## ) Linear Regression

Simply put we have data and we try to learn a line which passes through it in best fit possible.

### Supervised Learning :-

| Testing Set |
|---|

↓

Learning Algo ← Input

↓ Output

IP-→ | function h (hypothesis) | → Predicts → Price of house
Size of house

→ When designing a learning algo we need to ask how do you represent h?

In Linear Regression :-

$$h(x) = \theta_0 + \theta_1 x$$

We input size x and it outputs a function which is a linear combination of size x.

More generally :-

$x_1 = $ Size , $x_2 = $ No. of bedrooms

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

In other words :-

$$h(x) = \sum_{j=0}^{2} \theta_j x_j$$

where $x_0 = 1$

---

Here,

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$\theta \rightarrow$ Parameters

$M = $ # Training samples

$x = $ "Inputs" / feature

$y = $ "Output" / Target Variable

$(x,y) \Rightarrow$ One training example

$(x^i, y^i) \Rightarrow$ ith training example

$n \Rightarrow$ # features

→ Choose $\theta$ such that $h(x) \approx y$ for training examples.

$h_\theta(x)$ depicts that h depends both on params and input features x

In Linear Regression algos aim is to minimise the least squares error :-

$$\text{Minimise} ( h_\theta(x) - y)^2$$

Basically we need to choose values of $\theta$ that minimises least squares error

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

we add $\frac{1}{2}$ to make math simpler.

$$\therefore \text{Minimise} (J(\theta))$$
$$\theta$$

### * Gradient Decent :-

① Start with some $\theta$ (say $\theta = \vec{0}$)

Keep changing $\theta$ to minimise $J(\theta)$

One step of gradient decent is implemented as follows:-

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

Learning Rate

:= depicts assignment

Some math:-

$$\boxed{\frac{\partial J(\theta)}{\partial \theta_j}} = \frac{\partial}{\partial \theta_j} \left( \frac{1}{2} \left( h_\theta(x) - y \right)^2 \right)$$

$$= 2 \times \frac{1}{2} \frac{\partial}{\partial \theta_j} \left( h_\theta(x) - y \right)$$

$$= \left( h_\theta(x) - y \right) \frac{\partial}{\partial \theta_j} \left( \theta_0 x_0 + \theta_1 x_1 + \ldots + \theta_n x_n - y \right)$$

$$= \boxed{\left( h_\theta(x) - y \right) x_j}$$

Hence one step of gradient decent is the following:-

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

∴ The gradient decent algorithm is to Repeat eq ① until convergance. for

$$j = \{1, 2, \ldots, n\}$$

where n is the no. of features.

We find $\frac{\partial}{\partial \theta_j} J(\theta)$ by summing RHS over m — Thats why it is also called Batch Gradient decent.

Efficient algo than this is stochastic gradient decent.

Repeat {
   for j=1 to m {                    for every j
      $$\theta_j := \theta_j - \alpha \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)}$$
   }
}

$$\nabla_\theta J(\theta) = \begin{pmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{pmatrix}$$

Hence shorter way of computing global minima which is computed in gradient decent is by taken derivative of $J(\theta)$ wrt $\partial\theta$ and equating it to $\partial$.

J is a function which maps vector $\theta$ to a real value.

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

→ If A is n×n matrix
   tr A = Sum of primary diagonal elements

$$tr A = tr A^T$$

If $f(A) = tr AB$
Then
$$\nabla_A f(A) = B^T$$

$$\Rightarrow tr AB = tr BA$$

$$\Rightarrow tr ABC = tr CAB$$

$$\nabla_A tr A A^T C = CA + C^T A$$

Now,

we know that,

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

$$X\theta = \begin{bmatrix} - (x^1)^T - \\ - (x^2)^T - \\ \vdots \\ - (x^m)^T - \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} x^{1^T}\theta \\ x^{2^T}\theta \\ \vdots \\ x^{m^T}\theta \end{bmatrix}$$

