

Support Vector

Machine

→ well used classifier

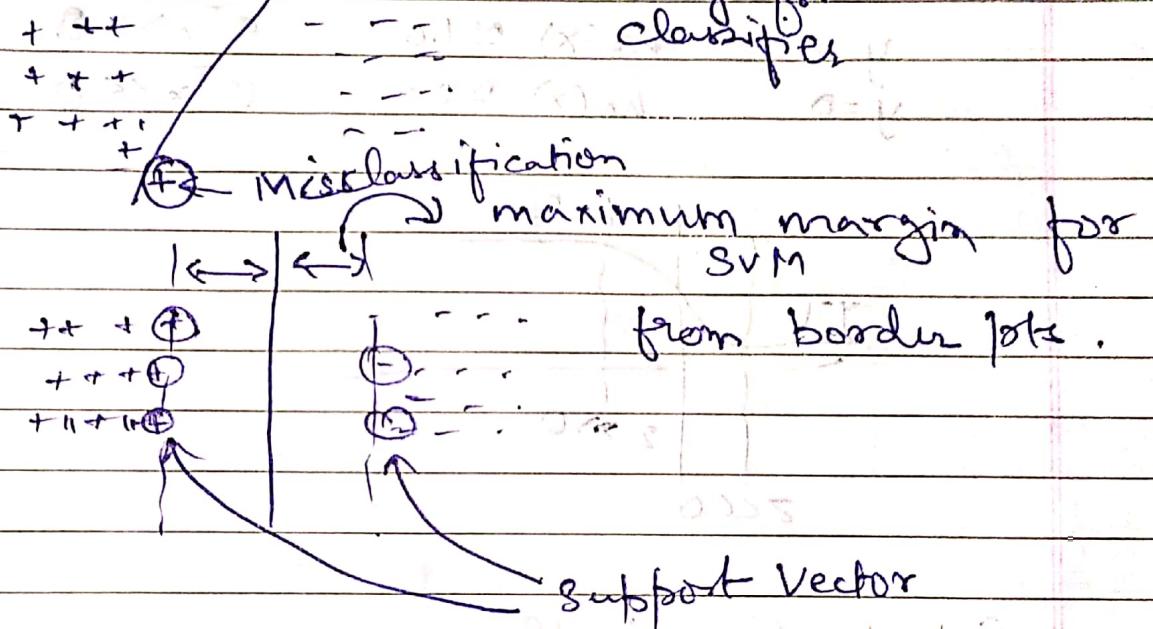
→ Binary Classifier

→ Large-margin classifier.

Margin b/w the +ve & -ve classes is larger than as compare to other binary classifier

Misclassification

e.g. decision boundary for some binary classifiers



Why SVM should perform better?

data \Rightarrow sample is not proper representation of whole data. if new point comes which at border b.c.z. first D.B. it is misclassified

SVM by choosing Maximum Margin tries to reduce this misclassification because of improper D.B.

Development of SVM from Logistic Regression

LR $\Leftrightarrow h_0(x) = \frac{e^{-\theta^T x}}{1+e^{-\theta^T x}}$
 $z = \theta^T x$

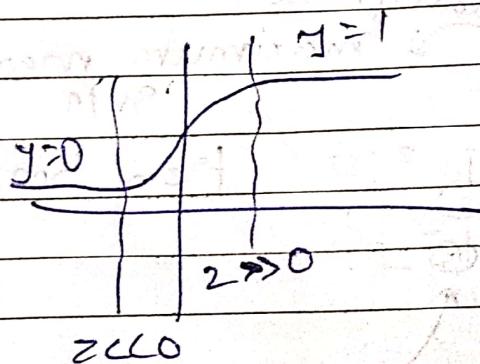
$y = 1 \quad h_0(x) > 0 \Rightarrow z > 0$

