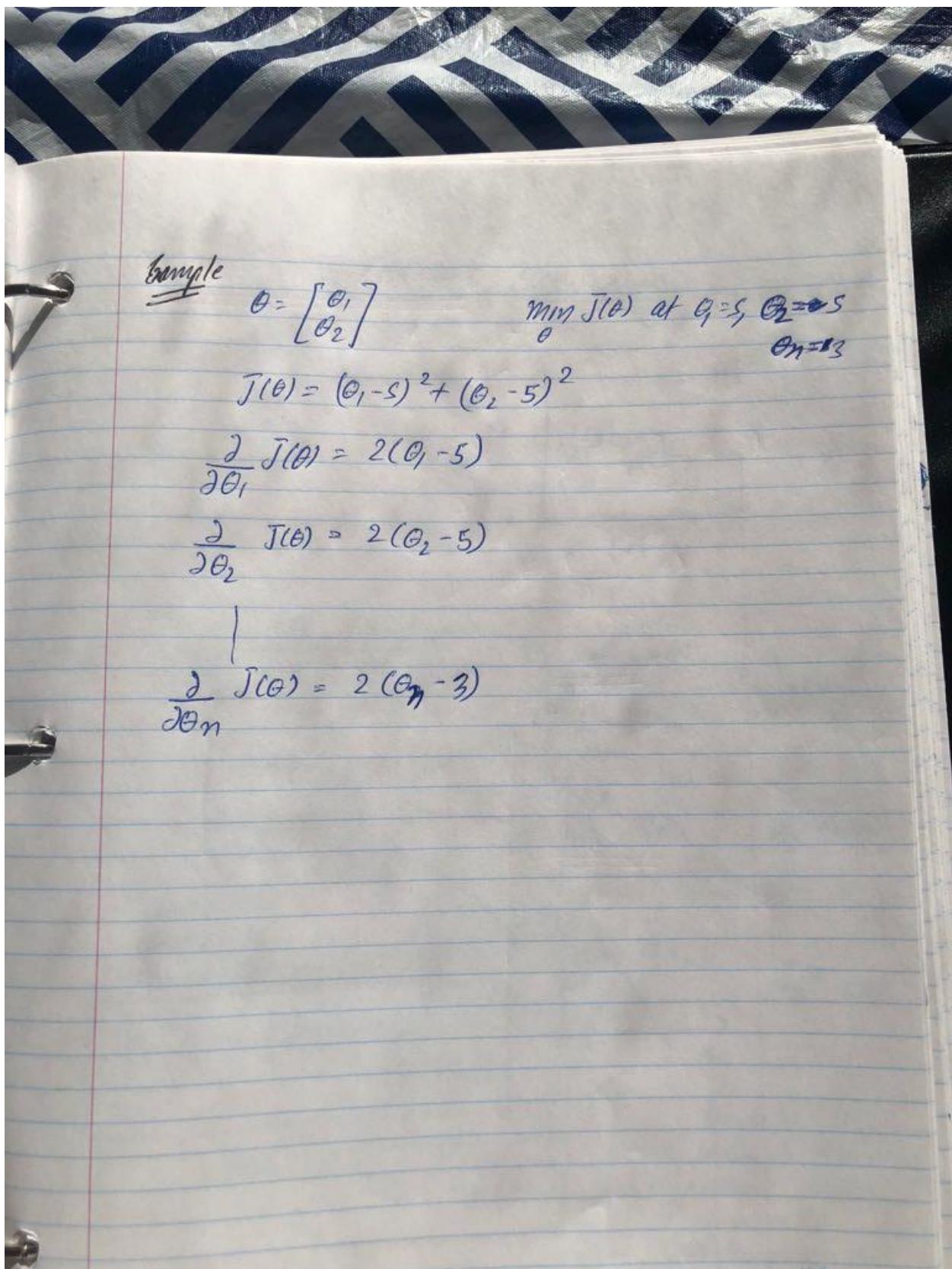
$$h_{\theta}^{(i)}(w) = P(y=i | w; \theta) \quad (i=1, 2, 3)$$

then,

$$\max_i h_{\theta}^{(i)}(w)$$



predict!
 We choose the hypothesis that
 is most enthusiastic!!
Hilary



Optimization Algorithm

Gradient descent:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

we found out two things

$$J(\theta) \text{ & } \frac{\partial J(\theta)}{\partial \theta_j}$$

Other Algorithms

- Conjugate gradient
- BFGS
- L-BFGS

Advantages

- No need to manually pick α !
- often faster than gradient descent

Disadvantages

- More complex

Algorithm

Repeat {

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_\theta(u^{(i)}) - y^{(i)}) u_j^{(i)}$$

}

Simultaneously update θ_j

But this is same as linear regression!!

In logistic regression

$$h_\theta(u) = \frac{1}{1 + e^{-\theta^T u}}$$

In linear regression:

$$h_\theta(u) = \theta^T u$$

The difference is that the definition of hypothesis
is different!!Tip: Feature Scaling can be used to make it
converge faster!

Cost function of logistic Regression

$$\Rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)})) \right]$$

To fit parameters θ :

$$\min_{\theta} J(\theta) \rightarrow \text{How? Gradient Descent!!}$$

Prediction:

$$h_\theta(w) = \frac{1}{1+e^{-\theta^T w}}$$

Gradient Descent

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

logistic cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(w), y^{(i)})$$

$$\text{cost}(h_\theta(w), y^{(i)}) = \begin{cases} -\log(h_\theta(w)) & ; y=1 \\ -\log(1-h_\theta(w)) & ; y=0 \end{cases}$$

Because $y=0$ or $y=1$, we can simplify!!

$$\text{cost}(h_\theta(w), y) = -y \log(h_\theta(w)) - (1-y) \log(1-h_\theta(w))$$

If $y=1$

$$\text{cost}(h_\theta(w), y) = -\log h_\theta(w)$$

If $y=0$

$$\text{cost}(h_\theta(w), y) = -\log(1-h_\theta(w))$$

} matches with above!!

So, Anyways.

Discrete Random Variable

- A discrete variable, that can have be present as a possible value in a instance of a random number!!

Continuous Random Variable

- Random Variable, in a continuous set of numbers, can have multiple possible values.

pdf (discrete Random Variable)

$$P(u) = P(X=u)$$

How many time 'u' occurs in all the possible values of 'X'.

i.e. frequency of 'u' occurring in 'X'.

Hilroy

Stochastic Model

- used to capture randomness/uncertainties in the world.
- Models with random variables

[Our concern here is to learn about the language & tools of probability theory or a tool for developing out methods.]

Probability theory is the theory of random numbers

- the frequency with which these numbers can occur is the most useful quantity to take into account.

Brain Thoughts

- Should random numbers be defined solely on basis of a frequency measurement?

OR
Should they be treated as a special kind of object with their inherent property?

FREQUENTISTS' OR

BAYESIANS' ?

Hilary

PROBABILITY THEORY

Random Variable?

Something that takes different values each instant
 (X)

Discrete Variable

Same value all the time.
 (a)

Probability Density Function

The complete knowledge of a random variable

Example: ~~Probability~~

let, $X \in x_i$ $\left[\begin{array}{l} x = 10, 14, 3, \dots \\ x_i \in \mathbb{R} \end{array} \right]$

then,

$$\text{pdf}(x_i) = ? \boxed{\frac{x_i}{X}}$$

i.e. how likely each value is for random variable x .