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

Also,

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

Then,

$$J(\theta) = \frac{1}{2} \left( X\theta - y \right)^T \left( X\theta - y \right)$$

Hence,

$$\nabla_\theta J(\theta) = \nabla_\theta \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

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

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

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

$$= x^T x \theta - x^T y \overset{set}{=} \vec{0}$$

$$\Rightarrow x^T x \theta = x^T y \quad \text{"Normal eq"}$$

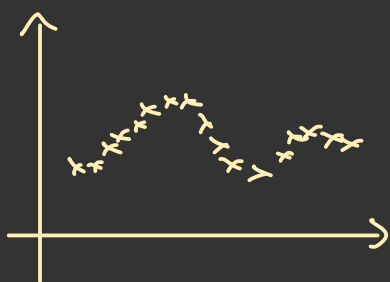and,

$$\theta = (x^T x)^{-1} x^T y \quad - ②$$

Hence, if we compute ② then we get the value of $\theta$ that corresponds to global minimum in one single step.

# Lecture - 3

## Locally weighted Regression :-

Note that - Parametric learning algos - Fit fixed set of parameters ($\theta_i$) to data.

Locally weighted regression is a 'non-parametric' learning algo as data/params you need to keep grows linearly with the size of data.



To evaluate $h$ at certain $x$:

Linear Reg: Fit $\theta$ to minimise

$$\frac{1}{2} \sum_i \left( y^{(i)} - \theta^T x^{(i)} \right)^2$$

Return $\theta^T x$

In Locally weighted regression :-

Fit $\theta$ to minimise

$$\sum_{i=1}^{m} \omega^{(i)} \left( y^{(i)} - \theta^T x^{(i)} \right)^2$$

where $\omega^{(i)}$ is a weight function

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

If $|x^{(i)} - x|$ is small, $\omega^{(i)} \approx 1$

$x \to$ location where you want to make a prediction

$x^{(i)} \to$ Input $x$ for your $i$th training example.

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

$\to$ We use a hyperparameter $\tau$ which is the bandwidth and is used to define the width of the neighbourhood we are using.
So,

$$\omega^{(i)} = \exp\left( -\frac{(x^{(i)} - x)^2}{2\tau^2} \right)$$

## *Probabilistic interpretation of linear regression :-

### Why least squares?

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

$$\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

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

We make the assumption that the $\varepsilon$ error terms are Independently and Identically distributed.

This implies that :-

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

$\to$ Read as - parametrised by

This means that

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

$$\mathcal{L}(\theta) = p(\bar{y}|x;\theta)$$

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

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

$\mathcal{L}(\theta) \rightarrow$ Likelihood of $\theta$

## Log likelihood :-

$$\ell(\theta) = \log \mathcal{L}(\theta)$$

$$= \log \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi}\sigma} \exp(\cdots)$$

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

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

## Maximum likelyhood Estimation:-

## (MLE)

Choose $\theta$ to maximise $\mathcal{L}(\theta)$
i.e choose $\theta$ to minimise:-

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

## Classification Problem :-

## Binary Classification:-

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



## Logistic Regression

Most commonly used classification algorithm

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

---

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

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

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

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

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

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

So, likelihood function now becomes:-

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

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

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

and,

$$\ell(\theta) = \log \mathcal{L}(\theta)$$

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

Choose $\theta$ to maximise $\ell(\theta)$,

Batch gradient Ascent:-

$$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} \ell(\theta)$$

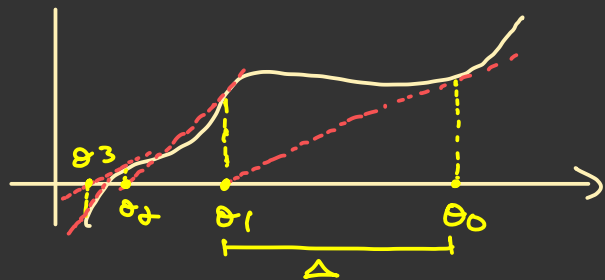$$\theta_j := \theta_j + \alpha \sum_{i=1}^{m} \left(y^i - h_\theta(x^i)\right) x_j^i$$

## Newton's Method :-

Much faster than gradient ascent for optimising the value of $\theta$.

