

Vid I

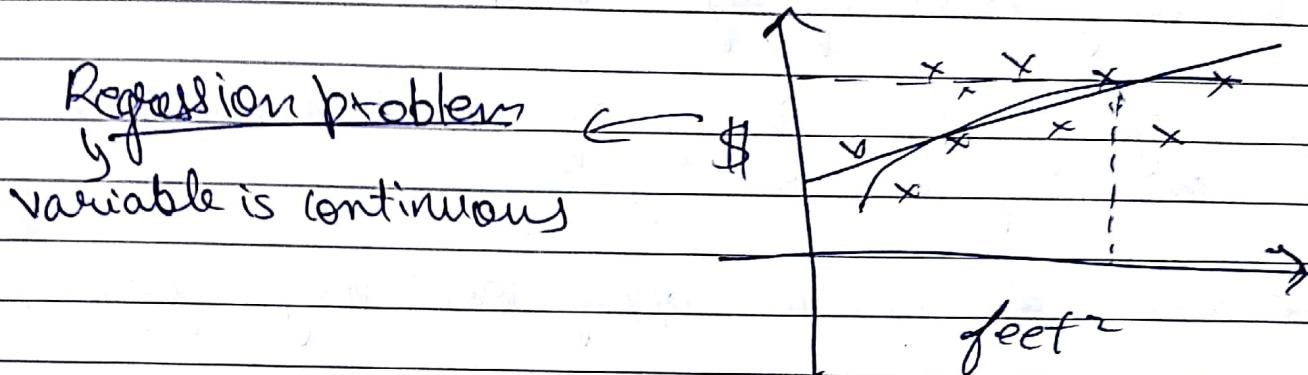
Machine Learning

Defn:- Arthur Samuel (1959)

Ability of computers to learn without being explicitly programmed

Tom Mitchell (1998) - Learn from exp E, ~~Problem~~
Measure P, P increases for T with E.