$y = 0 \quad h_0(x) < 0 \Rightarrow z < 0$

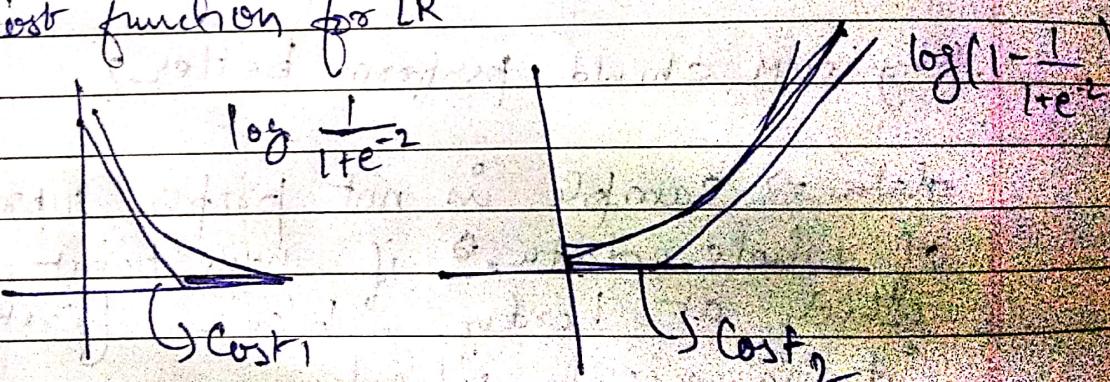
SVM

$y=1 \quad h_0(x) \approx 1.0 \quad \theta^T x \gg 0$

$y=0 \quad h_0(x) \approx 0.0 \quad \theta^T x \ll 0$



Cost function for LR



piece wise linearization for SVM

Cost₁

SVM

$$\min \sum_{i=1}^m y_i \text{cost}_1(\theta^T x_i) + (1-y_i) \text{cost}_0(\theta^T x_i) + \frac{\lambda}{2} \sum_{j=0}^n \theta_j^2$$

A ↓ B ↓

Regularization

Regularization is done so that there is no overfitting on training data of Machine Learning Parameters.

ex: $\min ((u-5)^2 + 1) \times 10 \quad u = 5$
 $\min 10(u-5)^2 + 10 \rightarrow u = 5$

$$A + \lambda B$$

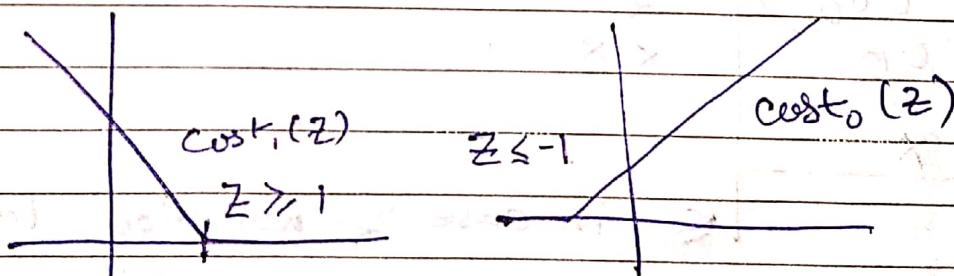
Minima
domain same

$$\min A + \lambda B \geq \frac{1}{\lambda} A + B$$

When multiplied by

$$\equiv C(A+B)$$

Large Margin



$$y = 1 \quad \theta^T x \geq 1$$

$$y = 0 \quad \theta^T x \leq -1$$

when C is large

$$C \cdot A + B$$

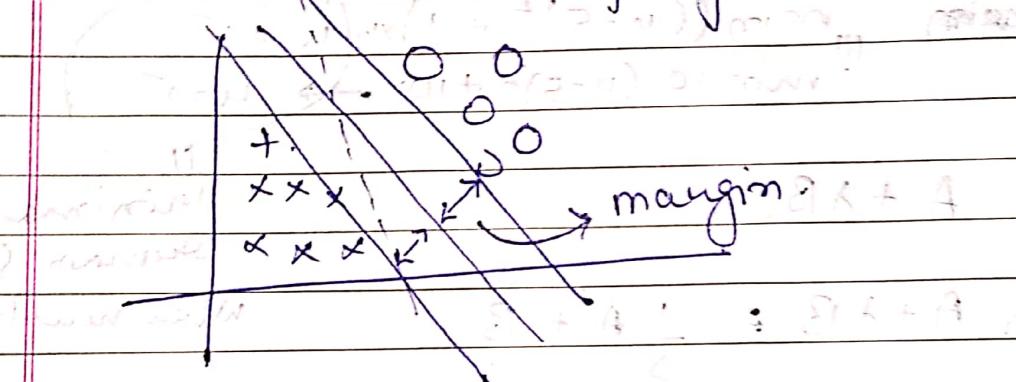
so for minimization.

$$\Theta^T x \geq 1$$

$$\text{min } C \cdot 0 + \frac{1}{2} \|y\|^2$$

$$\text{s.t. } \Theta^T x^i \geq 1 \text{ if } y^i = 1$$

$$\Theta^T x^i \leq -1 \text{ if } y^i = 0$$



large margin. in presence of outliers.

