

Topics to cover.

- 1) Gaussian function.
- 2) Bayes's Theorem.
- 3) Naive Bayes Classifier.
- 4) Non-Linearity
- 5) Logistic Regression

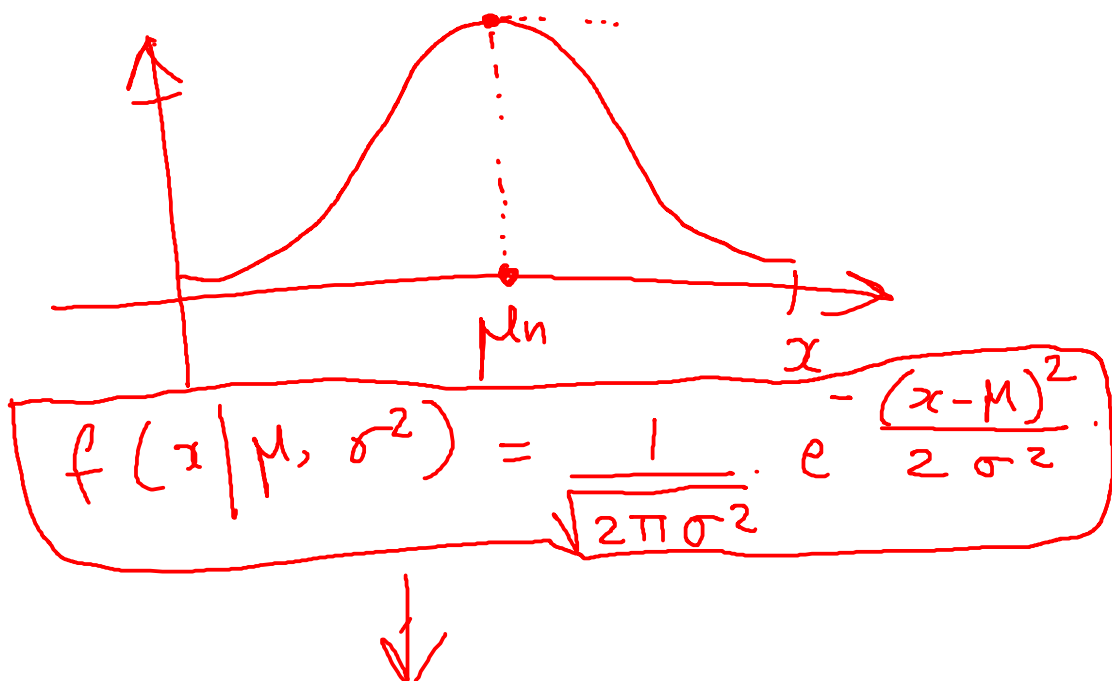
Gaussian Function.

1) Classification :-

Inputs = Features.

Outputs = Labels.

Gaussians



Relative Spread :-

if a r.v. is extracted from a sample that follows this curve, there is high probability that it belongs to the "middle-band".

If we claim that the input data follows a Gaussian curve what parameters do we need to define the data?

- 1) $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \rightarrow$ data points.
- 2) μ, σ values.

Now, how to find the μ and σ ??

For any candidate parameter, we can associate Likelihood function \rightarrow "support" provided by the input data for the given parameter.

Likelihood function is a Joint PMF/PDF.



So to find the "right" parameter, it needs to be the maxima of the Likelihood function. WHAT TO DO??

MLE (Max. Likelihood Estimation).

Let us say uni-variate Gaussian

$$Y = \{y_1, y_2, \dots, y_n\}.$$

then, the maximum likelihood estimates for μ and σ^2 are...

$$\mu = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$$

$$\bar{y} = \mu.$$

Naive Bayes Classifier.

