

WEEK - 7CLASSIFICATION

⇒ Supervised learning task

$$(x_i^{(i)}, y_i) \quad y_i \Rightarrow \text{class label}$$

Goal: Group of i/p sample according to class label.

$$\begin{matrix} y_i=0 \\ 0 \ 0 \ 0 \\ 0 \ 0 \ 0 \\ 0 \ 0 \ 0 \\ 0 \ 0 \ 0 \end{matrix} \left\{ \right. \rightarrow 0$$

find out the characteristic,

which makes those

Samples from a particular class

$$\begin{matrix} y_i=1 \\ 1 \ 1 \ 1 \\ 1 \ 1 \ 1 \\ 1 \ 1 \ 1 \\ 1 \ 1 \ 1 \end{matrix} \left\{ \right. \begin{matrix} \text{use this} \\ \text{to classify unseen sample.} \end{matrix} \rightarrow 1$$

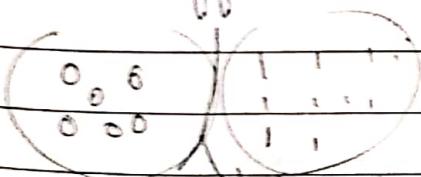
2 Ways of Classification

1) Discriminative

2) Generative

Discriminative

How different are samples of particular class
with respect to other class

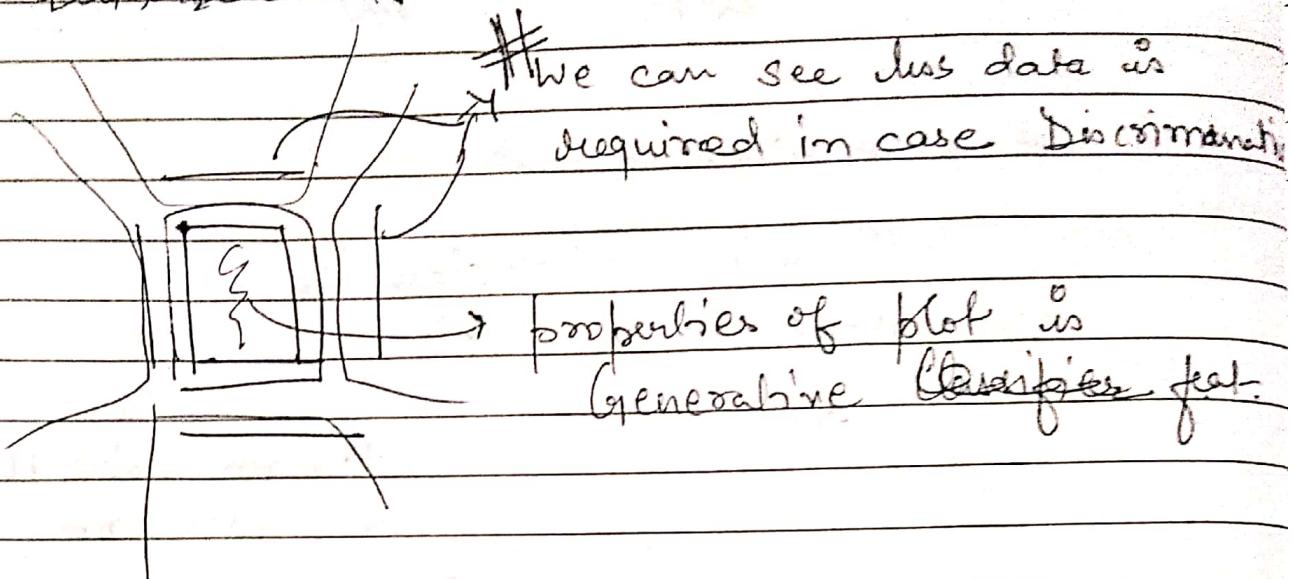


Discrimination of Samples.

Generative

Try to learn Generation characteristic of Samples.

try to figure out actual Generation process, distribution of Class.

Generative vs Discriminative

By this eg. we can see border fits are main focus to create a differentiated / Discriminative classifier

Whereas Generative tries to find features of whole Sample space. So to develop that we require more amount of data.

2) Num classes

Binary v/s Multiclass. (Several binary and true multiclass)

3) Parametric v/s Non-parametric.