Task T

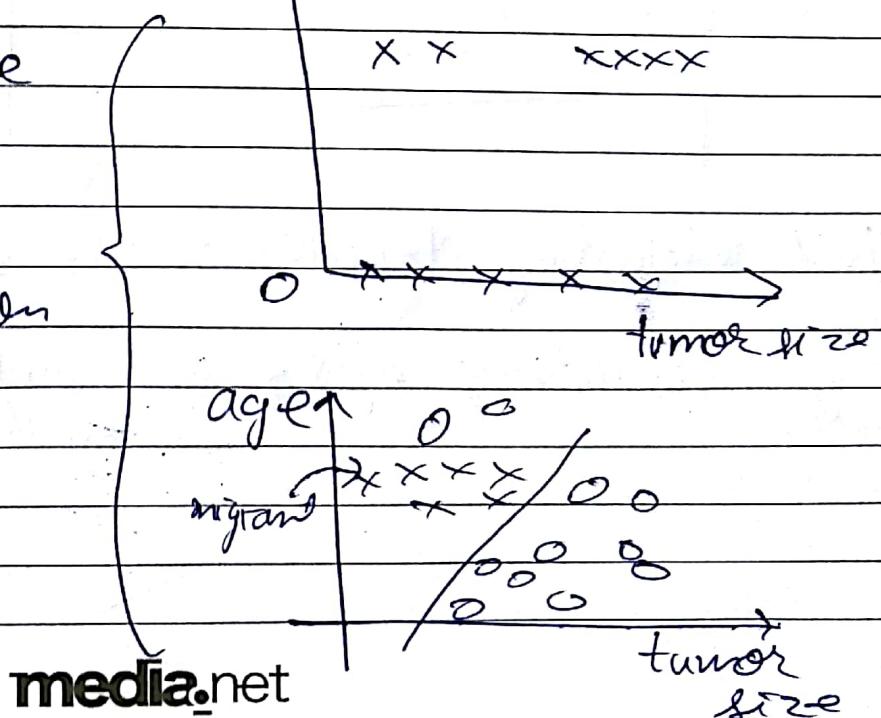


Supervised learning
provided data with labels

Classification Problem malignant

variable is discrete

Supervised classification



media.net

Support Vector Machines

↳ infinite dimensional vector
data changed to T ↳ classification

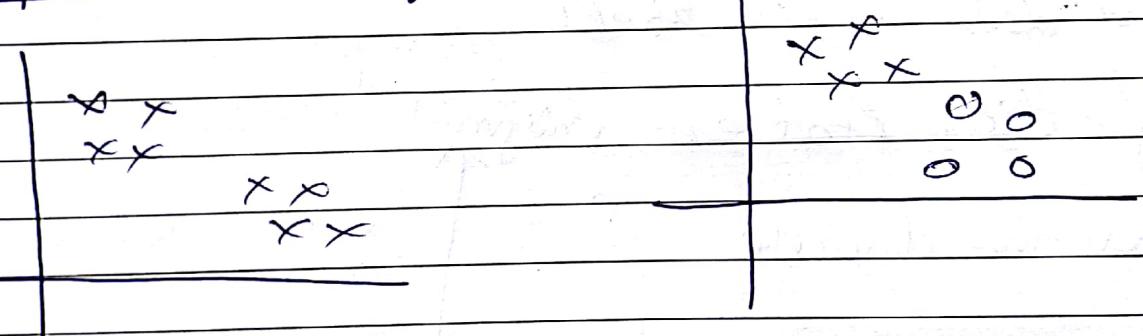
Learning Theory

how & why learning algo. works?

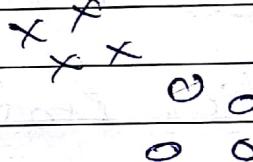
E.g.
Theorem:- guarantees learning algo. with predict 99.9%
correctly on reading of code

Take tools of ML, and apply tools very craftfully.

Unsupervised learning



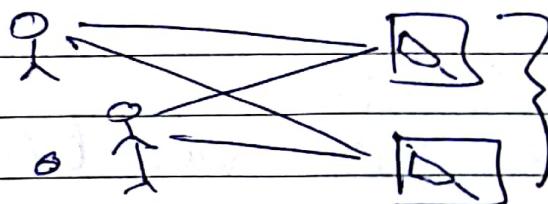
Supervised



Find interesting structure in the data

E.g. clustering algs in computing clusters
social network clusters

Cocktail party problem



Take these two audios &
separate the audios

Independent Component Analysis

Reinforcement Learning

We have reward function. Give positive reward for desirable action else negative.

Vid 2

1994 - Autonomous driving

Alwin - Jeep driver - Network returns steering angle and confidence is returned.

↳ Regression Problem

| Living area (feet ²) | \$ (1000s) |
|----------------------------------|------------|
| 2104 | 400 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| 1940 | 240 |

Notation

$m = \# \text{ training examples}$

$x = \text{"input" variables / features}$

$y = \text{"output" variable / "target" variable}$

$(x, y) - \text{training e.g.}$

$\text{; training example} = (x^{(i)}, y^{(i)})$

Training set

Learning algo.

new

living
area

\downarrow
hypothesis

\boxed{h}

estimated
price

$$h(x) = \cancel{b_0} + \cancel{b_1} \theta_0 + \theta_1 x$$

$x_1 = \text{size}, x_2 = \# \text{bedrooms}$

$$h(x) = h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

For conciseness, $x_0 = 1$

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

$n = \# \text{features}$

θ 's are called parameters

Learning algos & find the value of θ .

Can θ be function of other variables? Yes

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

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

minimize $J(\theta)$

Start with some θ . (Say $\theta = \vec{\theta}$)

Keep changing θ to reduce $J(\theta)$

Gradient Descent

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

$$\frac{\partial J(\theta)}{\partial \theta_i} = \frac{\partial}{\partial \theta_i} \left[\frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right]$$

$$= 2 \cdot \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_i} (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$= (h_{\theta}(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_i} (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\hookrightarrow (h_{\theta}(x^{(i)}) - y^{(i)}) \theta_i x_i$$

for one training e.g.

Repeat till convergence :-

$$\theta_i = \theta_i - \alpha \sum_{j=1}^m (h_\theta(x^{(j)}) - y^{(j)}) \cdot x_i^{(j)}$$

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

For sum of squares, $J(\theta)$ is a bowl shaped with one local minima

Batch Gradient Descent

Look at entire training set, all m instances are used for every update

Stochastic Gradient Descent

Repeat {

for. $j=1$ to m {

$$\theta_i := \theta_i - \alpha (h_\theta(x^{(j)}) - y^{(j)}) \cdot x_i^{(j)}$$

3

→ Works faster than ~~batch~~ gradient descent

$$\nabla_{\theta} J = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix} \in \mathbb{R}^{n+1}$$