Hilroy

Regularised Logistic Regression



$$h_{\theta}(u) = g(\theta_0 + \theta_1 u_1 + \theta_2 u_1^2 + \theta_3 u_1^2 u_2 + \theta_4 u_1^2 u_2^2 + \theta_5 u_1^2 u_2^3 + \dots)$$

Cost Function

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(u^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(u^{(i)})) \right]$$

$$+ \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad | \quad \theta_1, \theta_2, \dots, \theta_n \text{ don't become large}$$

Gradient Descent

$$\text{Repeat} \{ \quad \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(u^{(i)}) - y^{(i)}) u_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(u^{(i)}) - y^{(i)}) u_j^{(i)} - \frac{\lambda}{m} \theta_j$$

$j = 1, 2, 3, \dots, n$
 $\theta_1, \theta_2, \dots, \theta_n$

updated equation

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(u^{(i)}) - y^{(i)}) u_j^{(i)} - \frac{\lambda}{m} \theta_j$$

$\frac{\partial J(\theta)}{\partial \theta_j}$

Hilary

Non-invertibility

Suppose $m \leq n$,
then,

$$\theta = \underbrace{(x^T x)}_{\text{non invertible (singular) degenerate always}}^{-1} x^T y$$

\downarrow non invertible (singular) degenerate always

Regularisation takes care of it for us.

$$y \lambda > 0,$$

then,

$$\theta = \left(x^T x + \lambda \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right)^{-1} x^T y$$

\downarrow invertible always

Regularised Normal Equation

$$X = \begin{bmatrix} (x^{(0)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix} \quad Y = \begin{bmatrix} y^{(0)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\rightarrow \min_{\theta} J(\theta)$$

$$\rightarrow \theta = (X^\top X)^{-1} X^\top Y$$

$$\theta = (X^\top X + \lambda \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix})^{-1} X^\top Y$$

(n+1) x (n+1)

Example

$$n=2 \quad \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Regularised Linear Regression

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$$\min_{\theta} J(\theta)$$

Gradient Descent

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$j = 0, 1, 2, \dots, n$

Repeat {

$$\Rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left[\underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j}_{j=1, 2, \dots, n} \right]$$

$$\frac{\partial}{\partial \theta_j} J(\theta)$$



$$\boxed{\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}} \quad \text{Hilroy}$$

\downarrow
 $< 1, \text{ say } 0.99 \Rightarrow 0.99 \theta_j$

If (λ) is very large? (too large for our problem)
say $\lambda = 10^{10}$

We would end up penalizing $\theta_1, \theta_2, \theta_3, \theta_4$ heavily

then they will be all close to zero

which means we are getting rid of them

which then gives a straight line/it

fails to fit the data, because it has

a too high bias!

So, some care should be taken to choose (λ)

Regularised Cost Function

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

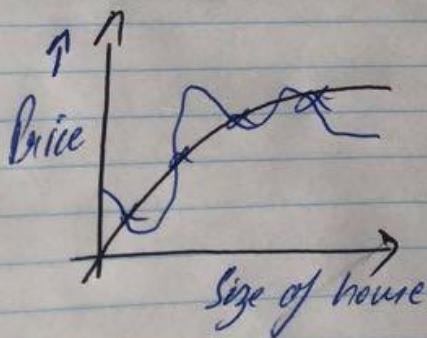
Goal 2
 regularization parameter

Regularisation parameter (λ)

Controls the trade-off b/w two different goals.

- Goal 1: Fit the training set well

- Goal 2: Keep parameter small, to get simple hypothesis, to prevent overfitting



Regularisation in brief

- Small values for parameter $\theta_0, \theta_1, \theta_2, \dots, \theta_n$
- "Simpler" hypothesis
 - less prone to overfitting!

We basically made the parameters sooo small that their effect on the overall cost function is also very small.

Which then helps it generalise better.

Example

Housing:

- Features: $x_0, x_1, x_2, \dots, x_{100}$

- Parameters: $\theta_0, \theta_1, \theta_2, \theta_3, \dots, \theta_{100}$

Which parameter less relevant?

We don't know!!

So we use regularisation.

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$\theta_0, \theta_1, \theta_2, \theta_3, \dots, \theta_{100}$ don't take θ_0 *library*

Bernoulli Distribution [Binary] two outcomes.

A random variable that has two possible outcomes from an experiment.

$$\begin{aligned} \text{Yes} &:= p \\ \text{No} &:= (1-p) \end{aligned}$$

Probability function

$$P(\text{Success}) = p \Rightarrow P(\text{Failure}) = 1-p$$

$$\text{Mean} = p$$

$$\text{Variance} = (1-p)$$

Multinomial Distribution

- 'n' trials, have 'k' possible ~~outcomes~~ outcomes

- the probability of each outcome is p_i

Probability function

$$P(x_i) = \frac{n!}{\prod_{j=1}^k (x_j!)^{\nu_j}} \left(\prod_{j=1}^k p_j^{\nu_j} \right)$$

$$\text{Mean} = np_i$$

$$\text{Variance} = np_i(1-p_i)$$

Hilary

What if we don't know the distribution function?
 - Estimate them from data!

Sample mean [Unbiased]

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Sample Variance [Biased]

$$s_1^2 = \frac{1}{n} \sum_{i=1}^n (\bar{x} - x_i)^2$$

These are the appropriate maximum likelihood estimates of these parameters.

Unbiased Sample Variance

$$s_2^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{x} - x_i)^2$$

useful for gaussian distributed data as all higher moment than variance is zero.

There are infinite pdfs available

Hilary

a.u.a. Standard deviation (std).

Variance is the square of the above std!

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

bump
arbitrary func.
mean
that Variable value

Variance can also be called the second moment
& mean the first moment.

\Rightarrow We thus need ~~at least~~ more than two moments
to characterise a probability function uniquely.

So let's say, we are trying to find the n^{th}
moment out of all moments.

$$m_n = \int_{-\infty}^{\infty} (x - \mu)^n f(x) dx$$

The Variance is the second moment about the mean.

Higher moments define more characteristics like third moment relate to quantity called skewness.
of fourth is Kurtosis.

Hilary

Example (mean isn't the only interesting thing)

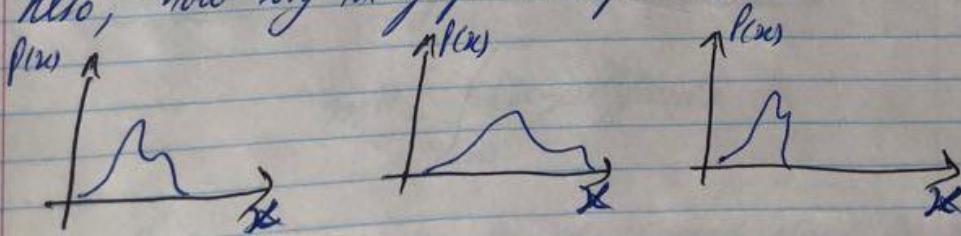
Ask what median is?

Median is some value
for which it is
equally likely to
find a value lower or higher
than itself!!

$$\int_{-\infty}^{\text{median}(x)} p(x) dx = \int_{\text{median}(x)}^{\infty} p(x) dx$$

Thanks to maths mean = median in
a symmetric distribution!

Also, how big the graph is (spread)



gives us an intuition of how far the values are.

But what is μ ?

-we usually don't know, ~~rather we don't~~
~~can't~~

If we knew what μ is, then we could
stop right here but we don't know
what it is!!

That is why we try to approximate the value
of μ .

Example:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N u_i p_i$$

Estimate of the mean!
 probability of the value
 value
 sum of all such values make it easy mathematically

BRAIN THOUGHT

The mean of a distribution is not the only quantity that characterizes a distribution.

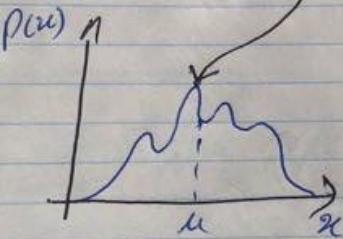
But, think about it!

it would be difficult to store that many values.