$h_w(x)$

$$W = \{w_1, w_2, \dots, w_n\}$$

parameters.

No-prior model.

we try to learn actual distribution

Eg. Regression, GMM

NN is eg. Non-parametric.

We estimate model

k-means is eg.

Logistic Regression

a) Starting state for this ~~Regressi~~ classifier was regression.

b) ~~Linear~~ Regression using sigmoidal function.

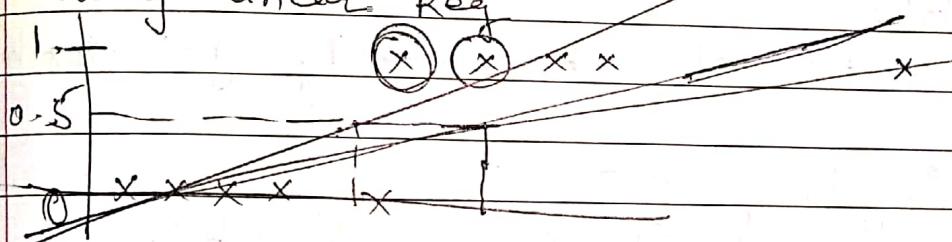
Binary Classification

Spam / Not Spam

Benign / Malignant

and many more.

Using linear Reg



If we keep threshold as 0.5

$$h_0(x) \geq 0.5 \equiv 1$$

$$h_0(x) < 0.5 \equiv 0$$

In this case this Linear Regression seems to work fine

but if there is one more point.

The circled one will give wrong result

Some data points can be misclassified in linear regression.

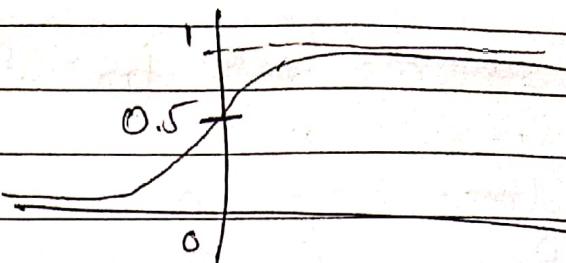
for linear reg

$$h_0(x) < 0 \quad \text{and} \quad h_0(x) > 1$$

Logistic Regression

$$h_0(x) = g(\theta^T x)$$

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



Sigmoid function.

$$0.5 \leq h_0(x) \leq 1$$

Interpretation

$h_0(x)$ = estimated probability $y=1$ for input x

$$\text{If } h_0(x) = 0.7$$

there is 70% chance of y being 1

$$P(y=0/x) = 1 - P(y=1/x)$$

$$g(z) \geq 0.5$$

$$z \geq 0$$

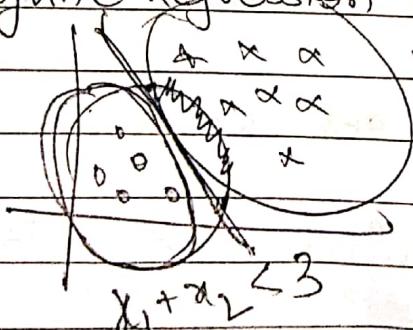
$$z = \theta^T x$$

$$\theta^T x \geq 0$$

$$g(z) < 0.5$$

$$z < 0$$

Logistic Regression is Discriminative classifier



$$x_1 + x_2 \geq 3$$

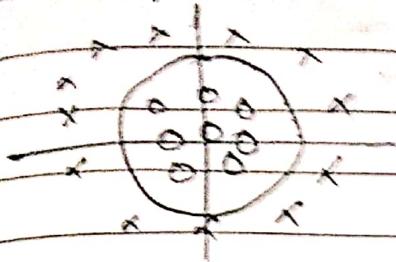
$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