Gradient Descent:

$$\theta := \theta - \alpha \nabla_{\theta} J$$

$$y \quad R^{n+1} \quad \quad \quad R^{n+1}$$

$$f: \mathbb{R}^{m \times n} \mapsto \mathbb{R} \quad f(A) \quad A \in \mathbb{R}^{m \times n}$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \cdots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$$

$$\text{If } A \in \mathbb{R}^{n \times n} \quad \text{tr} A = \sum_{i=1}^n A_{ii}$$

Fact:

$$\text{tr} AB = \text{tr} BA$$

$$\text{tr} ABC = \text{tr} CAB = \text{tr} BCA$$

$$\text{If } f(A) = \text{tr}(AB) \quad \nabla_A \text{tr} AB = B^T$$

$$\text{If } a \in \mathbb{R}, \quad \text{tr} a = a$$

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

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

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

$$X\theta - y = \begin{bmatrix} h(x^{(1)}) - y^{(1)} \\ \vdots \\ h(x^{(m)}) - y^{(m)} \end{bmatrix}$$

$$\frac{1}{2} (X\theta - y)^T (X\theta - y) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 = J(\theta)$$

$$\nabla_{\theta} J(\theta) \stackrel{\text{set}}{=} \bar{0}$$

$$\nabla_{\theta} \frac{1}{2} (x\theta - y)^T (x\theta - y)$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr} \left(\theta^T x^T x \theta - \theta^T x^T y - y^T x \theta + \cancel{y^T y} \right)$$

So trace \rightarrow Real number
is introduced

$$= \frac{1}{2} \left[\nabla_{\theta} \text{tr} \underbrace{\theta \theta^T x^T x}_{\substack{\text{applied cyclic} \\ \text{prop. of trace}}} - \nabla_{\theta} \text{tr} \underbrace{y^T x \theta}_{\substack{\text{transpose of} \\ \text{real is} \\ \text{same real}}} \right]$$

Independent
of θ

$$= \nabla_{\theta} \text{tr} \cancel{y^T x \theta}$$

$$\nabla_{\theta} \text{tr} \underbrace{\theta I \theta^T x^T x}_{\substack{\text{www} \\ \text{A B A}^T \\ \text{C}}} = \underbrace{x^T x \theta I}_{\substack{\text{www} \\ \text{C A B}}} + \underbrace{x^T x \theta^T I}_{\substack{\text{www} \\ \text{C}^T A B^T}}$$

$$\nabla_{\theta} \text{tr} \underbrace{y^T x \theta}_{\substack{\text{www} \\ \text{B A}}} = \underbrace{x^T y}_{\substack{\text{B}^T}}$$

$$\nabla_{\theta} J(\theta) = \frac{1}{2} \left[x^T x \theta + x^T x \theta - x^T y - x^T y \right]$$