So logically we try to take the probability that represents the function the most

i.e. we take the most frequent ~~most random~~ variable/value.

or graphically. Peak! it is the largest peak value of the distribution.



$$\Rightarrow \rho^{\text{more}} = \arg \max_x p(x)$$

More commonly, we ask the average value

which is found by calculating the mean!

$$\Rightarrow \mu = \int_{-\infty}^{\infty} x p(x) dx$$

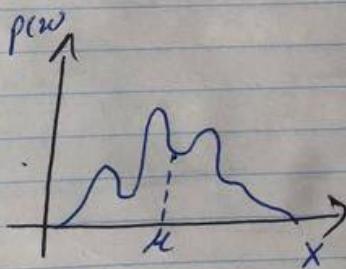
The calculation of adding all the numbers, ~~their frequency~~ together with

State Probability Density Function

- only consider independent random values
- State pdf() describes all that we know about the corresponding random variable.

Example

A poly~~pdf~~: $p(u)$ has values.



Such a distribution is called multimodal!

- Multimodal - having several peaks!

If we overall, pdf() is area under the graph!

- The area under the graph must be 1!

Example

Throw a dice!

$$X = \omega \quad [\omega = \{0, 1, 2, \dots, 6\}]$$

$$P(\omega) = ?$$

| |
|----------|
| ω |
| X |

$$P(\omega) = \frac{1}{6}$$

What is a density function :=

$$P(\omega) = \frac{1}{6} \delta(\omega = \omega_i) \quad [\omega_i = \{1, 2, 3, \dots, 6\}]$$

useful for showing the discrete random process in a continuous form!

as if we see it as a number on paper!

In reality it's different every time you look at it!

It is continuously changing!!

We don't want to take infinite considerations
So, we limit it!

$$P(a < u < b) = \int_{u=a}^b p(u) dx$$

Thankfully, we have a mathematical notation!

DELTA-
function δ -function: write discrete random process
in a continuous form!

\Rightarrow think of it as a density function.

\Rightarrow which is zero except for its arguments.

finite values.

$$\Rightarrow \int_{-\infty}^{\infty} \delta(u - u_0) du = 1$$

introducing
function can then be used as an integrating
function for other functions!!

$$\Rightarrow \int_{-\infty}^{\infty} \delta(u_0) f(u) du = f(u_0)$$

Then, for all values of X . $\underbrace{\quad \quad \quad}_{n}$ my deduction

$$\sum_{x_i} p(x_i) = 1 \Rightarrow p(x_i) = \frac{p(x_i) dx}{\text{prod } x_i}$$

pdf (~~discrete~~ continuous Random Variable)

- we have infinite number of possible values

- which means

$$\frac{x}{X} \ll \ll \ll 0$$

infinitesimally small.

- So! we use calculus to show this
small change due to each ~~bump~~
bump (Value change).

$$\Rightarrow p(x) = \underbrace{p(x) dx}_{\text{pdf}(x)}$$

Then, for all values of X :

$$\int_X p(x) dx = 1$$

we have to sum these infinitesimally small approximations!

calculus gave us integration!!

... find the area under the graph! :)

use log to reduce computation!

$$l(w) = \log L(w) = \sum_{i=1}^m \log(p(w_i; y^{(i)}, x^{(i)}))$$

Since, log function increases monotonically,

$$\max(l) = \max(L)$$

Thus, maximum (log-) likelihood

$$\text{can thus be: } w^{MLE} = \underset{w}{\operatorname{argmax}} l(w)$$

Gradient Descent

$$p(w, \sigma) = \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right)$$

$$= \sum_{i=1}^m \left(\log \frac{1}{\sqrt{2\pi\sigma}} - \frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2} \right)$$

OR

$$L(w, \sigma) = -\frac{m}{2} \log 2\pi\sigma - \underbrace{\sum_{i=1}^m \frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}}_1$$

maximize this

But, one data point isn't enough to make good predictions! we need more.

let's say, we have two data points.

$$(x^{(1)}, y^{(1)}) \Delta (x^{(2)}, y^{(2)})$$

Well, now we need the joint distribution function.

But let's just assume that each feature is independent of each other.

$$\begin{aligned} P(Y_1 = y^{(1)}, Y_2 = y^{(2)} | x^{(1)}, x^{(2)}, w_0, w_1, \sigma) &= \\ &= p(Y_1 = y^{(1)} | x^{(1)}, w_0, w_1, \sigma) p(Y_2 = y^{(2)} | x^{(2)}, w_0, w_1, \sigma) \end{aligned}$$

If we have 'm' samples then

$$P(Y_1 = y^{(1)}, Y_2 = y^{(2)}, \dots, Y_m = y^{(m)} | x_1, x_2, \dots, x_m; w) \boxed{\prod_i^m p(Y_i | x_i, w)}$$

We can now add specific observations (training data) into resulting function.

~~likelihood~~ We then have a function that is a function of parameters.

~~function~~

$$\boxed{L(w) = \prod_i^m p(w; y^{(i)}, x^{(i)})}$$

History

MAXIMUM LIKELIHOOD ESTIMATE

u

Given a parameterized hypothesis function $p(y|w; v)$, we will choose as parameters the values which make the data y most likely under assumption "u".

Consider only one dimension with only one variable. *

$$P(Y_1 = y^{(1)} | x^{(1)}, w_0, w_1, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(1)} - w_0 - w_1 x^{(1)})^2}{2\sigma^2}\right)$$

Now ask which parameter to choose so that the in order to maximise the probability that $y = y^{(1)}$

- How?

- By choosing values!

$\sigma = \text{arbitrary}$

$w_1 = \text{arbitrary}$

$$w_0 = y^{(1)} - w_1 x^{(1)}$$

This is our first MLE! of the parameters
 w_0, w_1, σ .

Hilroy

we keep ' σ ' a specific value

only consider ' w ' as free.

Helps us make it easy

"

So, in this type of learning we ~~keep or~~ consider that
the data is stochastic ie different in every
situation we deem identical."

Revising our hypothesis, according to the central limit theory stochastic nature of data.

Conditional Density Functions

$$p(y|x; w) = p(y|x; w)p(x)$$

Assume feature value are equally likely

so that $p(x)$ is constant &

we only have to focus on $p(y|x; w)$

So

previously, hypothesis = $\hat{y}(x; w)$

now hypothesis = $\hat{y}(x; w) + \text{additive Gaussian noise}$

$$p(y|x; w, \sigma) = N(\mu = w^T x, \sigma^2)$$

$$\boxed{= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - w^T x)^2}{2\sigma^2}\right)}$$

i.e. Gaussian distributed with mean $\mu = \hat{y}(x)$
that depends linearly on the value of x .

So, a Linear Regression

$$\hat{y}(u; w) = w_0 + w_1 u \quad \begin{matrix} \text{how long motor} \\ \text{turned on.} \end{matrix}$$

one variable

Bilinear Hypothesis

$$\hat{y}(u_1, u_2; w) = w_0 + w_1 u_1 + w_2 u_2 \quad \begin{matrix} \text{power setting} \\ \text{two variables} \end{matrix}$$

Sometimes we don't know one of the features.

These features are therefore unknown knowledge

rather than noise? Basically, the point is

We don't care if the noise is either

- A randomness that comes from irreducible indeterminacy, i.e. true randomness in the world that can't be penetrated further for knowledge.

OR

- Noise might represent epistemological limitations such as lack of knowledge of hidden process or limitation in observing states directly.

Hibroy

Probabilistic Models

A simple stochastic generalization of linear regression example to introduce formalism.

Example from Robotics

↳ model how far robot moves

- When wheel are turning after given number of seconds after activating the corresponding motor with certain power.