0 0 0

0 0 0

0 0 1

x

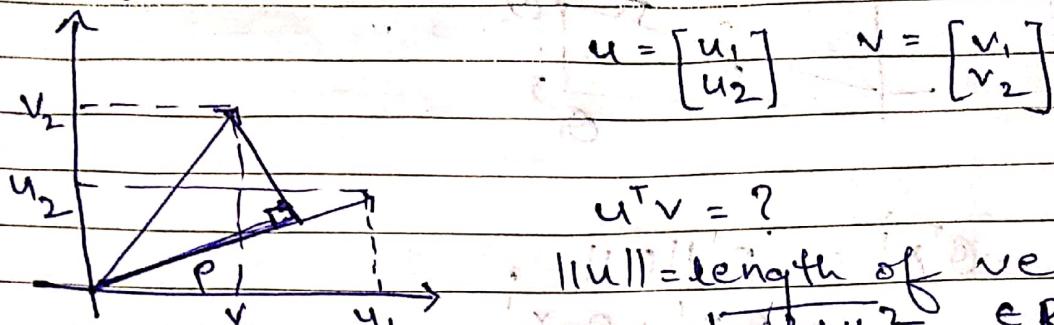
when C is set too large

in case of not too large.

C

it doesn't make sense to have dotted line. just b'cos we have one outlier.

how Optimization function leads to large Margin Classifier.



p is length of Projection of v on u .

$$\begin{aligned} u^T v &= p \cdot \|u\| \\ &= u_1 v_1 + u_2 v_2 = v^T u \end{aligned}$$

* p is signed scalar.

As $C\alpha + \beta$ is set high which will cause final Cost function to be

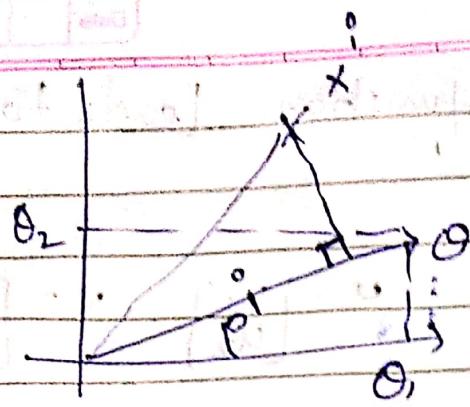
$$\min_{\theta} \frac{1}{2} \sum \theta_j^2$$

$$\theta^T x^i \geq 1 \text{ if } y^i = 1 \text{ and}$$

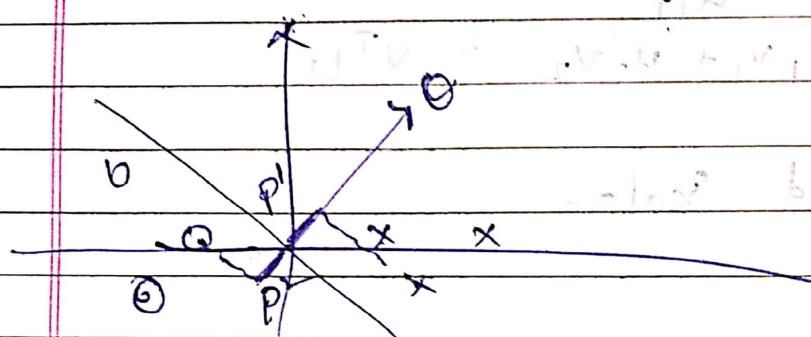
$$\theta^T x^i \leq -1 \text{ if } y^i = 0$$

Simplification: $\theta_0 = 0$, $n=2$.

$$\min_{\theta} \frac{1}{2} \sum \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \|\theta\|^2$$