Predict $y = 1$

$$-3 + \alpha_1 + \alpha_2 \geq 0$$

$$x_1 + x_2 \geq 3$$

Non-linear decision boundary



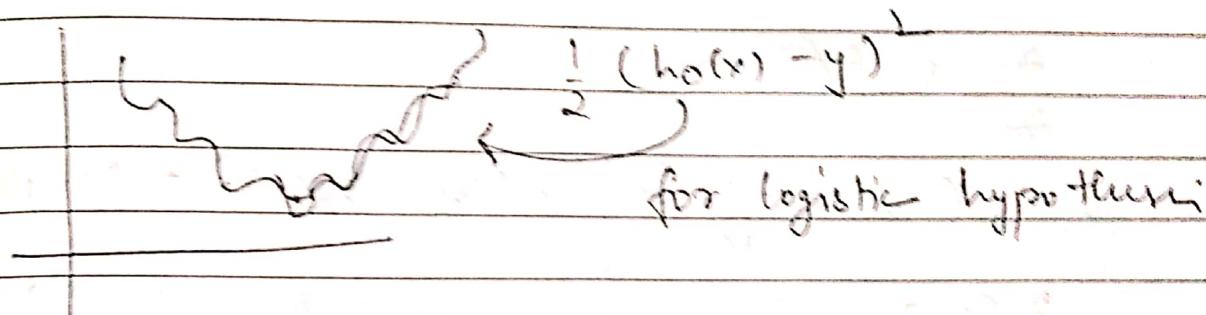
$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2$$

Predict

$$y=1 \quad -1 + \theta_1^2 + \theta_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 1$$

How to estimate θ values.



MSE may lead to local minimum as nature of $h_0(x)$ is not linear but sigmoidal.

New cost function

$$\text{Cost} = \begin{cases} -\log(h_0(x)) & \text{if } y=1 \\ -\log(1-h_0(x)) & \text{if } y=0 \end{cases}$$

piece-wise convex nature is exploited

$$\text{Cost} = -y \log(h_0(x)) - (1-y) \log(1-h_0(x))$$

when $y=0$

$y=1$

$$\text{Cost} = -\log(1-h_0(x))$$

$$\text{Cost} = -\log(h_0(x))$$

Gradient descent

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

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

Derivation

$$\frac{\partial}{\partial w_j} h_w(x^i) =$$

$$C(w) = -\frac{1}{M} \sum_{i=1}^M y^i (\log(h_w(x^i))) + (1-y^i) \log(1-h_w(x^i))$$

$$\frac{\partial C(w)}{\partial w_j} = -\frac{1}{M} \sum_{i=1}^M y^i \frac{\partial \log(h_w(x^i))}{\partial w_j} + (1-y^i) \log(1-h_w(x^i))$$

$$\frac{\partial}{\partial w_j} \left(\log(h_w(x^i)) \right) = \frac{1}{h_w(x^i)} \times \frac{\partial}{\partial w_j} (h_w(x^i))$$

$$\frac{\partial}{\partial w_j} \left(\log(1-h_w(x^i)) \right) = \frac{1}{1-h_w(x^i)} \times \frac{\partial}{\partial w_j} (-h_w(x^i))$$

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

$$\frac{\partial h_w(x^i)}{\partial w_j} = \frac{-1}{(1 + e^{-w^T x})^2} \underbrace{(e^{-w^T x})}_{(-1)^j} \underbrace{(x^i)}_{(x^i)_j} e^{-w^T x}$$

$$\frac{\partial C}{\partial w_j} = \underbrace{M}_{M \times 1} \underbrace{y_i^o}_{1 \times 1} \frac{1}{h_w(x_i^o)} \times + \frac{e^{-w^T x_i^o}}{(1+e^{-w^T x_i^o})^2} + \cancel{(C \times \delta)}(x_i^o)$$

$$\frac{y_i^o}{h_w(x_i^o)} \times \frac{1}{x_i^o - e^{-w^T x_i^o}} \times (h_w(x_i^o))^2 \times x_i^o$$

$$f x_j^o y_i^o \times e^{-w^T x_i^o}$$

$$(1-y) \frac{1}{1-h_w x} \times \frac{e^{-w^T x_i^o}}{(1+e^{-w^T x_i^o})^2} x_j^o$$

$$\frac{1}{1-h_w x} = \frac{1}{1 - \frac{1}{1+e^{-w^T x_i^o}}} = \frac{1+e^{-w^T x_i^o}}{e^{-w^T x_i^o}}$$

$$(1-y) \times \cancel{\frac{1+e^{-w^T x_i^o}}{e^{-w^T x_i^o}}} \times \frac{e^{-w^T x_i^o}}{(1+e^{-w^T x_i^o})^2} x_j^o$$

$$- \frac{(1-y) x_j^o}{(1+e^{-w^T x_i^o})}$$

$$x_j^o y_i^o \times \frac{e^{-w^T x_i^o}}{1+e^{-w^T x_i^o}} \quad f x_i^o y_i^o \times \frac{-x_j^o}{1+e^{-w^T x_i^o}}$$