Then input would be $(x^{(i)}, y^{(i)})$

$x^{(i)}$ is time we let motors run

$y^{(i)}$ distance robot traveled

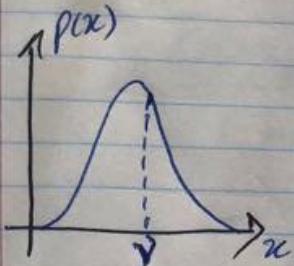
y has to be provided by teacher. We use ultrasonic sensor to automate the collection of data.
(how close the robot is to wall)

We don't need the teacher later on!

Hilary

Chi-Square Distribution

- sum of squares of normally distributed random numbers.

Probability density function

$$p(u) = \frac{u^{(\nu-2)/2} e^{-u/2}}{2^{\nu/2} \Gamma(\nu/2)}$$

mean = v

Scipy.stats.gm

Variance = $2v$

Multivariate Normal Distribution

- several random variables $u_1, u_2, u_3, \dots, u_n$

$$p(u_1, u_2, \dots, u_n) = p(u) = \frac{1}{(2\pi)^n \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

mean = $\vec{\mu}$

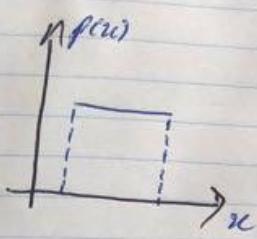
Variance $\Rightarrow \Sigma = \begin{bmatrix} \text{cov}[x_i, x_j] \end{bmatrix}_{i=1,2,\dots,n, j=1,2,\dots,n}$

covariance matrix

Niloy

Uniform Distribution

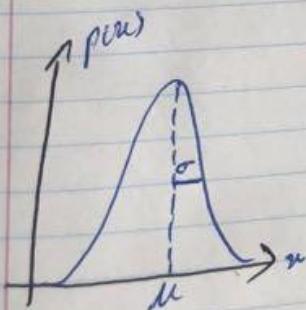
- Equally distributed random numbers
in the interval $a \leq u \leq b$

Probability density function

$$P(u) = \frac{1}{b-a}$$

$$\text{mean} = \frac{a+b}{2}$$

$$\text{Variance} = \frac{(b-a)^2}{12}$$

Gaussian DistributionProbability density function

$$P(u) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(u-\mu)^2}{2\sigma^2}}$$

$$\text{mean} = \mu$$

$$\text{variance} = \sigma^2$$

Hilroy

Binomial Distribution ($k=2$)

Binomial coefficient := $\binom{n}{w} = \frac{n!}{w!(n-w)!}$

- Describes the number of 'successes' w ,
in ' n ' Bernoulli trials,
with probability of success p .

Probability function

$$P(w) = \binom{n}{w} p^w (1-p)^{n-w}$$

$$\text{Mean} := np$$

$$\text{Variance} := np(1-p)$$

GENERATIVE MODELS

- High dimensional data
- First learn ~~out~~ ^{table} about nature of specific
- Predict using this knowledge
- Learning about classes & representing them in neural form is called representational learning.
- Representational models can be used to 'generate' examples of class objects. This model is therefore also called generative model.

$$P(X|Y; \theta)$$

Example = learn hair, learn table
knowledge
use to predict.

What is L^1 ?

L^1 is regularization with a bias in the parameters.

- Assume $\boxed{\text{they}}$ should be normally distributed around zero.
- Value of parameters then,

$$w^{MAP} \rightarrow \arg \max_w p(x, y|w) N(0, \sigma^2)$$

Take log!

$$\boxed{w^{MAP}, \arg \max_w \log(p(x, y|w)) + \alpha \|w\|^2}$$

$$\alpha = \frac{\log 1}{2\sigma^2}$$

thus

LASSO $\Rightarrow L^1$ regularization for Gaussian prior
 Ridge $\Rightarrow L^2$ for non-negative Laplace distribution
 forces more weights toward zero than L^2 !

Hilary

- which can then be used to find the most likely value for the parameters!
- partition function doesn't depend on w . (the denominator)

thus,

$$w^{\text{MAP}} = \underset{w}{\operatorname{argmax}} p(x, y|w) p(w)$$

↑ ↑ ↑
Max approach likelihood of that value prior distribution
More a w
So that it returns
max. value

- treat probability function as function of parameters.

- for uniform $p(w) = \text{const}$

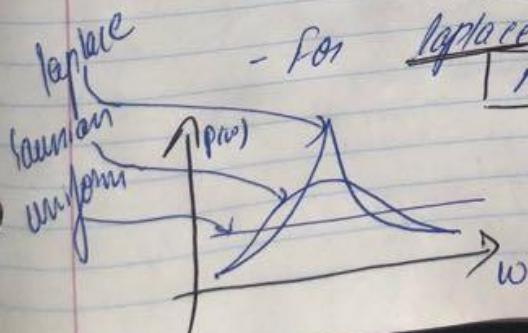
$$w^{\text{MAP}} = w^{\text{MLE}}$$

- For Gaussian distributed

$$\text{MAP} = \text{MLE} + L^2$$

- For Laplace

$$\text{MAP} = \text{MLE} + L'$$



The main problem,

is that we don't know the distribution of w .

- what if feature data is not given?

$$P(y, x|w) = \underbrace{P(y|x; w)}_{\text{conditional density}} \underbrace{P(x|w)}_{\text{marginal density}},$$

- Now choose the feature data uniformly.
- So, that the marginal distribution is a constant

$$P(x|w) = \text{constant}$$

Hence,

$$P(y, x|w) \propto P(y|x; w)$$

Use Bayes theorem,

$$P(w|x, y) \propto \frac{P(y, x|w) P(w)}{\int_{w \in W} P(x, y|w) P(w) dw}$$

Partitioning factor
 domain of possible parameter values

MAXIMUM \hat{P} POSTERIORI

Assume we know $p(w|x,y)$, we can then find the most probable value, w^{MAP} , given data set.

$$w^{\text{MAP}} = \underset{w}{\operatorname{argmax}} p(w|x,y)$$

- MLE is work horse of probabilistic supervised learning.
- For probabilistic models, parameters should be selected based on data.

$$p(w|x,y)$$

Assume, we know the conditional probability
we can then choose a parameter value
we like.

$$\int w^{\text{MAP}} = \underset{w}{\operatorname{argmax}} p(w|x,y)$$

Value w , which is the most probable value

My $x < u_0$

$$e_2 = \frac{1}{2} - \frac{1}{2} \operatorname{erf}\left(\frac{u_0}{\sqrt{2}}\right) P(c_2)$$

thus, we have

$$P(c_1) = P(c_2) = 1/2$$
$$\sigma^2 = 1, \mu = 2$$
$$Acc = \frac{\left(1 - \operatorname{erf}\left(\frac{1}{\sqrt{2}}\right)\right)}{2} \approx 0.84$$

Take $\mu = 2$ & $\sigma = 1$

The bayesian decision point is where the two ~~post~~ posteriors are equal.

$$e^{\frac{-x_0^2}{2\sigma^2}} p(c_1) = e^{\frac{-(x_0 - \mu)^2}{2\sigma^2}} p(c_2)$$

$p(c_1)$ & $p(c_2)$ are marginal class probabilities

$$u_0 = \frac{1}{2} \mu + \frac{\sigma^2}{\mu} (p(c_2) - p(c_1))$$

To find the error, find the area under the graph with wrong classification.