We have $f$

We want to find $\theta$ such that:-
$$f(\theta) = 0$$

[ want to maximise $L(\theta)$
i.e. want $l'(\theta) = 0$ ]



$$\theta_1 = \theta_0 - \Delta$$

$$f'(\theta_0) = \frac{f(\theta_0)}{\Delta}$$

$$\therefore \Delta = \frac{f(\theta_0)}{f'(\theta_0)}$$

$$\therefore \theta_{t+1} = \theta_t - \frac{f(\theta_t)}{f'(\theta_t)}$$

Now,
Let $f(\theta) = l'(\theta)$

So,
$$\theta_{(t+1)} = \theta_t - \frac{l'(\theta_t)}{l''(\theta_t)}$$

when $\theta$ is a vector

$$\theta_{t+1} = \theta_t + H^{-1} \underbrace{\nabla_\theta l}_{\to \text{vector } \mathbb{R}^{n+1}}$$

where $H$ is Hessian matrix

$$H^{-1} \to \mathbb{R}^{n+1 \times n+1}$$

$$H_{ij} = \frac{\partial^2 l}{\partial \theta_i \partial \theta_j}$$

---

We saw that logistic regression used sigmoid function :-

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

The perceptron function is the hard version of sigmoid :-

$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$
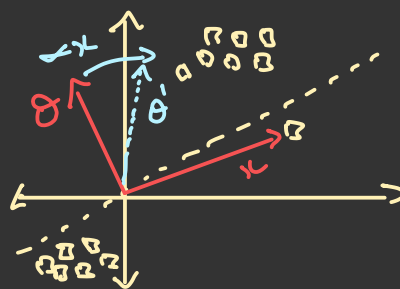
$$h_\theta(x) = g(\theta^T x)$$

The gradient ascent is the same but the $h_\theta(x)$ changes in perceptron function:-

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

Note that, $(y^i - h_\theta(x)^i)$ is a scalar. So, the values will be

$0 \to$ algo got it right

$\pm 1 \to 1$ if wrong $y^i = 1$
$\quad\quad -1$ if wrong $y^i = 0$



$\theta \approx x | y = 1$
$\theta \not\approx x | y = 0$

## Exponential families:-

The one whose PDF:-

$$P(y; \eta) = b(y) \exp\left[\eta^T T(y) - a(\eta)\right]$$

$y$ — data

$\eta$ — natural parameter

$T(y)$ — sufficient statistics

$b(y)$ - Base measure

$a(\eta)$ - log-partition

Note that the dimensions of $\eta$ and $T(y)$ has to match

\*) Bernoulli dist (To map binary data)

$\phi$ = Probability of event

$P(y; \phi) = \phi^{y}(1-\phi)^{(1-y)}$

$= \exp\left[\log\left(\phi^{y}(1-\phi)^{(1-y)}\right)\right]$

$= \exp\left[\log\left(\frac{\phi}{1-\phi}\right)y + \log(1-\phi)\right]$

Here,
$b(y) = 1$

$\eta = \log\left(\frac{\phi}{1-\phi}\right) \Rightarrow \phi = \frac{1}{1+e^{-\eta}}$

$T(y) = y$

$a(\eta) = -\log(1-\phi)$

$= -\log\left(1 - \frac{1}{1+e^{-\eta}}\right) = \log(1+e^{\eta})$

This proves that Bernoulli is a member of exponential family

\* <u>Gaussian Dist (with fix variance)</u>

Assume $\sigma^2 = 1$

$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right)$

$= \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} \exp\left(\mu y - \frac{1}{2}\mu^2\right)$

Say $b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-y^2}{2}\right)$

$T(y) = y, \quad \eta = \mu$

$a(\eta) = \frac{\mu^2}{2} = \frac{\eta^2}{2}$

This proves that Gaussian is also a member of exponential family.