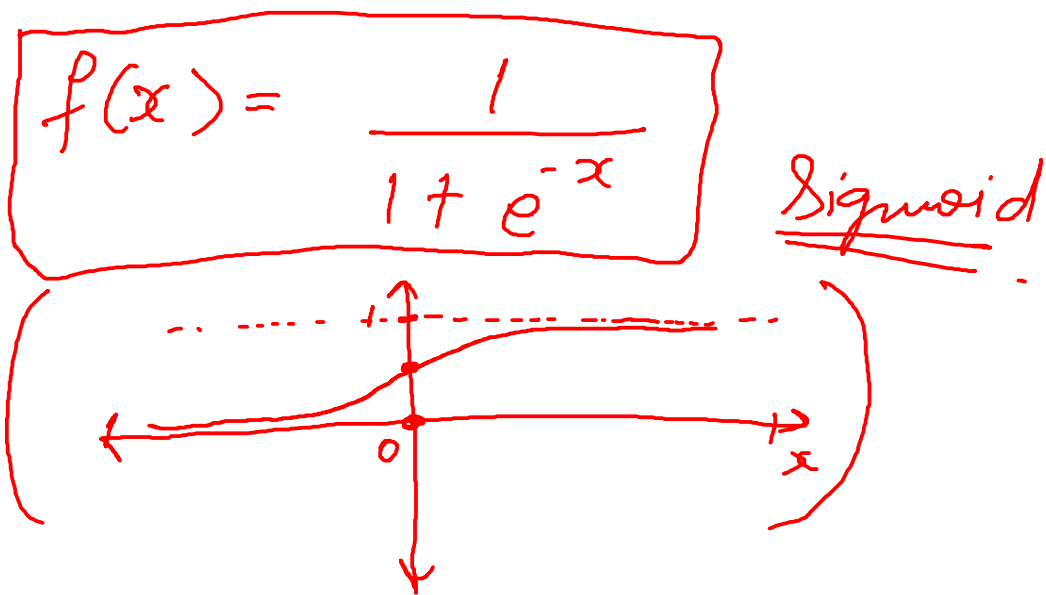
- Let's say input vector is x .
- Set of K classes C_1, C_2, \dots, C_K are there.
- We need to find K where $P(C_K|x)$ is maximum.

(class K has highest p. of accommodating x).

- Let's assume all features are independent of each other.
- Let's assume that each feature follows a Gaussian.

$$P(\text{class}|\text{data}) = \frac{P(\text{data}|\text{class}) \cdot P(\text{class})}{P(\text{data})}$$

Logistic Regression



- Just like Linear Regression; we get a line.
- Now, what we do is we feed this line into the Sigmoid function to get a value between 0 and 1.

Then we threshold the value to particular class (0/1).

$$y = b_0 + b_1 x \leftarrow \text{Linear Model.}$$

$$P = \frac{1}{1 + e^{-(b_0 + b_1 x)}} \left\} \text{logistic Model.}$$

Q. How is this approach different from the Naive Bayes Model?

→ We are directly calculating $p(y|x)$

* Logistic Regression (Mustafa)

- 1) Logistic Idea + Linear Regression.
- 2) Discriminative / Generative Model.
- 3) Parameters - feature weights.
- 4) Estimation - Max. Likelihood Esti...

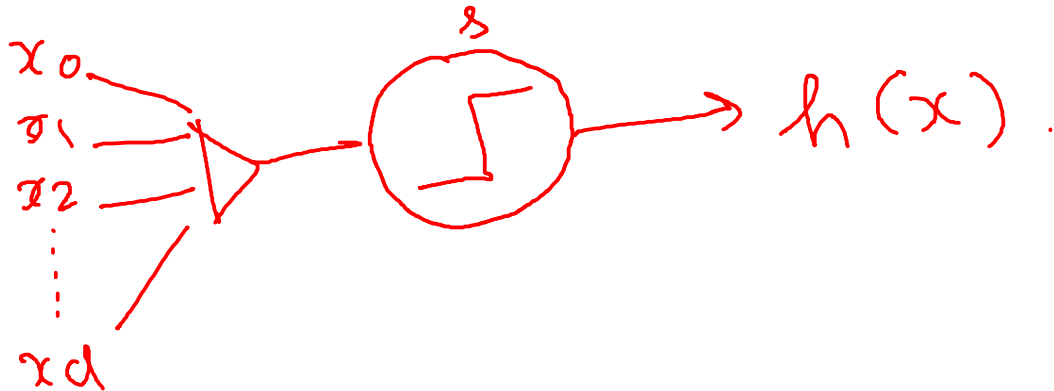
↓

Linear Model

$$s = \sum_{i=0}^d w_i x_i$$

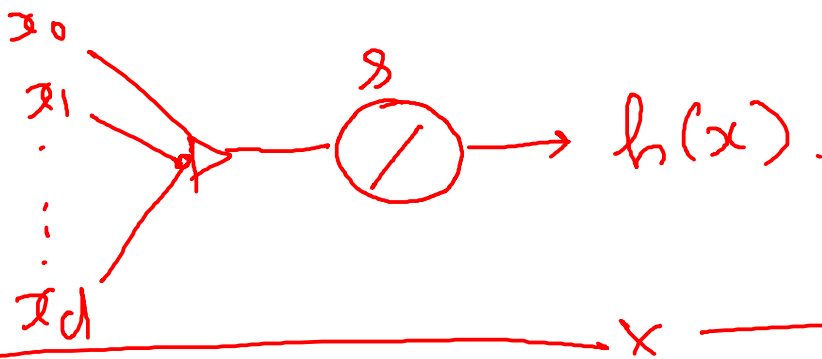
1) Linear Classification.