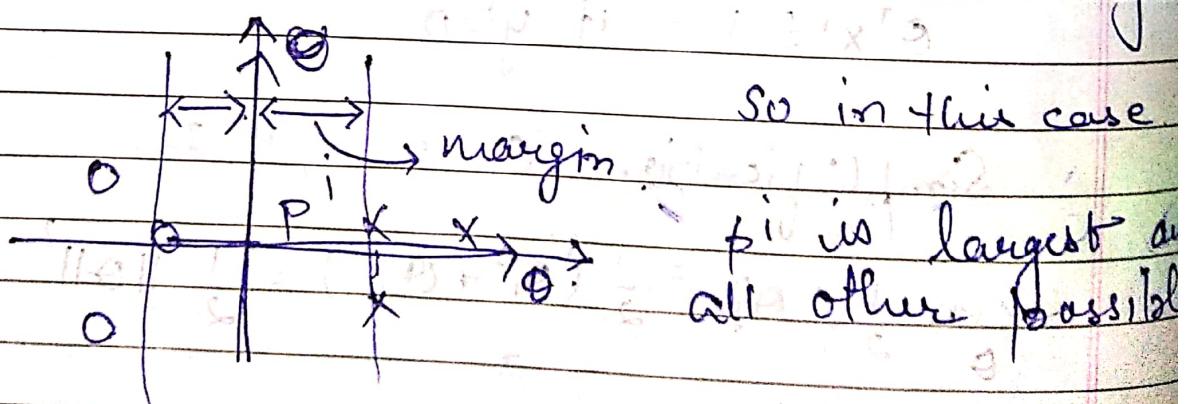
$$\theta^T x^i = p^i \|\theta\| > 1 \text{ if } y^i = 1 \\ \theta^T x^i = p^i \|\theta\| < -1 \text{ if } y^i = -1$$



So now let's say $\|\theta\|$ is small

So now $\|\theta\| > 1$

\downarrow \downarrow
Small $\|\theta\|$ This needs to be large



So in this case

$\|\theta\|$ is largest among all other possible

implies

$\|w\|_2^2 \geq 1$
 \downarrow can be small.
 large margin.

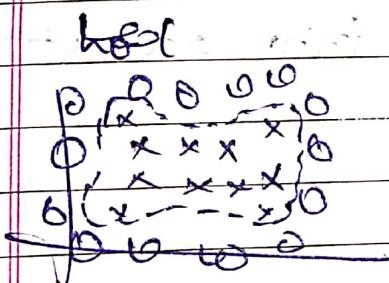
As $\theta_0 = 0$, Decision boundary always passes through origin

This way we have large margin after optimization.

ii

kernel

Non-linear Decision Boundary



$h_0(x) = \text{some higher degree polynomial}$

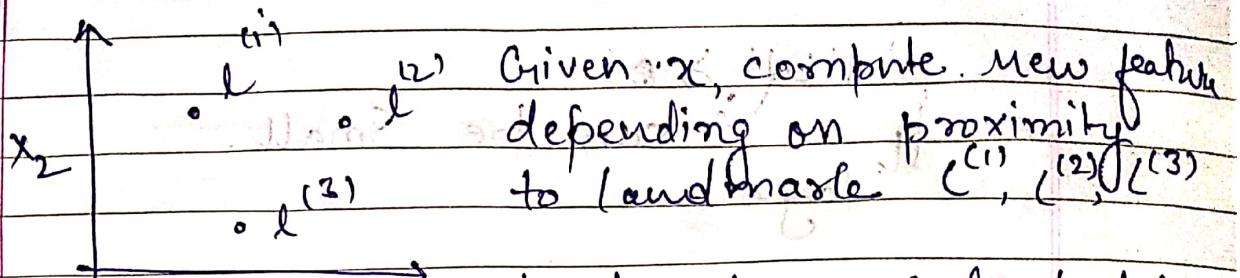
which makes it computationally expensive

expensive.

Is there is better way?

Will try to generate new feature from old ones.

