



linear regression

Supervised Learning

θ = # parameters

m = # training examples

n = # features

h = hypothesis function.

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

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

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

$$\underset{\theta}{\text{minimize}} \cdot J(\theta)$$

↓
objective.

→ start with some θ ($\theta = \vec{\theta}$) keep changing

θ to reduce $J(\theta)$
(minimize).

→ ~~gradient~~ gradient descent

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

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_j} (h_\theta(x^{(i)}) - y^{(i)}) \\
 &= (h_\theta(x) - y) \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\
 &= x_j (h_\theta(x) - y).
 \end{aligned}$$

$$\theta_j := \theta_j - \alpha (h_\theta(x) - y) x_j$$

↓
learning rate { magnitude of step
in direction of steepest descent }.

for 'm' training examples

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

Batch Gradient Descent

takes more time to converge to the minimum.

Stochastic Gradient Descent

Repeat of

for $j=1$ to m do

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

(for all i).

converges to θ
which is near
to minimum but
not

exactly
minimum.

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

derivative function of
gradient of
parameters

Gradient Descent

$$\theta := \underbrace{\theta}_{\mathbb{R}^n} - \alpha \underbrace{\nabla_{\theta} J}_{\mathbb{R}^n}$$

Suppose

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

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

$$\rightarrow g f A \in \mathbb{R}^{n \times n}$$

$$\text{tr}(A) = \sum_{i=1}^n A_{ii}$$

$$\boxed{\text{tr}(AB) = \text{tr}(BA)}$$

$$\rightarrow f(A) = \text{tr}(AB)$$

$$\boxed{\nabla_A \text{tr}(AB) = B^T}$$

$$\boxed{\text{tr } A = \text{tr } A^T}$$

$$\text{if } a \in \mathbb{R} \quad \text{tr}(a) = a$$

$$\rightarrow \boxed{\nabla_A \text{tr } ABA^TC = CAB + C^TAB^T}$$

$\rightarrow x \rightarrow \text{design matrix}$

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

$$x_0 = \begin{bmatrix} (x^1)^T \\ \vdots \\ (x^m)^T \end{bmatrix} \quad \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} (x^1)^T 0 \\ \vdots \\ (x^m)^T 0 \end{bmatrix}$$

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

$$x\theta - y = \begin{bmatrix} h(x^1) - y^{(1)} \\ \vdots \\ h(x^m) - y^{(m)} \end{bmatrix} \in \mathbb{R}^m.$$

$$z^T z = \sum_{i=1}^m z_i^2$$

$$\begin{aligned} & \frac{1}{2} (x\theta - y)^T (x\theta - y) \\ &= \frac{1}{2} \sum_{i=1}^m (h(x^i) - y^{(i)})^2 = J(\theta). \end{aligned}$$

$$\rightarrow \nabla_{\theta} J(\theta) = \vec{0}$$

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

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

$$= \frac{1}{2} \nabla_{\theta} [\cancel{\theta^T x \theta} - \theta^T x^T y - y^T \theta + y^T y]$$

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

$$\star \nabla_{\theta} \text{tr} \frac{y^T x \theta}{B A} = x^T y$$

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

$$= \frac{1}{2} [x^T x \theta + x^T x \theta - x^T y - x^T y]$$

$$= x^T x \theta - x^T y = 0$$

$$[x^T x \theta = x^T y] \rightarrow \text{normal equ^n.}$$

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

Proofs (Algebra).

$$\nabla_A \text{tr } AB = \begin{bmatrix} 1 & & & & 1 \\ c_1 & \dots & c_n \\ 1 & & 1 \end{bmatrix} \begin{bmatrix} 1 & & & & 1 \\ b_1 & \dots & b_n \\ 1 & & 1 \end{bmatrix}$$

assume $A \times B$
are $n \times n$
 $A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times n}$

$$a_i \cdot b_j = \begin{bmatrix} a_{11} \cdot b_{11} & a_{12} \cdot b_{12} \\ \vdots & \vdots \\ a_{in} \cdot b_{in} \end{bmatrix} = \begin{bmatrix} a_{11} \cdot b_{11} & \dots & a_{in} \cdot b_{in} \\ a_{12} \cdot b_{12} & \dots & a_{in} \cdot b_{in} \\ \vdots & \ddots & \vdots \\ a_{n1} \cdot b_{n1} & \dots & a_{nn} \cdot b_{nn} \end{bmatrix}$$

$$\nabla_B \text{tr } AB = \begin{bmatrix} b_{11} & & & & b_{n1} \\ b_{12} & \dots & & & b_{n2} \\ \vdots & & \ddots & & \vdots \\ b_{1n} & & & & b_{nn} \\ B^T & & & & \end{bmatrix}$$