<u>Properties of **Exponential** families:-</u>

(a) MLE wrt $\eta \Rightarrow$ concave

   negative log likelihood $\Rightarrow$ convex

(b) $E[y; \eta] = \frac{\partial a(\eta)}{\partial \eta}$

(c) $V[y; \eta] = \frac{\partial^2}{\partial \eta^2} a(\eta)$
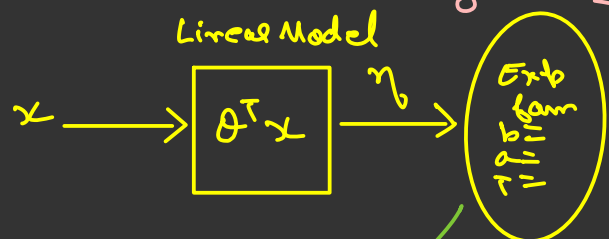
<u>GLM :-</u>

< Generalized Linear Models >

<u>Assumptions / Design Choices :-</u>

(i) $y|x; \theta \sim$ Exponential Family

(ii) $\eta = \theta^T x \qquad \theta \in R^n, x \in R^n$

(iii) Test time: Output $E[y|x;\theta]$
   $\Rightarrow h_{\theta}(x) = E[y|x;\theta]$

Linear Model

$x \longrightarrow \boxed{\theta^T x} \xrightarrow{\eta} \bigcirc$ Exp fam $b = $ $a = $ $T = $

Distribution

$E[y|\eta] = E[y|\theta^T x]$

$= h_{\theta}(x)$

It is very clear here that we are training $\theta$ to predict the params of exponential family dist. whose mean is the prediction we are going to make for $y$.

All this is during Test time.
During Train time :-

The program we are learning
is θ using gradient descent.

During learning we perform the
MLE over θ :

$$\text{Max } \log P\left(y^{(i)}; \theta^T x^{(i)}\right)$$

## GLM Training :-

Learning update Rule :-

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

We can also use Newton's method
to learn θ, as long as the
dimensionality of our features
is less than a few thousand.

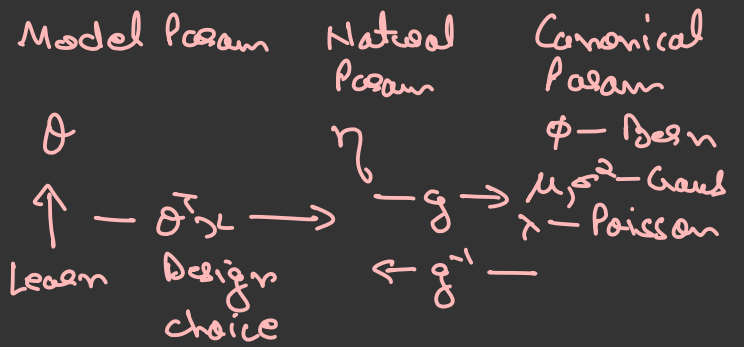| Data Type | Dist. used to model |
|---|---|
| Real | Gaussian |
| Binary | Bernoulli |
| Positive Integers | Poisson |
| $R^+$ | Gamma, Exponential |

## Terminology :-

η - Natural parameter

$$\mu = E[y; \eta] = g(\eta) \Rightarrow \text{Canonical Response function}$$

$$\eta = g^-(\mu) \Rightarrow \text{Canonical Link function}$$

$$g(\eta) = \frac{\partial a(\eta)}{\partial \eta}$$

## We have 3- parameterizations :-

| Model Param | Natural Param | Canonical Param |
|---|---|---|
| θ | η | φ - Bern |

$$\uparrow - \theta^T x \longrightarrow \quad -g \to \quad \mu, \sigma^2 - \text{Gauss}$$
$$\lambda - \text{Poisson}$$

Learn   Design          $\leftarrow g^{-1} -$
        choice

So, now it becomes evident that,
for example, if we use both
to be the dist. of our output
then the logistic regression
becomes :-

$$h_\theta(x) = E[y | x; \theta]$$

$$= \phi$$

$$= \frac{1}{1 + e^{-\eta}}$$

$$= \frac{1}{1 + e^{-\theta^T x}} \left.\begin{array}{c} \end{array}\right\} \text{Sigmoid func.}$$

## Softmax Regression :-

### Cross Entropy interpretation



Here we are talking about
multi class classification.