Kernels



x_1 and landmark can randomly taken

$$\text{Given } x_1 : f_1 = \text{similarity}(x_1, l^{(1)}) = e^{-\frac{\|x_1 - l^{(1)}\|^2}{2\sigma^2}}$$

$$f_2 = \text{similarity}(x_1, l^{(2)}) = e^{-\frac{\|x_1 - l^{(2)}\|^2}{2\sigma^2}}$$

$$f_3 = \text{similarity}(x_1, l^{(3)})$$

↑
Kernel

Gaussian

also denoted by kernel.

$$k(x, u)$$

Taking Gaussian kernel into account.

$$x_1 \approx l^{(1)}$$

$$f_1 = \exp\left(-\frac{\|x_1 - l^{(1)}\|^2}{2\sigma^2}\right) \approx 1$$

When x_1 is far from $l^{(1)}$,

$$f_1 = \exp\left(-\frac{\text{large no.}}{2\sigma^2}\right) \approx 0$$

New hypo thesis function can be given in terms of new feature vector
 ~~f_1, f_2, f_3~~ f_1, f_2, f_3

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 > 0$$

Diff

When f_1 is near to 1

$$f_1' \approx 1, f_2 = 0, f_3 = 0$$

$$\theta_0 + \theta_1 \approx$$

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$

$$-0.5 + f_1 + 1.f_2 > 0$$

$$0.5 > 0$$

So

SVM parameters

$C(= \frac{1}{2})$ large C : lower bias, high variance

small C : higher bias, low variance

Other choices of kernel:

- Polynomial kernel | $k(x, l) = (x^T l)^L$
or
 $(x^T e)^N$
- string kernel
- chi-square kernel
- histogram intersection kernel.
- linear / with no given kernel

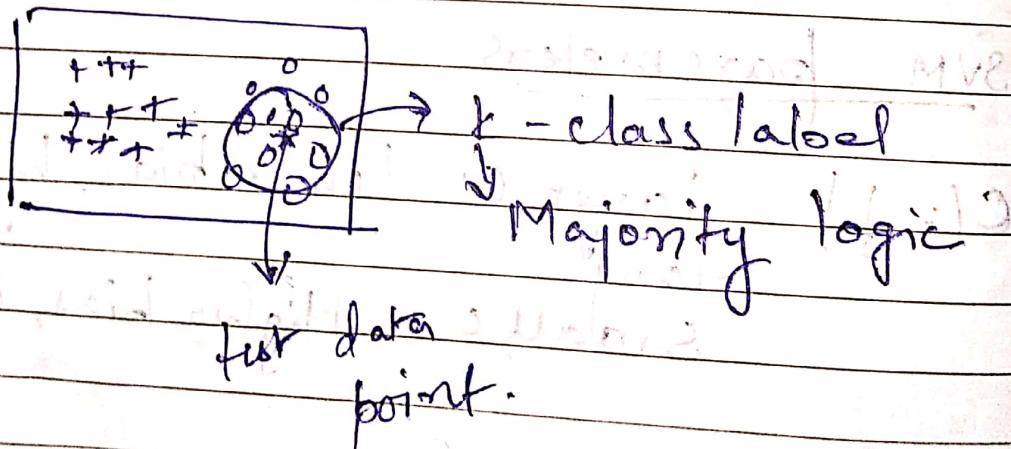
k-nearest neighbour Classifier.

knn.

There is no training involved

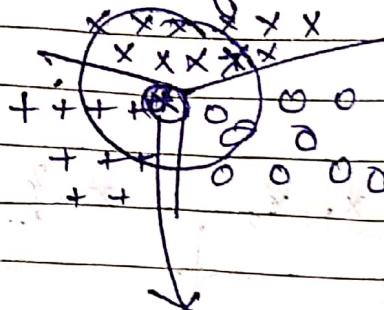
→ Prediction of class label is based on k nearest neighbour