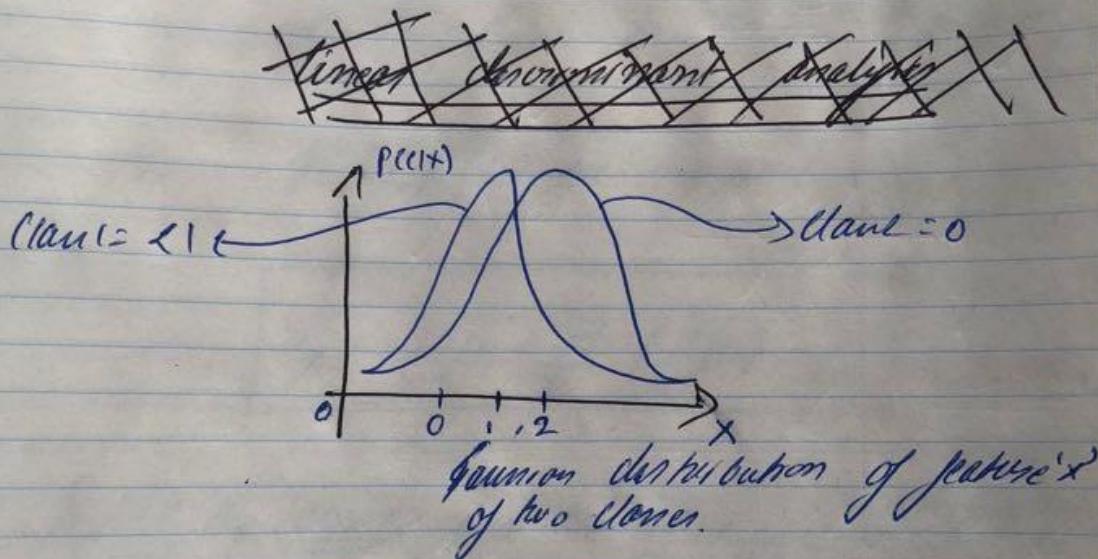
Example: we predict c_1 wrongly for $x > u_0$ in the following fraction of cases

$$e_{c_1} = \frac{1}{\sqrt{2\pi}\sigma} \int_{u_0}^{\infty} e^{\frac{-x^2}{2\sigma^2}} dx p(c_1)$$

$$\Rightarrow e_{c_1} = \frac{1}{2} - \frac{1}{2} \operatorname{erf}\left(\frac{x_0}{\sqrt{2\sigma}}\right) p(c_1)$$

Supervised Gaussian Generative Models

- 1D Gaussian Example



- The features are overlapping i.e.
the classes aren't fully discriminated
by this feature.

\Rightarrow we can only classify values
less than 0.1 as C1 & > 0.1 as C2

$$p(u|y=C_1) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{u^2}{2\sigma^2}}$$

$$p(u|y=C_2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(u-\mu)^2}{2\sigma^2}}$$

In order to make a discriminative model from generative model, we need (class priors).

What is the relative frequency of the classes is.

Then,

$$P(y|x; \theta) = \frac{P(x|y; \theta) P(y)}{P(x)}$$

- The Bayesian way of deciding which class has the target posterior probability

$$\operatorname{argmax}_{(y)} P(y|x; \theta) = \operatorname{argmax}_{(y)} \frac{P(x|y; \theta) P(y)}{P(x)}$$

Since, denominator has no 'y'
we get \Rightarrow

$$\boxed{\operatorname{argmax}_{(y)} P(x|y; \theta) P(y)}$$

1/y - Binary Classification

$$\operatorname{argmax}_{(y)} P(y|x; \theta) = \operatorname{argmax}_{(y)} P(x|y=0, \theta) P(y=0), \\ P(x|y=1, \theta) P(y=1)$$

Bayesian Decision Point / Dividing Hyperplane

$$P(y=1/x^T) = P(y=0/x^T) = 0.5 \rightarrow x^T w + b = 0$$

- Values higher than this would result in positive output.

Sigmoid function

$$h(x; w) = \frac{1}{1 + e^{-(w^T x)}}$$

$$\begin{aligned}\frac{\partial h}{\partial w} &= \frac{2}{2w} \frac{1}{1 + e^{-wx}} \\ &= \frac{1}{(1 + e^{-wx})^2} e^{-wx}(-x) \\ &= \frac{1}{(1 + e^{-wx})} \left(1 - \frac{1}{(1 + e^{-wx})}\right)(-x) \\ &= -h(1-h)x\end{aligned}$$

using this in ② & putting it in ①

$$\begin{aligned}y\left(\frac{1}{h} - (1-y)\frac{1}{1-h}\right)h(1-h) &= y(1-h) - 1(y-h)h \\ &= y - yh - h + yh \\ &= y - h\end{aligned}$$

Given learning rule,

$$w_j \leftarrow w_j + \alpha (y^{(i)} - h(x^{(i)}; w)) x_j^{(i)}$$

Same as linear regression, but hypothesis is different!

Hilary

Logistic Regression

$$P(y=1|x; \omega) = h(x; \omega)$$

$$P(y=0|x; \omega) = 1 - h(x; \omega)$$

hypothesis

Combined,

$$P(y|x; \omega) = (h(x; \omega))^y (1 - h(x; \omega))^{1-y}$$

Log-Likelihood

$$\ell(\omega) = \sum_{i=1}^m y^{(i)} \log(h(x_i; \omega)) + (1 - y^{(i)}) \log(1 - h(x_i; \omega))$$

- To find maximum,

Gradient Ascent!

$$\omega \leftarrow \omega + \alpha \nabla_{\omega} \ell(\omega) \quad \text{--- (1)}$$

We calculate partial derivative of log-likelihood function w.r.t. parameters.

$$\frac{\partial \ell(\omega)}{\partial \omega_j} = \left(\frac{y}{h} - \frac{(1-y)}{1-h} \right) \frac{\partial h(\omega)}{\partial \omega_j} \quad \text{--- (2)}$$

Let's try to apply it!

MLE of Bernoulli model

$$P(Y) = \underbrace{\phi^y}_{\rightarrow P(Y=1)} \underbrace{(1-\phi)^{1-y}}_{\rightarrow P(Y=0)}$$

Bernoulli variable \Rightarrow characterized by only one parameter ' ϕ '.

Let's consider ' m ' tosses in which ' h ' heads have been found. [Coin Toss]

Log-Likelihood of m trials is:

$$\begin{aligned} l(\phi) &= \log \prod \phi^{y^{(i)}} (1-\phi)^{1-y^{(i)}} \\ &= \log (\phi^h (1-\phi)^{m-h}) \end{aligned}$$

To find max ' ϕ ' we set the derivative of l to 0.

$$\frac{d l}{d \phi} = \frac{h}{\phi} - \frac{m-h}{1-\phi}$$

$$\frac{h}{\phi} - \frac{m-h}{1-\phi} = 0$$

$$\Rightarrow \boxed{\phi = \frac{h}{m}} \quad \begin{array}{l} \rightarrow \text{Fraction of } h \text{ in trials!} \\ \text{This proves our intuition!} \end{array}$$

Hurray

That quadratic error term is

$$E = \frac{1}{2} (y - h(x; w))^2 \Leftrightarrow p(y|x; w) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - h(x; w))^2}{2\sigma^2}\right)$$

Square error loss function

least mean square (LMS) error if we consider batch algorithm

Polynomial loss function

$$E = \frac{1}{p} \|y - h(x; w)\|^p \Leftrightarrow p(y|x; w) = \frac{1}{2\Gamma(1/p)} \exp\left(-\frac{\|y - h(x; w)\|^p}{2\sigma^p}\right)$$

ϵ loss function

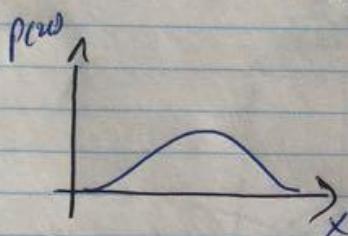
Errors less than ϵ don't contribute to error measure

$$E = \|y - h(x; w)\|/\epsilon \Leftrightarrow p(y|x; w) = \frac{1}{2(1-\epsilon)} \exp\left(-\frac{\|y - h(x; w)\|/\epsilon}{2(1-\epsilon)}\right)$$

Example Consider two gaussian distributions having same likelihood, generate data from both having equal likelihood \Leftrightarrow to be from either.