\rightarrow to prove

$$\nabla_{A^T} f(A) = (\nabla_A f(A))^T$$

$$A = \begin{bmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & \ddots & \vdots \\ A_{m1} & \cdots & A_{mn} \end{bmatrix}$$

$$A^T = \begin{bmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & \ddots & \vdots \\ A_{m1} & \cdots & A_{mn} \end{bmatrix}$$

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

by defⁿ of
 $\frac{\partial f(A)}{\partial A}$.

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

from observation:- \rightarrow and by transposing

$$\nabla_{A^T} f(A) = (\nabla_A f(A))^T$$

$$\textcircled{1} = [\textcircled{2}]^T$$

$$\frac{\partial}{\partial x_1}$$

$$\nabla_x x^T A x = Ax + A^T x$$

$x \in \mathbb{R}^{n \times 1}$, $A \in \mathbb{R}^{n \times n}$

$$Ax = \begin{bmatrix} \xrightarrow{\quad} & \hat{a}_1^T \\ \vdots & \vdots \\ \xrightarrow{\quad} & \hat{a}_n^T \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} \hat{a}_1^T x \\ \vdots \\ \hat{a}_n^T x \end{bmatrix}$$

$$= \begin{bmatrix} \cancel{\hat{a}_1^T x} \\ \vdots \\ \cancel{\hat{a}_n^T x} \end{bmatrix} = \begin{bmatrix} \hat{a}_1^T x & \dots & \hat{a}_n^T x \end{bmatrix}^T$$

$$\nabla_x^T A x$$

$$= [x_1 \dots x_n] \begin{bmatrix} \hat{a}_1^T x \\ \vdots \\ \hat{a}_n^T x \end{bmatrix}$$

$$= x_1 \hat{a}_1^T x + \dots + x_n \hat{a}_n^T x$$

$$\frac{\partial}{\partial x_i} (x^T A x) = \frac{\partial}{\partial x_i} (x_1 \hat{a}_1^T x + \dots + x_n \hat{a}_n^T x)$$

$$= (\hat{a}_1^T x + x_1 \hat{a}_1^T \frac{\partial x}{\partial x_i} + \dots + x_n \hat{a}_n^T \frac{\partial x}{\partial x_i})$$

$$= \hat{a}_1^T x + (x_1 \hat{a}_{1i} + \dots + x_n \hat{a}_{ni}) \frac{\partial x}{\partial x_i}$$

$$\frac{\partial (iA)}{\partial x_i} = \hat{a}_1^T x + (x_1 a_{1i} + \dots + x_n a_{ni})$$

$$\frac{\partial x}{\partial x_i} = \begin{bmatrix} \frac{\partial x_1}{\partial x_i} \\ \frac{\partial x_2}{\partial x_i} \\ \vdots \\ \frac{\partial x_n}{\partial x_i} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Only for all x_i^0 , $i \geq 2$. And $i \leq n$

$$\frac{\partial}{\partial x} (x^T A x) = \begin{bmatrix} \hat{a}_1^T x + (x_1 a_{11} + \dots + x_n a_{n1}) \\ \hat{a}_2^T x + (x_1 a_{12} + \dots + x_n a_{n2}) \\ \vdots \\ \hat{a}_n^T x + (x_1 a_{1n} + \dots + x_n a_{nn}) \end{bmatrix}$$

$\underbrace{A x}_{Ax}$ $\underbrace{A^T x}_{A^T x}$

$$\frac{\partial}{\partial x} (x^T A x) = Ax + A^T x.$$

$$\boxed{\nabla_x (x^T A x) = Ax + A^T x}.$$

\rightarrow Derivative of a trace product in

$$\nabla_A \text{tr}(AB) = B^T$$

$$\text{tr}(AB) = \text{tr} \begin{bmatrix} \vec{a}_1 \\ \vec{a}_2 \\ \vdots \\ \vec{a}_n \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

$$= \text{tr} \begin{bmatrix} \vec{a}_1^T b_1 & \vec{a}_1^T b_2 & \cdots & \vec{a}_1^T b_m \\ \vec{a}_2^T b_1 & \vec{a}_2^T b_2 & \cdots & \vec{a}_2^T b_m \\ \vdots & \vdots & \ddots & \vdots \\ \vec{a}_n^T b_1 & \vec{a}_n^T b_2 & \cdots & \vec{a}_n^T b_m \end{bmatrix}$$