Goal - Learn a model which
when given a new datapoint
can make a prediction
to which class the new
point belongs to.

K - # classes

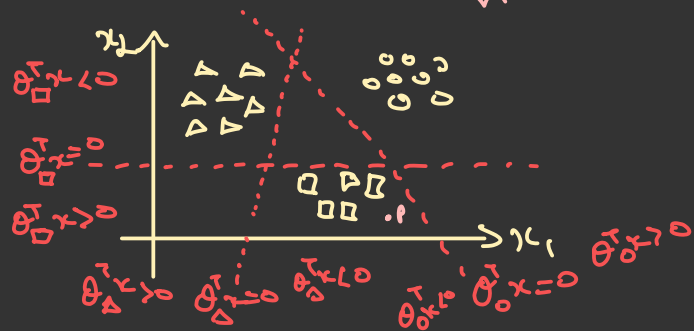$x^{(i)} \in \mathbb{R}^n$

Label $y = [\{0,1\}^K]$ eg $[0,0,1,0]$

$\theta_{class} \in \mathbb{R}^n$

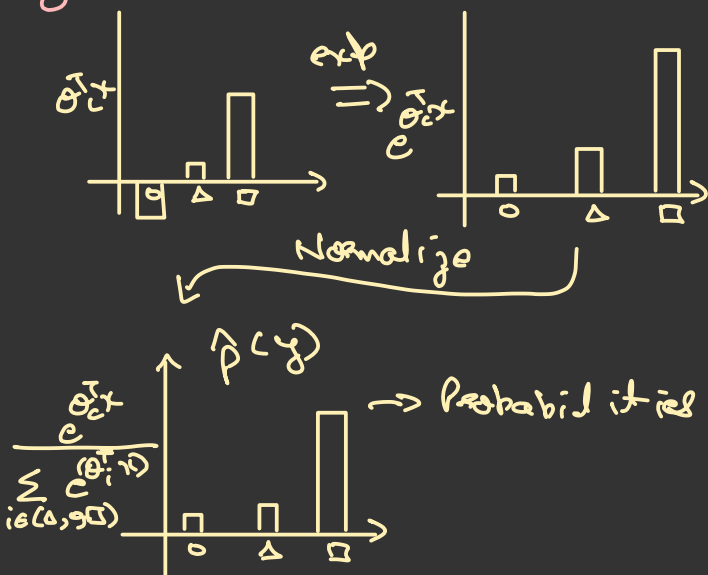class $\in \{\triangle, \circ, \square \cdots\}$

We have K such $\theta s$, one for each class

We can also represent it as a matrix :-

$$K \begin{bmatrix} — \theta_{c_1} — \\ — \theta_{c_2} — \\ \vdots \end{bmatrix}$$

$n$



So, for a new point P we might get :-



Normalize



$\rightarrow$ Probabilities

Let's assume that the new point is $\triangle$. In that case P(y) would look like :-

---



$p(y)$

So, the learning approach we need to do is to minimize the distance between $\hat{p}(y)$ and $p(y)$

Technically the term is to minimize the cross entropy between $\hat{p}(y)$ and $p(y)$

$$CrossEnt(P, \hat{P}) = \sum_{y \in \{0, \triangle, \square\}} p(y) \log \hat{p}(y)$$

$$= -\log \hat{p}(y_\triangle)$$

$$= -\log \frac{e^{\theta_\triangle^T x}}{\sum_{c \in \{0, \triangle, \square\}} e^{\theta_c^T x}} \quad — Ⓐ$$

we treat Ⓐ as loss and do gradient decent wrt the parameters.