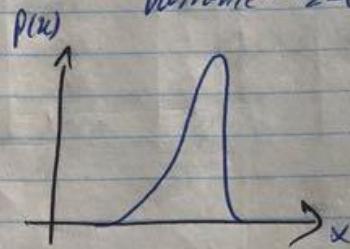
data1 mean = $\mu_1 = -1$

Variance $\sigma_1^2 = 2$



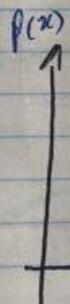
data2 mean = $\mu_2 = 4$

Variance $\sigma_2^2 = 0.5$



↓ Equal Likelihood

Initial condition



generative models

Real data with gaussian distribution

We could do this,

choose the heuristics to match the actual data-generating system (world).

E-Step: Make Assumption of label from data

M-Step: Use this hypothesis to update the parameters of the model to maximise the probability of observations.

Hilroy

9

Algorithm (Estimation, maximization)① Initialize parameters ϕ, μ, Σ randomly!

② Repeat until convergence!

{

Estimation step:

For each data point 'i' & class 'j'
"soft-classify" data as

$$w_j^{(i)} = p(z^{(i)}=j | x^{(i)}; \phi, \mu, \Sigma)$$

Maximization step:

Update the parameter according to

$$\phi_j = \frac{1}{m} \sum_{i=1}^m w_j^{(i)}$$

$$\mu_j = \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}}$$

$$\Sigma_R = \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}}$$

γ convergence

Hilroy

But we ~~don't~~ know label! AKA don't know $\hat{z}^{(i)}$

- what if we knew? Then,

log-likelihood would be

$$\ell(\Phi, \mu, \sigma) = \sum_{i=1}^m \log p(z^{(i)}; \mu, \sigma)$$

& its maximum likelihood of parameters would be

$$\mu_k = \frac{1}{m} \sum_{i=1}^m \log \mu_{z^{(i)}=j}$$

$$\mu_{ik} = \frac{\sum_{i=1}^m \log \mu_{z^{(i)}=j} x^{(i)}}{\sum_{i=1}^m \log \mu_{z^{(i)}=j}}$$

$$\Sigma_{ik} = \frac{\sum_{i=1}^m \log (\mu_{z^{(i)}=j}) (x^{(i)} - \mu_j) (x^{(i)} - \mu_j)^T}{\sum_{i=1}^m \mu_{y^{(i)}=k}}$$

- The function ~~$\mu_{z^{(i)}=j}$~~ is simply 1 if $x=y$, otherwise zero.

So, to solve our problem of not knowing the labels
JUST ASSUME ALL & let algorithm improve it!!

Hilary

7

MIXTURE OF GAUSSIAN AND EM ALGORITHM

- try to combine k-means with generative models, so that each class is then a specific generative model.

Example (a gaussian model)

we can use any ~~one~~ type of distribution for this.

- Assume 'K' class, each class chosen randomly from a multinomial distribution.

$\rho(z^{(i)} = j) \propto \text{multinomial}(\bar{\Phi}_j)$

$\boxed{\rho(x^{(i)} | z^{(i)} = j) \propto N(\mu_j, \Sigma_j)}$

Gaussian mixture model →

- log-likelihood function

$$\ell(\bar{\Phi}, \mu, \Sigma) = \sum_{i=1}^m \log \sum_{z^{(i)}=1}^K \rho(x^{(i)} | z^{(i)}; \mu, \Sigma) \rho(z^{(i)}; \bar{\Phi})$$

Hilroy

K-MEAN CLUSTERING

- no labels
- learning is still guided though it follows specific principles that are used to organise the data based on characteristics provided by the data.
- Name random k means
- Sort all to clusters from each k .
- find new means
- repeat ② & ③ until you get the same means as the first set of means.

Algorithm

- ① Initialize mean randomly
- ② Repeat until convergence

{ Model prediction :

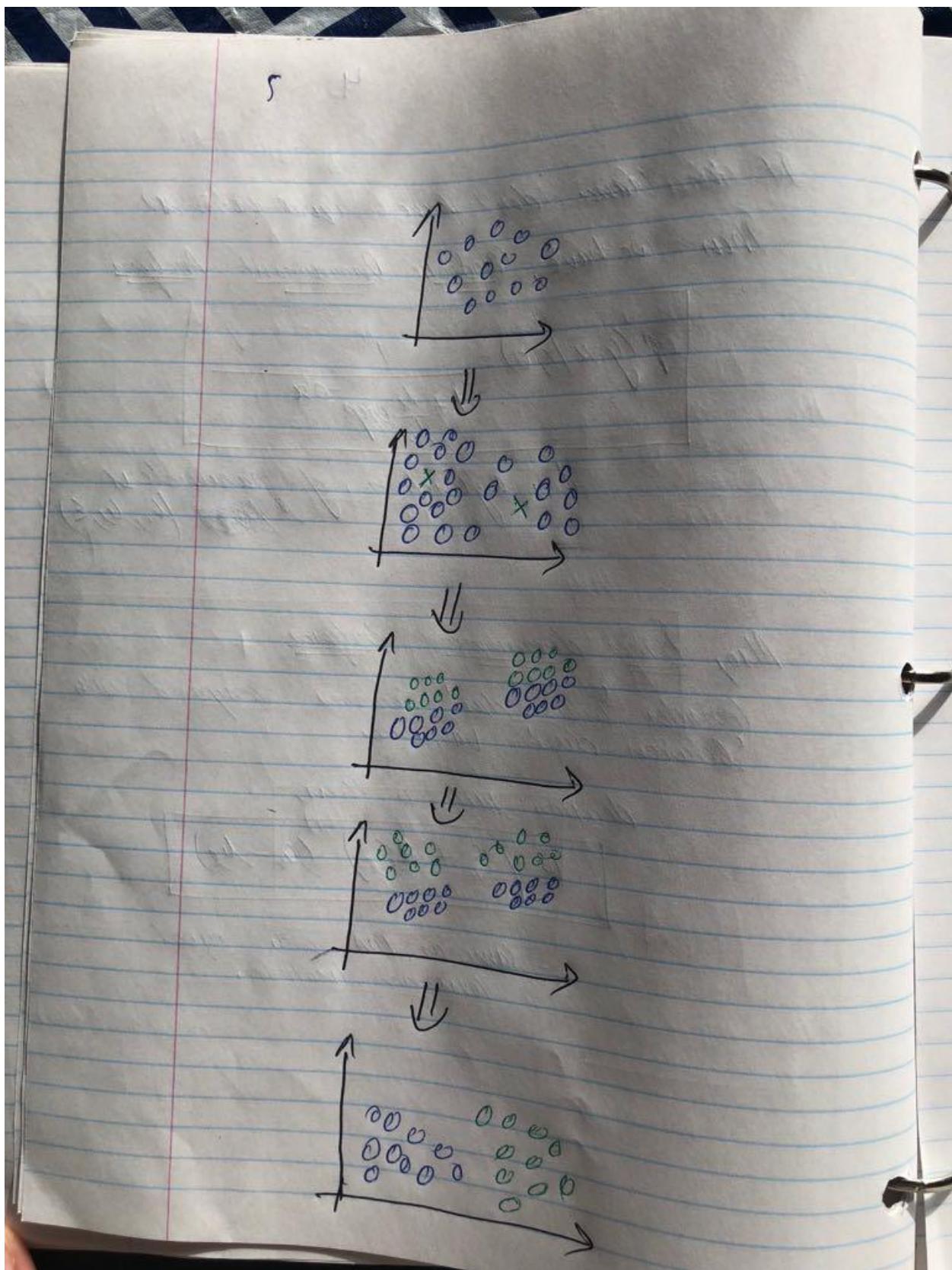
for each data point 'i', classify data with closest mean

$$c^{(i)} = \underset{j}{\operatorname{arg\,min}} \|x^{(i)} - \mu_j\|$$

Model Refinement :

calculate new mean for each class

$$\mu_j = \frac{\sum_{c^{(i)}=j} x^{(i)}}{\sum_{c^{(i)}=j}}$$



4

If the classes don't have equal variances?