$$\sum_{i=1}^n a_{ii} b_{ii} + a_{12} b_{12} + \dots + a_{nn} b_{nn}$$

$$\frac{\partial \text{tr}(AB)}{\partial a_{ij}} = b_{ji}$$

$$\nabla_A \text{tr } AB = B^T.$$

Lecture - 3

Linear Regression

↓ locally weighted regression.
probabilistic interpretation

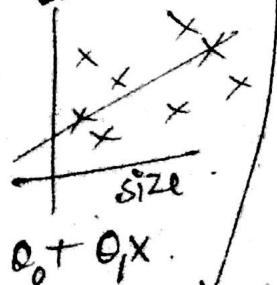
↓ Logistic regression.

Newton's method
↓ regression → perceptron learning alg.

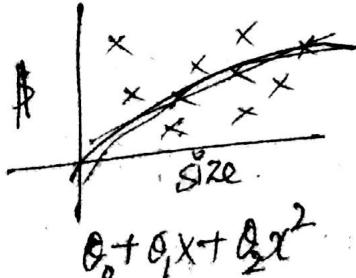
$(x^{(i)}, y^{(i)}) \rightarrow$ i^{th} training example
 $h_0(x^{(i)}) = \sum_{j=0}^n \theta_j x_j^{(i)} = \theta^T x. (x_0 = 1)$ convention.

Housing prices

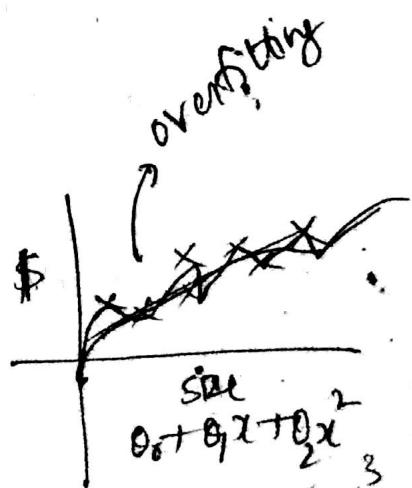
underfitting



patterns in data
that alg. fails
to model



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



algorithm. $\theta_3 x^3$.
models the $\theta_4 x^4$
idiosyncrasies
of this particular
dataset. $\dots \theta_6 x^6$.

→ "Non-parametric" Learning Algorithm

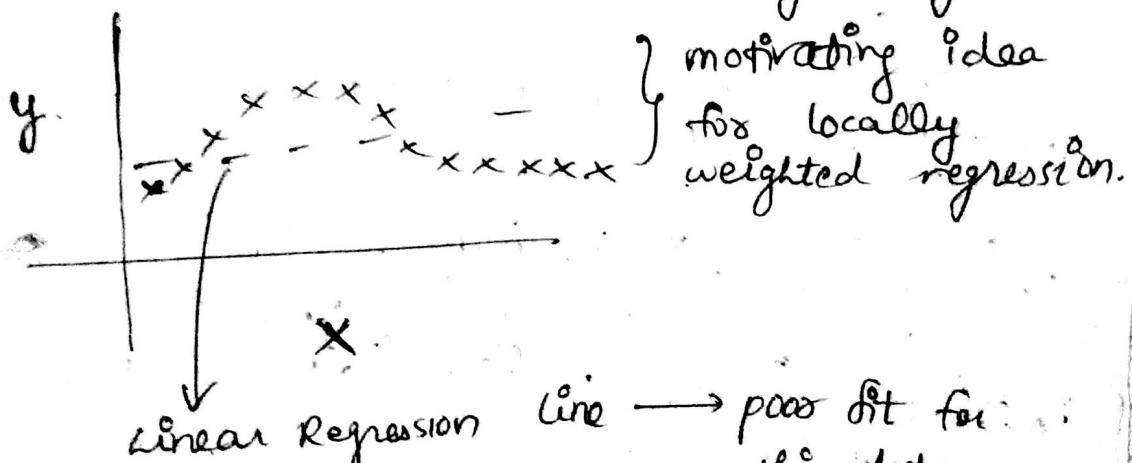
Lin. Regression is one example of parametric learning alg.

→ "Non-parametric" Learning algorithm.

- no. of parameters grows with training examples (m).

Locally weighted regression (Loess).