# Lecture - 5

All algos studied so far are called discriminative learning algos.

Now we will look into Discriminative learning algos.

A discriminative learning algo learns $p(y/x)$

or, learns $h_\theta(x) = \begin{cases} 0 \\ 1 \end{cases}$ directly

where $h_\theta(x)$ is a mapping from $x$ to $y$

A generative learning algo learns

$p(x | y)$
$\quad\quad\searrow$ class
feature

It also learns $p(y)$ which is called class prior.

So by using Bayes rule :-

$$p(y=1|x) = \frac{p(x|y=1)\ p(y=1)}{p(x)}$$

where,

$$p(x) = p(x|y=1)\,p(y=1) +$$
$$p(x|y=0)\,p(y=0)$$

## *Gaussian Discriminant Analysis :-

Suppose $x \in \mathbb{R}^n$ (drop $x_0 = 1$ convension)

Assume $p(x|y)$ is Gaussian

$z \sim \mathcal{N}(\vec{\mu}, \Sigma) \quad\quad z \in \mathbb{R}^n$
$\quad\quad \searrow \mathbb{R}^n \searrow \mathbb{R}^{n\times n}$

$E[z] = \mu$

$$Cov(z) = E\left[(z-\mu)(z-\mu)^T\right]$$
$$= Ezz^T - (Ez)(Ez)^T$$

I am writing $E[z]$ as $Ez$ for easier writing

$$p(z) = \frac{1}{(2\pi)^{n/2}\,|\Sigma|^{1/2}}\exp\left(-\frac{1}{2}(x-\mu)^T\,\Sigma^{-1}(x-\mu)\right)$$

## GDA Model:

$$p(x|y=0) = \frac{1}{(2\pi)^{n/2}\,|\Sigma|^{1/2}}\exp\left(-\frac{1}{2}(x-\mu_0)^T\,\Sigma^{-1}(x-\mu_0)\right)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{n/2}\,|\Sigma|^{1/2}}\exp\left(-\frac{1}{2}(x-\mu_1)^T\,\Sigma^{-1}(x-\mu_1)\right)$$

Parameters are - $\mu_0, \mu_1, \underset{\mathbb{R}^{n\times n}}{\overset{\mathbb{R}^n}{\Sigma}}, \phi - \mathbb{R}$

Note that we are using same cov matrix $\Sigma$ for both classes but diff $\mu - \mu_1$ and $\mu_0$

$$p(y) = \phi^y(1-\phi)^{1-y}$$

$y$ is a bern random variable as it takes values 0 and 1. That's why we have used bern dist^d for $p(y)$

## * How to fit the parameters?

Training set - $\{x^i, y^i\}_{i=1}^m$

In Generative algos, we define Joint likelihood where,

$$\mathcal{L}(\phi, \mu_1, \mu_2, \Sigma) = \prod_{i=1}^m p(x^i, y^i; \phi, \mu_1, \mu_2, \Sigma)$$
$$= \prod_{i=1}^m p(x^i|y^i)\,p(y^i)$$

Note that here we are maximizing the joint likelihood.

In case of discriminative algo we maxing the conditional likelihood

$$\mathcal{L}(\theta) = \prod_{i=1}^{m} p(y^i \mid x^i; \theta)$$

Now in order to learn the params which maximises $\ell(\phi, \mu_1, \mu_0, \Sigma)$ we use max likelihood estimation we get :—

$$\phi = \frac{\sum_{i=1}^{m} y^{(i)}}{m} = \frac{\sum_{i=1}^{m} 1\{y^{(i)} = 1\}}{m}$$

$1$ is called indicator notation where,

$$1\{true\} = 1$$
$$1\{false\} = 0$$

$$\mu_0 = \frac{\sum_{i=1}^{m} 1\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^{m} 1\{y^{(i)} = 0\}}$$