$$= x^T x \theta - x^T y \stackrel{\text{set}}{=} 0$$

$$\Rightarrow x^T x \theta = x^T y$$

$$\Rightarrow \theta = (x^T x)^{-1} x^T y$$

$$\Rightarrow \theta = (x^T x^{-1}) x^T y$$

~~#Vid 2~~ # Vid 3

Outline

Linear regression

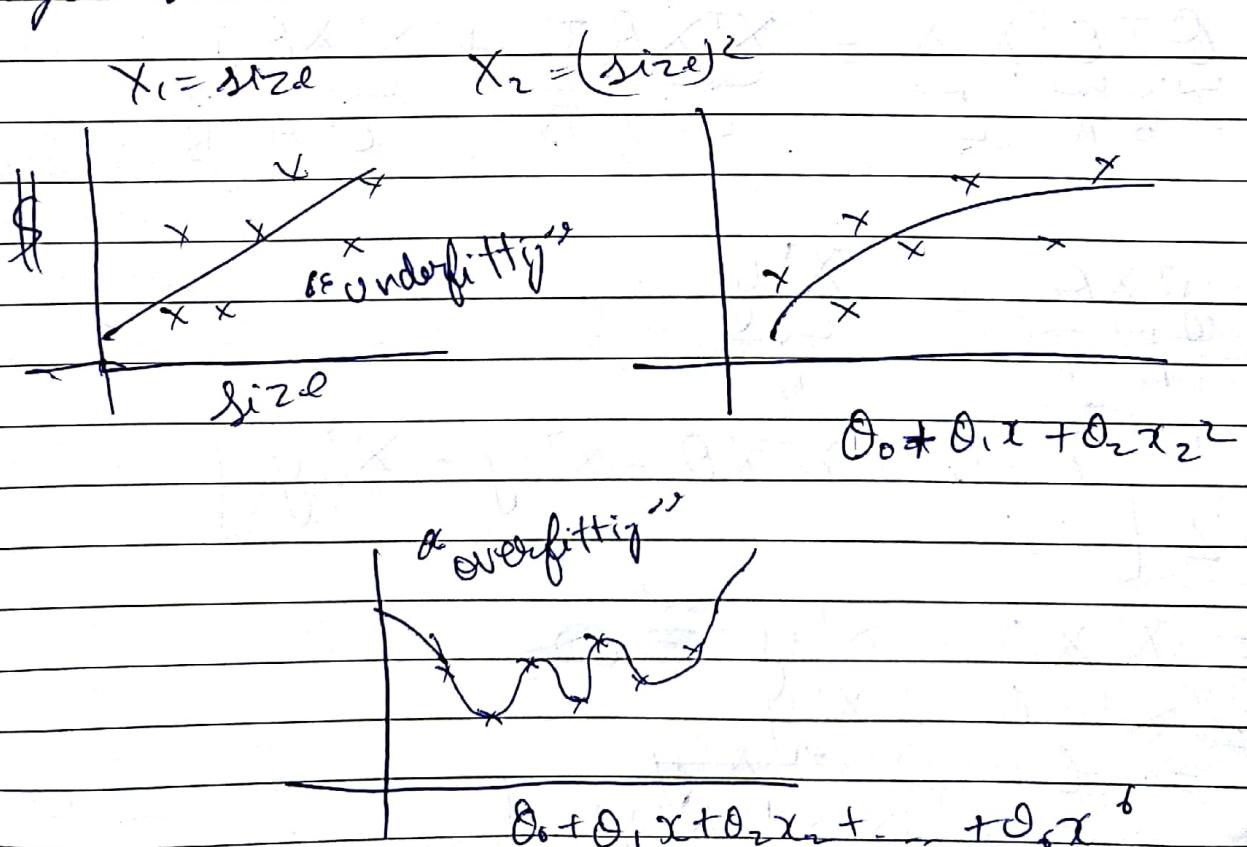
↓ → locally weighted regression
probabilistic interpretation

logistic regression

↓ → Decision Perceptron
Perception

$f(x^{(i)}, y^{(i)})$ - i^{th} training example

Choice of features - has large impact on learning algorithms



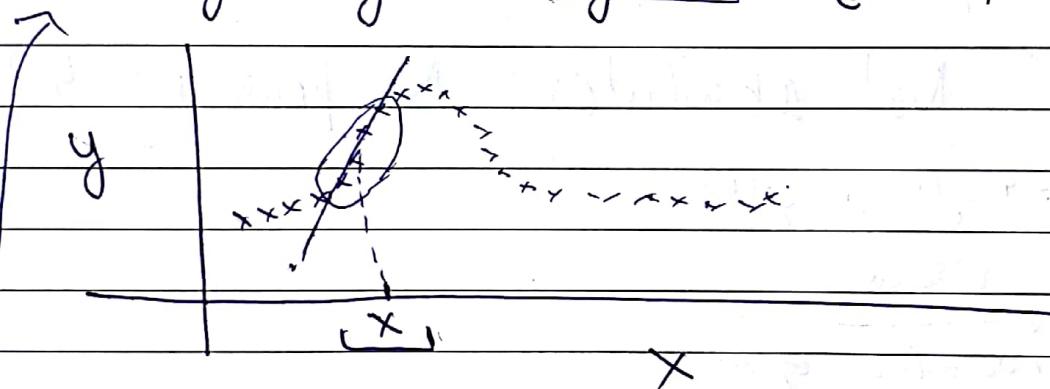
"Parametric" learning algorithm

number of parameters - fixed e.g. linear regression

"Non-parametric" learning algorithm

- no. of parameters grows with m

Locally weighted regression (Loess/Lowers)