We worry less about feature engineering.

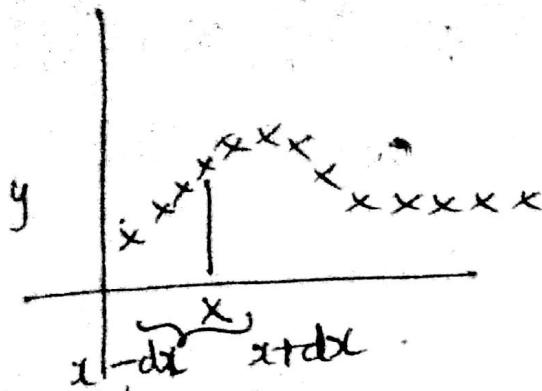


Linear Regression line → poor fit for this data.

* You have to guess experiment with features to fit this data.

* To evaluate h at a certain x .

LR → fit θ to minimise $\sum_{i=0}^T (y^{(i)} - \theta^T x^{(i)})^2$.
and output $\theta^T x^{(i)}$.



In LoWoRo, 1) check vicinity of x as well.

ii) You have a bunch of points, i.e. vicinity of x , then we apply Linear Regression (LoRo) to these points.

iii) You get a LoRo line, evaluate x on the L.R. line.



In LoWoRo:

Fit θ to minimize

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

where $w^{(i)} = e^{-\frac{(x^{(i)} - x)^2}{2\tau^2}}$

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

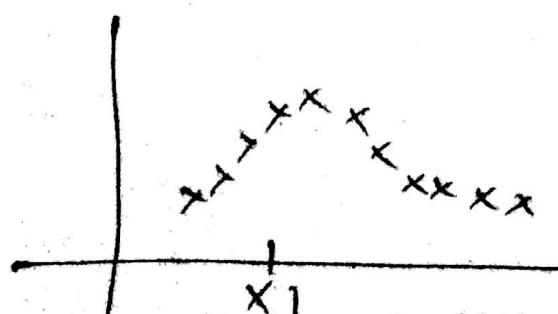
If $|x^{(i)} - x|$ large, then $w^{(i)} \approx 0$.

Not a Normal / Gaussian distribution.

τ_2 bandwidth

parameter

controls how fast weights decrease with distance.



query for x .

→ Then points close to x ($|x^{(i)} - x|$ small). get larger weights than ones further from x ($|x^{(i)} - x|$ large). Points further away from 'x' (query point) would have close to zero contribution.

Research different functions for Locally weighted regression. (Andrew Moore, K-D Trees for LWR)

Probabilistic Interpretation of L.R.

- why minimize least squares and not any other metric?

Assume.

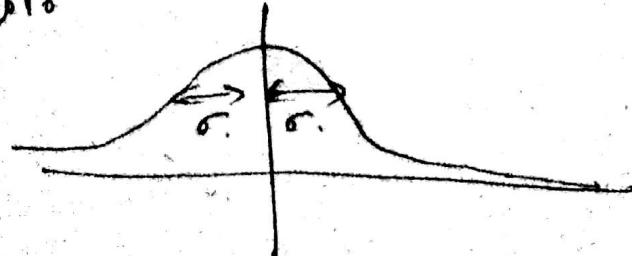
$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$\epsilon^{(i)}$ = error term \approx random noise.

$\epsilon^{(i)} \sim N(0, \sigma^2)$ → Normal distribution.

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

P.D.F.



h
i
i
o
n
a

→
 ϵ

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

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

$$\mu = \theta^T x^{(i)}$$

$$\sigma^2 = \sigma^2$$

Q) Why is error gaussian?

- i) mathematical convenience
- ii) distribution of errors is more often Gaussian.
- iii) Central limit theorem.

* We are also assuming that error has zero mean.

$y^{(i)} | x^{(i)}; \theta \rightarrow \theta$ is not a random variable
in this case.

We don't know
what θ is, but
 θ is also not a
random variable.

→ Read as:- $y^{(i)}$ given $x^{(i)}$ parametrized by θ .

→ $\epsilon^{(i)}$ are iid, independent and identically distributed.

$$\rightarrow L(\theta) = P(\vec{y} | \vec{x}; \theta)$$

of likelihood of θ .

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

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

over all training samples.

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

likelihood means we are viewing $P(\vec{y} | \vec{x}; \theta)$ as a function of θ .