Numerator is the sum of feature vectors for examples with $y = 0$

Denominator is no. of examples with $y = 0$

Similarly for $\mu_1$

$$\Sigma = \frac{1}{m} \sum_{i=1}^{m} (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

How to make prediction :—

$$\arg\max_y p(y \mid x) = \arg\max_y \frac{p(x \mid y) p(y)}{p(x)}$$

→ what is the argmin / argmax notation?

$$\min (2-5)^2 = 0$$
$$\arg\min (2-5)^2 = 5 \quad \text{as } 5 \text{ is}$$

the value of $z$ we need to plug in to get $\min (z-5)^2 = 0$

Same logic for argmax

$$\arg\max_y p(x \mid y) = \arg\max_y \frac{p(x \mid y)}{p(y)}$$

Hence, we need to output the value of $y$ which maximises $p(y \mid x)$

→ If we plot $p(x \mid y = 0)$ and $p(x \mid y = 1)$ we will get 2 gaussian curves for all $x^{(i)}$

If we try to plot $p(y = 1 \mid x)$ for all $x^{(i)}$ we will see that it gives sigmoid function plot similar to the logistic linear regression.

## GDA and Logistic LR side-to-side :—

GDA assumes
$$\left.\begin{array}{l} x \mid y = 0 \sim \mathcal{N}(\mu_0, \Sigma) \\ x \mid y = 1 \sim \mathcal{N}(\mu_1, \Sigma) \\ y \sim Bern(\phi) \end{array}\right\} \Rightarrow$$

Logistic Regr.
$$P(y = 1 \mid x) = \frac{1}{1 + e^{-\theta^T x}}$$
$$(i = 1^m)$$
logistic

GDA model
∴ Stronger assumption

Logistic LR model weaker assumption

$\neq$

If
$$\left.\begin{array}{l} x \mid y = 1 \sim Poisson(\lambda_1) \\ x \mid y = 2 \sim Poisson(\lambda_0) \\ y \sim Bern(\phi) \end{array}\right\} \Rightarrow P(y = 1 \mid x) \text{ is logistic}$$

# Naive Bayes :-

Another generative learning algo
( Segway into Natural language
Processing !!)

<u>Step 1</u> -> Take your sentence
and first map it to
the feature vector x

Let's say I make a dictionary
which contain top 10,000
words which appear in my
email.
Then if I get a new email
then I can map it to
a feature vector x where
$x \in R^n$ (n=10,000)
$x[i]=1$ if email contains
ith word of dictionary. (and
$x[i]=0$ otherwise.

We want to model $p(x|y), p(y)$

$2^{10,000}$ possible values of x!

So we need $2^{10,000}-1$
params to learn!!!

Assume $x_i$'s are conditionally
independent given y.

$p(x_1, ..., x_{10000}|y) =$

$\quad p(x_1|y)p(x_2|x_1,y)p(x_3|x_1,x_2,y)$
$\quad\quad ---p(x_{10000}|x_1...x_{9999},y)$

$\quad \overset{assume}{=} p(x_1|y) p(x_2|y) p(x_3|y)$
$\quad\quad --- p(x_{10000}|y)$

$\quad = \prod_{i=1}^{n} p(x_i|y)$

Parameters of this model are:-

$\phi_{j|y=1} = p(x_j=1|y=1)$

$\phi_{j|y=0} = p(x_j=0|y=1)$

$\phi_y = p(y=1)$

Joint likelyhood :

$L(\phi_y, \phi_{j|y}) = \prod_{i=1}^{m} p(x^{(i)}, y^{(i)}, \phi_y, \phi_{j|y})$

MLE :

$\phi_y = \sum_{i=1}^{m} \frac{1\{y^{(i)}=1\}}{m}$

$\phi_{j|y=1} = \frac{\sum_{i=1}^{m} 1\{x_{j=1}^{(i)}, y^{(i)}=1\}}{\sum_{i=1}^{m} 1\{y^{(i)}=1\}}$