$$h(x) = \text{sign}(s).$$



2) Linear Regression.

$$h(x) = s.$$



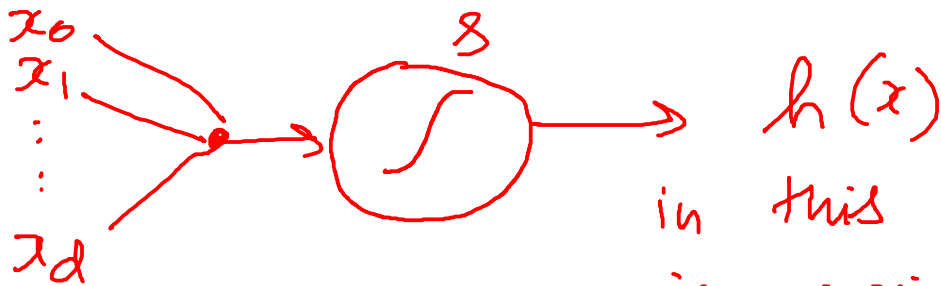
3) Logistic Regression

$$h(x) = \sigma(s).$$



We will take s and apply non-linearity to it.

$$h(x) = \sigma(s).$$



in this case, output is going to be

interpreted as a probability.

X

$$\sigma(s) = \frac{1}{1 + e^{-s}} = \frac{1}{1 + \frac{1}{e^s}}$$

$$\sigma(s) = \frac{e^s}{1 + e^s}.$$

X

- Hard-threshold will be decision making \rightarrow either 0 or 1.
- The sigmoid function brings in the soft threshold (uncertainty)

- Why is it important?

Sometimes Probability is more important than the direct decision 0/1.

- Let's say the problem is about what is the probability that a person will get a heart attack....

$$\boxed{z = w^T x} \quad \frac{[1 \times n][n \times 1]}{\downarrow [1 \times 1]}$$

Let's say (x, y) is the data we have. (Then y is binary)
 $y = \pm 1$.

$$P(y|x) = \begin{cases} f(x) & \text{for } y = +1. \\ 1 - f(x) & \text{for } y = -1. \end{cases}$$

* We are trying to learn the $f()$ here ...

$f: \mathbb{R}^d \rightarrow [0, 1]$ is the probability.

Learn $g(x) = \sigma(w^T x) \approx f(x)$.

Error Measure (Loss function)

for each (x, y) $y = \pm 1$.

y is generated by probability $f(x)$.

We have a plausible error measure based on likelihood.

If $h = f$.

What is the probability of generating this data if your assumption is true.

If that probability is small, the assumption is pure.

(Vice Versa).

if $h=f$, how likely is it to get y from x ?

If one chooses to use the probabilistic approach for choosing the hypothesis;


"What is the most probable hypothesis given the data?"

Here, we are asking ...

"What is the probability of the data given the hypothesis?"

WHAT IS BACKWARDS

$$P(y|x) = \begin{cases} f(x) & \text{for } y = +1 \\ 1 - f(x) & \text{for } y = -1. \end{cases}$$


$$q(y|x) = \begin{cases} h(x) & \text{for } y = +1 \\ 1 - h(x) & \text{for } y = -1. \end{cases}$$

Formula for Likelihood...

$$h(x) = \sigma(W^T x)$$

$$\text{Note:- } \sigma(-s) = 1 - \sigma(s)$$

$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f(-x) = \frac{1}{1 + e^x} = \frac{1}{1 + \frac{1}{e^{-x}}}$$

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

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

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

$$f(-x) = 1 - f(x)$$

_____ x _____

$$P(y|x) = \sigma(y W^T x).$$

Entire dataset

$$(x_1 y_1) (x_2 y_2) \dots (x_n y_n)$$

$$\prod_{n=1}^N \sigma(y_n W^T x_n).$$

Maximizing the Likelihood.

$$\rightarrow \prod_{n=1}^N \sigma(y_n W^T x_n).$$

$$\rightarrow \sum_{n=1}^N \ln(\sigma(y_n W^T x_n))$$

This we have to maximize.

$$-\frac{1}{N} \sum \ln \left(\frac{1}{\sigma(y_n W^T x_n)} \right) \quad \underline{\underline{\text{minimise}}}$$

$$o(\delta) = \frac{e^\delta}{1+e^\delta} \Rightarrow \frac{e^\delta/e^\delta}{\frac{1+e^\delta}{e^\delta}} \Rightarrow \frac{1}{1+e^{-\delta}}$$

$$Error = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n})$$

This we have
to minimize

$$c(h(x_n), y(n))$$

cross-entropy
error

WHAT TO DO ??