Then we have, Quadratic Discriminant Analysis

$$P(Y=k|x_i; \theta) = \frac{1}{1 + \frac{\rho_k}{\rho_n} \exp^{-\theta^T x}}$$

ϕ is a function of parameters ~~(μ_1, μ_2, Σ)~~ (μ_1, μ_2, Σ)
variance

Also, Fisher discriminant analysis

involves within-class variances compared with between-class variances.

decision rule is thus,

$$w = (\Sigma_k + \alpha \Sigma_l)^{-1} (\mu_k - \mu_l)$$

- same as LDA with equal covariance matrices.

3

b/w classes '1' & '2'

To calculate decision boundary b/w

$$\Sigma_k := \Sigma$$

- Assume covariance matrix of the classes are the same.



Now, the decision point b/w two classes with equal class priors is given by the point where probabilities for the two classes are equal!

Thus,

$$\log\left(\frac{\phi_1}{\phi_2}\right) - \frac{1}{2}(\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 + \mu_2) + \frac{1}{2} \Sigma^{-1} (\mu_1 - \mu_2) = 0$$

[Don't have 'x', so they are constant!]

$$\Rightarrow w = \Sigma^{-1} (\mu_1 - \mu_2)$$

Simplify \Rightarrow

$$a + w^T x = 0$$

It is linear!!

\Rightarrow This method with Gaussian class distribution with equal variances is called linear discriminant analysis (LDA)

Vector w is \perp to Gaussian surface

Hilary

2

By MLE of gaussian parameters

(mean) $\mu_K = \frac{1}{|K|} \sum_{i \in K} x^{(i)}$

(Variance) $\Sigma_K = \frac{1}{|K|} \sum_{i \in K} (x^{(i)} - \mu_K)(x^{(i)} - \mu_K)^T$

Now, we can calculate Bayesian decision rule
 $\Rightarrow g(x; \theta)$

$\Rightarrow g(x) = \operatorname{argmax}_K p(y=K|x)$

$\Rightarrow \operatorname{argmax}_K \left[p(x|y=K; \theta) p(y=K) \right]$

$\Rightarrow \operatorname{argmax}_K \left[\log(p(x|y=K; \theta)) + \log(p(y=K)) \right]$

$\Rightarrow \operatorname{argmax}_K \left[-\log(\sqrt{2\pi} \sqrt{\det(\Sigma_K)}) - \frac{1}{2} (x - \mu_K)^T \Sigma_K^{-1} (x - \mu_K) + \log(\phi_K) \right]$

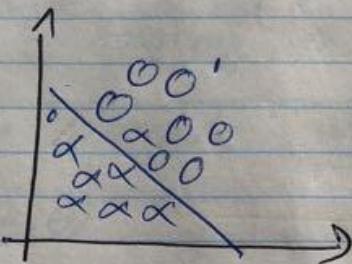
$\Rightarrow \operatorname{argmax}_K \left[\frac{1}{2} x^T \Sigma_K^{-1} x - \frac{1}{2} \mu_K^T \Sigma_K^{-1} \mu_K + x^T \Sigma_K^{-1} \mu_K + \log(\phi_K) \right]$

depend on class

We now have everything to make this decision.

- Linear Discriminant Analysis Gaussian class distribution with equal variances

We now have 'K' class with gaussian distribution over 'n' attributes.



- Each class has class prior

$$P(y=k) = \phi$$

- Each class is multivariate Gaussian distributed each with different μ_k & Σ_k

$$p(x|y=k) = \frac{1}{\sqrt{2\pi}^n \sqrt{|\Sigma_k|}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)}$$

- Use ME to ~~predict~~ estimate the value of parameters

$$\theta = (\phi_k, \mu_k, \Sigma_k)$$

class priors. $\phi_k = \frac{|T_k|}{M}$ no. of examples in this set
Set of examples of class k.

Hilary

20

Well then, probability of this word is 0.

$$\phi_{j,y=1} = 0$$

$$\phi_{j,y=0} = 0$$

Hence

$$P(y=1/x) = \frac{0}{0}$$

A trick! Laplace Smoothing.

- Add one occurrence of this word in every case which will add a small probability proportional to no. of training examples to estimate

$$\Rightarrow \phi_{j,y=1} = \frac{\sum_{i=1}^m I\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\} + 1}{\sum_{i=1}^m I\{y^{(i)} = 1\} + 2}$$

$$\Rightarrow \phi_{j,y=0} = \frac{\sum_{i=1}^m I\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\} + 1}{\sum_{i=1}^m I\{y^{(i)} = 0\} + 2}$$

19

$$\Rightarrow \phi_{j,y=1} = \frac{1}{|\{y=1\}|} \underset{i \in \{y=1\}}{\leq} x_j^{(i)}$$

$$\Rightarrow \phi_{j,y=0} = \frac{1}{|\{y=0\}|} \underset{i \in \{y=0\}}{\leq} x_j^{(i)}$$

$$\Rightarrow \phi_y = \frac{|\{y=1\}|}{m}$$

$$\boxed{\phi_{j,y=1}}$$

probability of being in spam-sent

$$\boxed{\phi_{j,y=0}}$$

probability of being in non-spam-sent

$$\boxed{\phi_y}$$

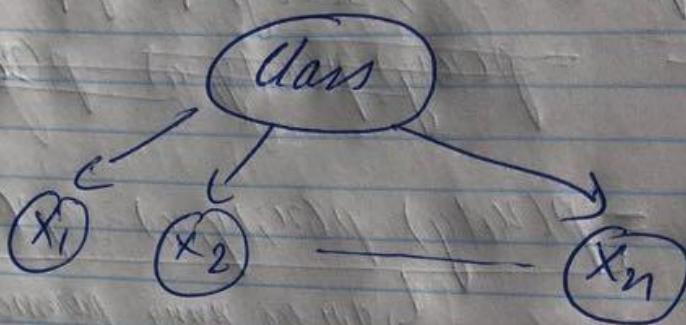
ϕ_y = frequency of spam examples in the data set.

- probability that email 'x' is spam.

$$\boxed{P(y=1|x) = \frac{\prod_{j=1}^{50000} \phi_{j,y=1} \phi_{y=1}}{\prod_{j=1}^{50000} \phi_{j,y=1} \phi_{y=1} + \prod_{j=1}^{50000} \phi_{j,y=0} \phi_{y=0}}}$$

\Rightarrow What if we get a word that is not in our vocabulary??

18



Naive Bayes Classifier!

17

- we can factorize it!!

$$P(u_1, u_2, \dots, u_{50000} | y) = P(u_1) P(u_2) \dots P(u_{50000})$$

$$= p(u_1 | y) p(u_2 | y, u_1) \dots p(u_{50000} | y, u_1, u_2, \dots, u_{49999})$$

$$= p(u_1 | y) p(u_2 | y, u_1) \dots p(u_{50000} | y, u_1, u_2, \dots, u_{49999})$$

↳ still too many parameters!

Naive Bayes assumption:

let's assume each word is conditionally independent to each other!!

then,

$$p(u_1 | y) p(u_2 | y, u_1, u_2) \dots p(u_{50000} | y, u_1, u_2, \dots, u_{49999})$$

$$= p(u_1 | y) p(u_2 | y) \dots p(u_{50000} | y)$$

Only 50000 parameters to deal with now -

so, conditional probability

$$P(x|y) = \prod_{j=1}^{50000} p(u_j | y)$$

& the parameters can then be found by maximization, likelihood estimation

Hilary

16

NAIVE BAYES

Example

Classify spam & not spam
 $(y=1)$ $(y=0)$

Store all the words in a vocabulary vector.

$$x = \begin{bmatrix} a \\ ab \\ c \\ d \\ e \\ i \end{bmatrix}$$

Then, text is represented by 1 if it has that word.

$$x = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

If this vocabulary has the count of each word then it is called a "bag of words".