To evaluate h at a certain x

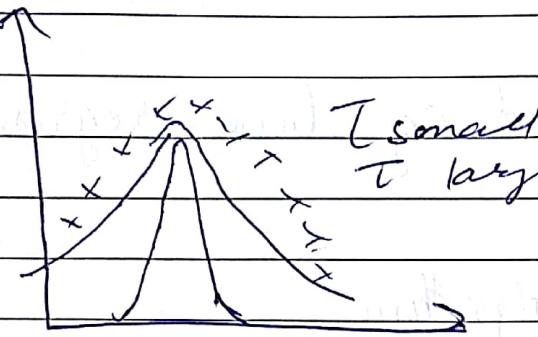
For predicting at x , we consider in vicinity of x only
LWR: Fit θ to minimize

$$\sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

where $w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2T^2}\right)$

If $|x^{(i)} - x|$ small, then $w^{(i)} \approx 1$
If $|x^{(i)} - x|$ large then $w^{(i)} \approx 0$

Don't associate semantic of Gaussian function
Don't associate semantic of Gaussian function



Andrew More on K0 Trees

Probabilistic interpretation

Why squares? Not absolutes? Not powers of 4?

Assume $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$
 $\epsilon^{(i)}$ = error

↳ unmodeled effects
 ↳ random noise

$\epsilon^i \sim \mathcal{N}(0, \sigma^2) \rightarrow$ Normal/Gaussian distributed

$$P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$y^{(i)} | x^{(i)}; \theta \sim \mathcal{N}(\theta^T x, \sigma^2)$$

Why gaussian?

- 1) Mathematically convenient
- 2) Central Limit Theorem

→ Sum of many RVs is approx. gaussian
independent

Assumptions we make are never absolutely true

$\underbrace{P(y^i | x^i; \theta)}$ → Notation :- θ is not a RV.
 $\therefore P(y^i | x^i, \theta)$ is wrong

Frequentist
Viewpt.

$\Sigma^{(i)}$'s are IID - independently identically distributed.

$$L(\theta) = \boxed{P(\bar{y} | X; \theta)}$$

Likelihood emphasis the boxed term as function of θ

Likelihood of parameter \leftrightarrow Probabilities of data

$$= \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

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

Maximum likelihood:
Choose θ to maximize $L(\theta) = P(\vec{y} | \vec{x}; \theta)$

$$l(\theta) = \log(L(\theta))$$

$$= \log \prod_i \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right]$$

$$= \sum_{i=1}^m \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp(-\cdot) \right]$$

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

So maximizing $l(\theta)$ is the same as minimizing

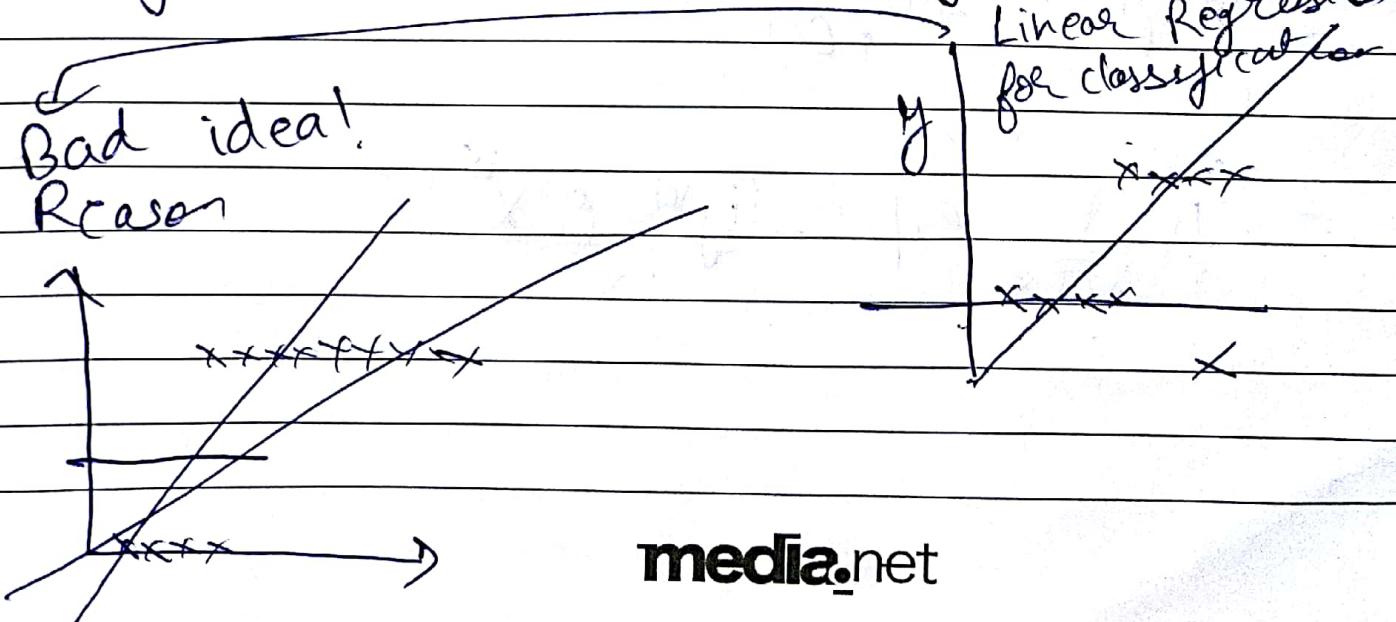
$$\sum_{i=1}^m \frac{(y^{(i)} - \theta^T x)^2}{2} = J(\theta)$$