→ Gradient Descent !!

Now, what if it is more than
2 classes ?? Now, something
called as Softmax Regression
comes into the picture. --

Logistic Regression

* 1 more attempt

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

$$\text{cost}(h_\theta(x), y)$$

$$= -y \cdot \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

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

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{d(\sigma(x))}{dx} = \frac{0 - (1) \cdot [e^{-x}(-1)]}{(1 + e^{-x})^2}$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$= \frac{1 + e^{-x} - 1}{(1 + e^{-x})^2} = \frac{1 + e^{-x}}{(1 + e^{-x})^2} - \frac{1}{(\quad)^2}$$

$$= \frac{1}{(1 + e^{-x})} - \frac{1}{(1 + e^{-x})^2}$$

$$= \left(\frac{1}{1 + e^{-x}} \right) \left[1 - \frac{1}{(1 + e^{-x})} \right]$$

$$\frac{d(\sigma(x))}{dx} = \sigma(x) \cdot (1 - \sigma(x)).$$

Chain Rule

$$\begin{aligned} \frac{\partial(J(\theta))}{\partial(\theta_j)} &= \frac{1}{m} \cdot \sum_{i=1}^m \left[y_i \left[\frac{1}{h_{\theta}(x_i)} \right] \left[\frac{\partial(h_{\theta}(x_i))}{\partial \theta_j} \right] \right] \\ &+ \sum_{i=1}^m (1 - y_i) \left[\frac{1}{[1 - h_{\theta}(x_i)]} \cdot \frac{\partial(1 - h_{\theta}(x_i))}{\partial(\theta_j)} \right] \end{aligned}$$

= 

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

$$\Rightarrow \frac{-1}{m} * \left(\sum_{i=1}^m \left[y^{(i)} * \frac{1}{h_{\theta}(x^{(i)})} * \sigma(z) * [1 - \sigma(z)] \frac{\partial(\theta^T x)}{\partial \theta_j} \right] \right. \\ \left. + \sum_{i=1}^m \left[(1 - y^{(i)}) * \frac{1}{(1 - h_{\theta}(x^{(i)}))} * (-\sigma(z) * [1 - \sigma(z)]) * \frac{\partial(\theta^T x)}{\partial \theta_j} \right] \right)$$

$$\Rightarrow \frac{\partial(J(\theta))}{\partial(\theta_j)} = \frac{-1}{m} * \left(\sum_{i=1}^m \left[y^{(i)} * (1 - h_{\theta}(x^{(i)})) * (x_j^{(i)}) \right. \right.$$

$$\left. - (1 - y^{(i)}) * h_{\theta}(x^{(i)}) * (x_j^{(i)}) \right]$$

Horrible
difficult to
understand

$$\frac{\partial(J(\theta))}{\partial(\theta)} = \frac{1}{m} * X^T [h_{\theta}(x) - y]$$

[This is what we finally
get !!]

$$\frac{\partial(J(\theta))}{\partial(\theta)} = \frac{1}{m} \cdot X^T [h_\theta(x) - y]$$

$$[X] = [\text{d-brain-points} \times \text{Features}]$$

$$[W] = [4 \times 1]$$

$$[W] = [W] - \alpha [dw]$$

$$[dw] = [4 \times 1]$$

$$[A] = [\text{d-brainpoints} \times 1]$$

$$A = \left[\begin{array}{cccc} p_1 & p_2 & \dots & p_N \end{array} \right]$$

← dbrain →

$$\Psi_{\text{temp}} = \left[\begin{array}{c} \text{d-brain points} \end{array} \right]$$

← →

$$\left[X^T \times (h(\theta) - y) \right]$$

↓

$$\left[\left[\text{Features} \times \text{dbrain} \right] \left[\text{dbrain} \times 1 \right] \right]$$

$$\Rightarrow [\text{Features} \times 1]$$

$$[\text{New Features}] = [\text{old}] - \alpha \left[\frac{\partial (J(\theta))}{\partial \theta} \right]$$

Naive Bayes

1) Assumption :- Gaussian Distribution of all features. And all features are considered independent.

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

This let's say 2 classes.

0
4 features
4 Gaussians

1
4 features
4 Gaussians

Product for given datapoint.

$[x_1 \ x_2 \ x_3 \ x_4]$

$(C_1 \ C_2 \ C_3 \ C_4)$

$(\text{Prior } P)^4$

$(C_1 \ C_2 \ C_3 \ C_4)$

$(\text{Prior})^4$

Whichever is
greater

data belongs to
that class

x