→ Maximum Likelihood Estimation -

* Choose θ which maximises ~~probabil~~ the likelihood i.e. which makes the data as much as possible.

i.e. choose θ to maximise $L(\theta) = P(\vec{y} | \vec{x}; \theta)$

$$l(\theta) = \log L(\theta) = \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

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

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

so, maximising $l(\theta)$ is the same

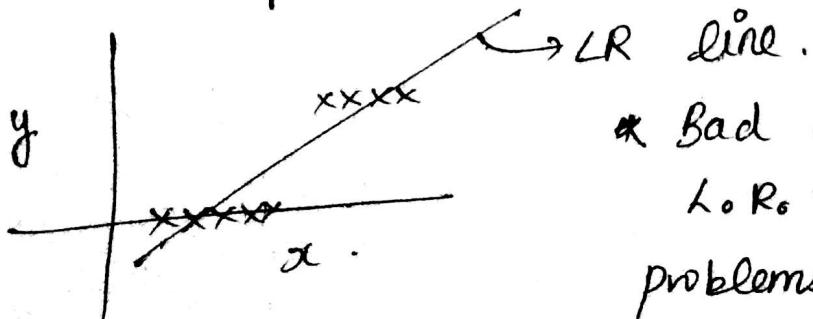
as minimising $\sum_{i=1}^m \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} = J(\theta)$.

→ So Least Squares regression assuming errors are i.i.d Gaussians is equivalent to M.L.E.

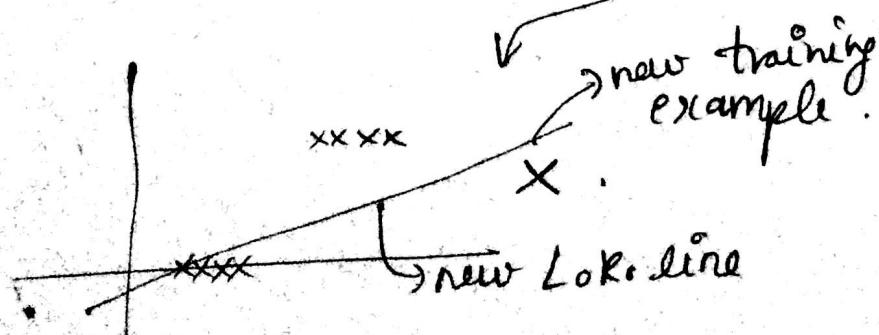
Classification

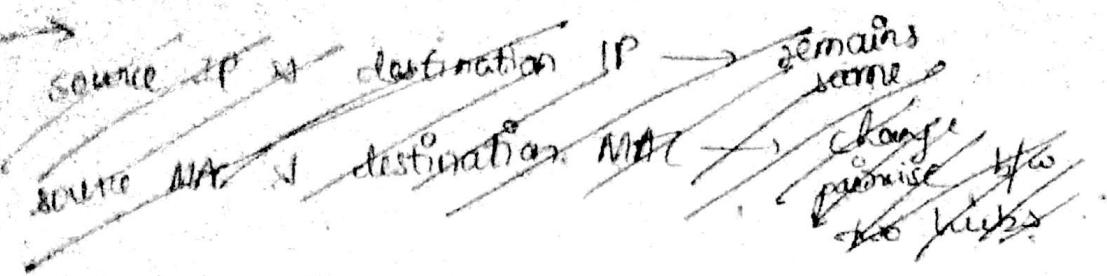
$$y \in \{0, 1\}$$

↳ to be predicted, discrete values.



* Bad idea to apply L.R. to classification problems because,



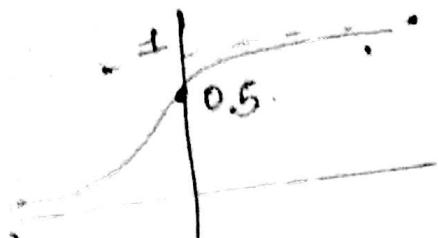


$$y \in \{0, 1\}$$

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

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

$$g(z) = \frac{1}{1 + e^{-z}} \rightarrow \text{sigmoid function}$$



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

Think of my hypothesis as estimating probability that $y=1$.

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

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

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

$$= \prod_{i=1}^m (h_{\theta}(x^{(i)}))^y (1 - h_{\theta}(x^{(i)}))^{1-y}$$

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

$$= \sum_{i=1}^m y^{(i)} (\log h_\theta(x^{(i)})) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))$$

maximise

find θ by gradient descent

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

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

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

different from one in least squares regression.

not a linear function

→ Perception

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

$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

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

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