Classification

target variable is discrete $y \in \{0, 1\}$

Bad idea!

Reason



$$y \in \{0, 1\}$$

$$h_{\theta}(x) \in [0, 1]$$

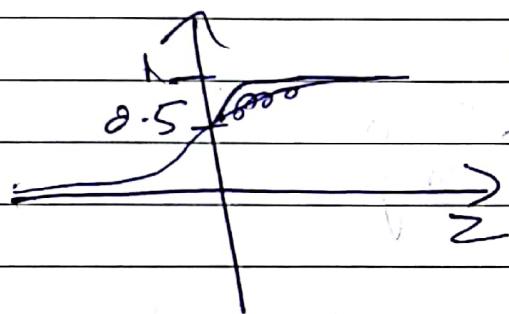
Choose

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

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

$g(z)$ - sigmoid function

logistic function



From where?

1) Generalized Linear Models

2)

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

$$P(y=0|x; \theta) = 1 - h_{\theta}(x)$$

$$P(y|x; \theta) = \underbrace{h_{\theta}(x)^y}_{\text{if } y=1} \underbrace{(1-h_{\theta}(x))^{1-y}}_{\text{if } y=0}$$

$\Rightarrow y=1, h_{\theta}(x)$

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

$$L(\theta) = P(\vec{y} | \vec{x}; \theta) = \prod_i P(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_i h(x^{(i)})^{y^{(i)}} (1 - h(x^{(i)}))^{1-y^{(i)}}$$

$$l(\theta) = \log L(\theta)$$

$$= \sum_{i=1}^m [y^{(i)} \log(h(x^{(i)})) + (1-y^{(i)}) \log(1-h(x^{(i)}))]$$

$$\theta := \theta + \alpha \nabla_{\theta} l(\theta)$$

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

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

Diff. from that of
linear regression

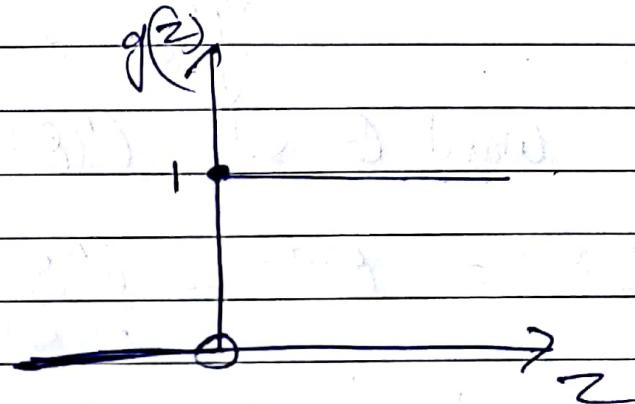
Generalized linear model will show
that even if we are using entirely diff.
model, update eqn. will look cosmetically
similar

Cover up the derivations & prove them
— Andrew Ng

Digression : Reception

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

$$h_{\theta}(x) = g(\theta^T x)$$



$$\theta_j = \theta_j + \alpha(y^{(i)} - h(\theta^T x^{(i)})) x_j^{(i)}$$

Vid 4

Today

Logistic Regression

- Newton's Method

Exponential Family

Generalized Linear Models (GLMs)

Find θ s.t. $f(\theta) = 0$

$$f'(\theta^{(0)}) = f(\theta^{(0)})$$

$$\Delta = \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