So, our vocabulary has 50,000 words now.

i.e. we want to model

$$P(x|y) = P(u_1, u_2, \dots, u_{50000} | y)$$

- This has $\star 2^{50,000}$ -1 parameters!!

~~$x_0 = 1$~~ [normalization condition]

Example

- jointed distribution of the five variables

$$P(B, E, A, J, M) = P(B|E, A, J, M) P(E|A, J, M) \\ P(A|J, M) \\ P(J|M) \\ P(M)$$

- However, we have a specific factorization

$$P(B, E, A, J, M) = P(B) P(E) P(A|B, E) P(J|A) P(M)$$

which is much easier to handle!!

- What if we want to deduce about parent from child nodes?

$$P(A=t | J=t) = \frac{P(J=t | A=t) P(A=t)}{P(J|A) \underset{\substack{t \\ \text{true}}}{P(A)} + P(J|A) \underset{\substack{\bar{t} \\ \text{false}}}{P(A)}}$$

OR

$$P(A|J) = \underbrace{P(J|A)}_{\substack{\text{true} \\ \text{true}}} \underbrace{P(A)}_{\substack{\text{true}}} \frac{\underbrace{P(J|A) P(A)}_{\substack{\text{true}}} + \underbrace{P(J|A) P(A)}_{\substack{\text{false}}}}{\underbrace{P(J|A) P(A)}_{\substack{\text{true}}} + \underbrace{P(J|A) P(A)}_{\substack{\text{false}}}}$$

Hilary

- For this, we have to add knowledge or hypothesis about causal relations!

Causal relation

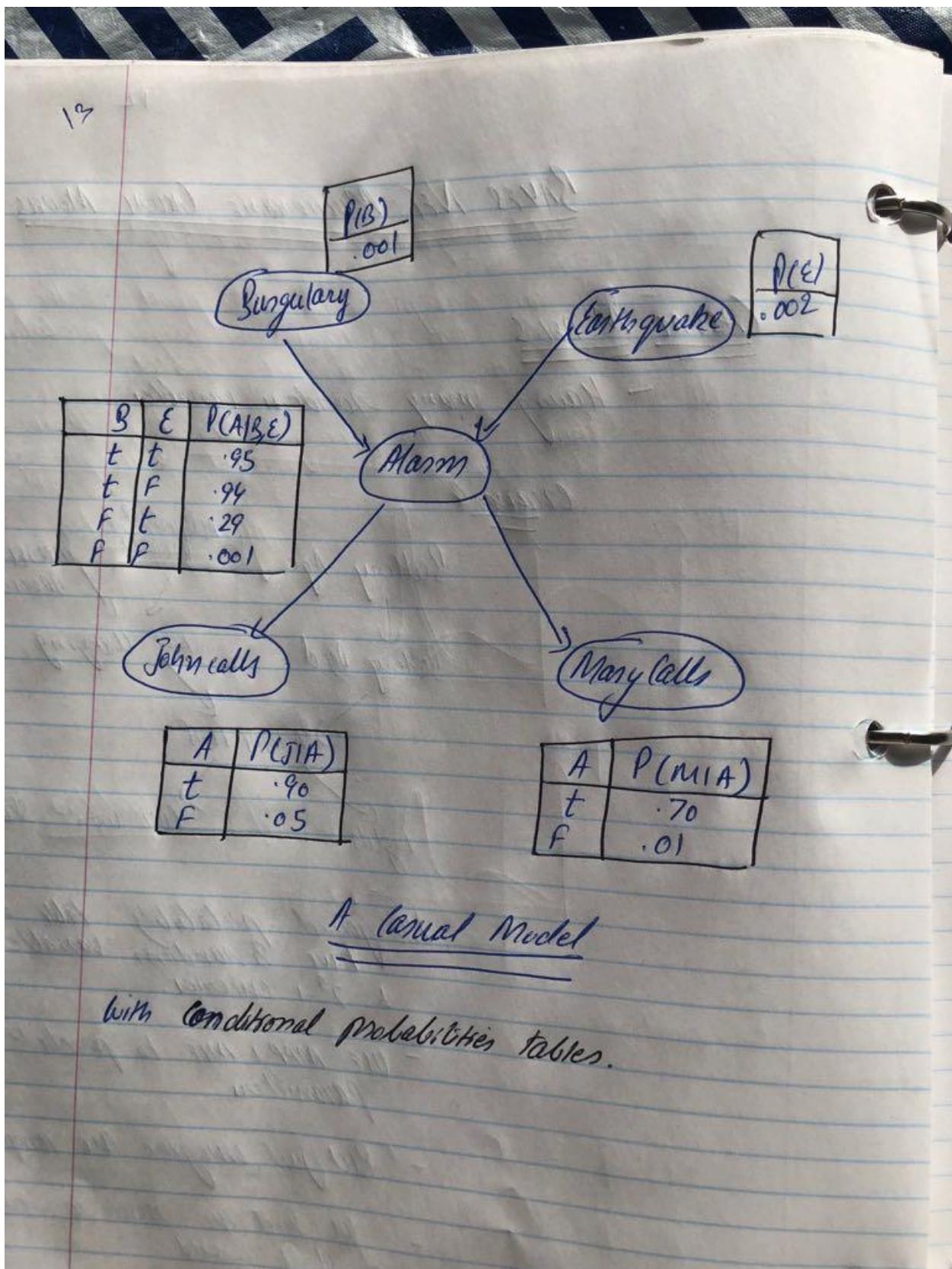
{ Suppose a fire alarm off everytime there is fire
 But ~~will~~ go off, there won't be fire?

Explainably



causal relation $P(A|B)$

- There are then graphical models
- We only consider graphs without loops for now.
- Thus, we are working with graphs with arrows + no loops called Directed acyclic graphs.

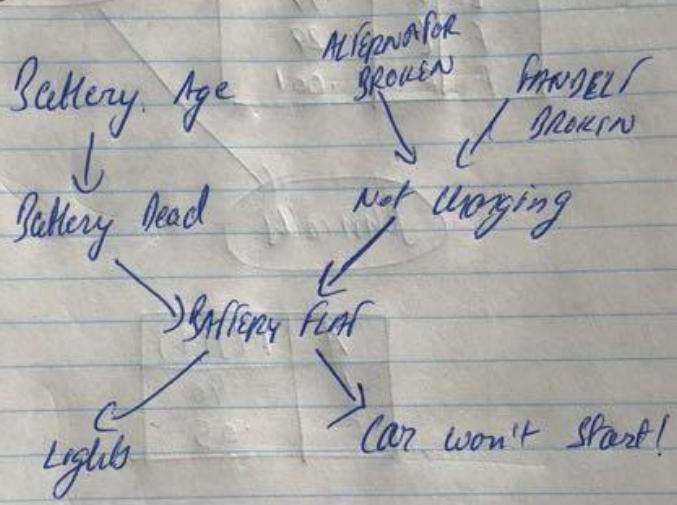


BAYES NET: MULTIVARIATE CAUSAL MODELS

CAUSAL MODELS

- Having many random variables

example



The joint probability table will have 8 random variables
 $\hookrightarrow 2^8 - 1$ entries.

- USE MAP & MLE to predict parameters.
- But just imagine the complexity!!
- We need the model to reason about world

Hilary

- We choose parameters randomly!
- We can assign each data point to the class which produces the larger probability, within the current model
- Thus, we are using ~~the hypothesis or recognition model~~ recognition model!