C. g.e

$$x_j^o y_i^o \left(\frac{1+e^{-w^T x_i^o}}{1+e^{-w^T x_i^o}} - \frac{x_j^o}{1+e^{-w^T x_i^o}} \right)$$

$$x_i^o \left(y_i^o - h_w(x_i^o) \right)$$

$$C(w) = -\frac{1}{M} \sum_{i=1}^M (h_w(x^i) - y^i) x^i \quad (i)$$

Multiclass Classification

- ⇒ Multiclass can be realized as one vs all sort binary classifier
- ⇒ Train a logistic Regression Classifier $h_w^i(x)$ for each class i to predict the prob y^i
- ⇒ During Testing, for given x , pick class that maximizes the $\max_i h_w^i(x)$

Use of Clustering for Classification.

1st method

In a common feature space clustering is done & label are assigned.

Disadvantage

! There may be overlaps in same feature space soft boundaries which may lead to misclassification

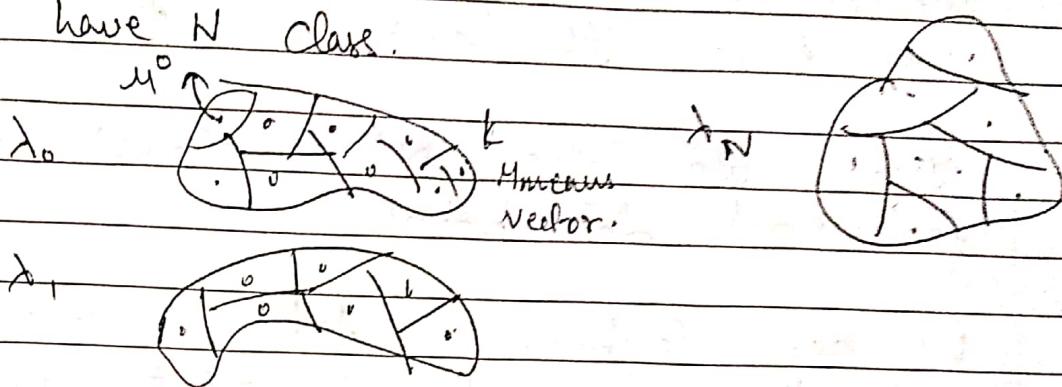


feature vector

cluster boundary is
not proper

Generative Classifier Approach.

we have N class.



Classification

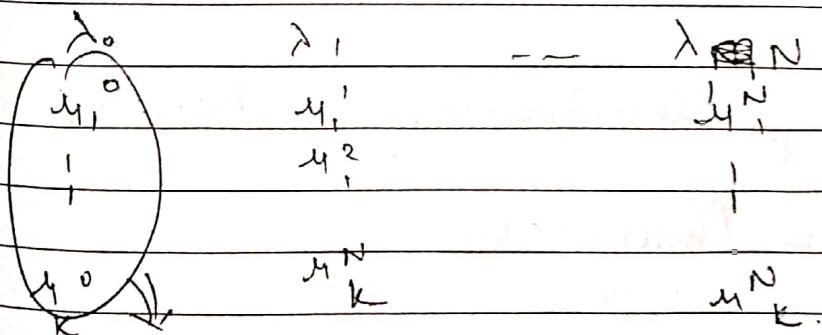
Method we have Separate feature space for each class label

we cluster that particular feature space into k -clusters.

using this mean vector we try to compare and estimate class of unknown data

Classification

$$\mathbf{x} = \{x_1, x_2, \dots, x_N\}$$



min of
all up dis

We calculate distance for all class cluster and select the least one.

Then we have N (num class) distance for each data point.

We again compare and take minimum of all those class, the class that has minimum ~~said to~~ we classify that point to that class.

$$\{d_0^1 \dots d_0^N\}$$

$$\min d_0^j$$

$i =$ classes

$j =$ data point.

Why cluster?

If we don't we have to compare to all data points in that class that is computational headache.

This way of classification is part of Generative Classification.

No need for other class data to build classifier.