for testing



~~Example~~Limitation -

- Misclassification of border pts / Thin borders



Now it may happen there are majority of k belongs to 'x' for given pt. P

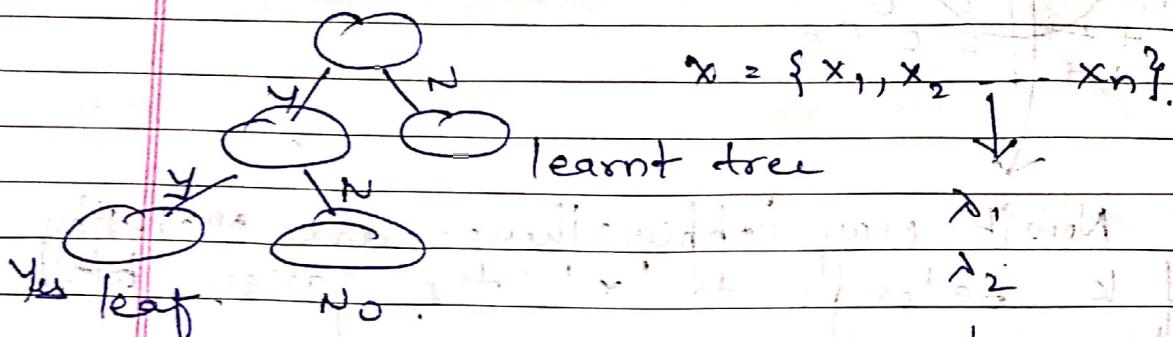
but it belongs to '+' side of decision boundary of SVM. so this way it may lead to wrong result.

- Large data, Prediction stage is slow
- Sensitive to scale of data and irrelevant features
- Requires high memory - need to store all of the training data
- Given that it stores all training, it can be computationally expensive.

Decision Tree

Random Forest Classification

Binary / Multiclass



Tree Structure can be different for each class.

In Random Forest we have multiple decision tree for one class

λ_1 - first Class λ_2 - 2nd \dots λ_k

Testing

$$\lambda_1 \rightarrow M \times N \leftarrow S_1$$

$$S_{11}, S_{12} \dots S_p$$

$s_{11}, s_{12}, \dots, s_{1P}$
 $(m_{xi}) \xrightarrow{\text{with replacement}} (n; j) \xrightarrow{(P, q)}$

Random Creation of 'Multiple' dataset from $M \times N$

$d_{t1}, d_{t2}, \dots, d_{tp}$

Majority logic or avg. logic.

Overfitting is related to depth of tree.

- Entropy
- Information Gain
- Gini index

Dimension Reduction

→ PCA

→ LDA: linear discriminant analysis

Given: $X = \{x_1, x_2, \dots, x_N\} \rightarrow N$ -dimensional

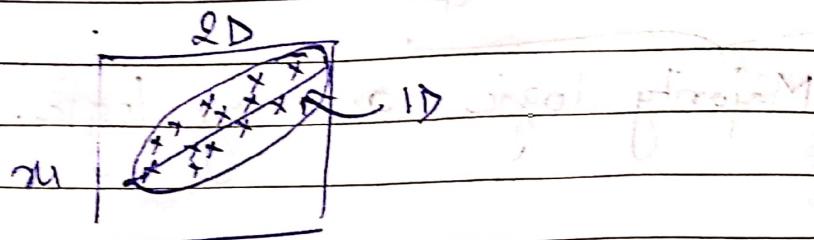
Obtain $Z = \{z_1, z_2, \dots, z_k\} \rightarrow k$ -dimensional

$X \Rightarrow Z$ where $k < n$

This process is called dimensionality reduction.

uses

- To Reduce Computational Cost.
- To remove redundant dimension.
- Some dimension, info is not significant
- Variance is low / negligible
- Data Visualization in lower dim.

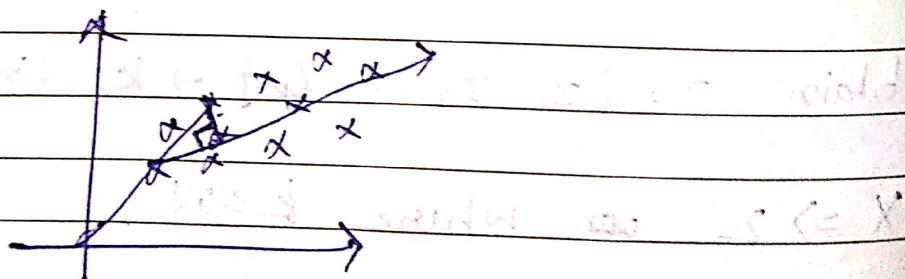


Principal Component Analysis (PCA)

$N \rightarrow k$
projection / Principal Component.

Project data only along the direction of maximum variance.

- Try to Capture maximum variation in N dimension while transforming to k-dimension.



Projection line is Vertical line from data pts to fitted straight line.

Angle to the fitted straight line that gives shortest Projection error.

PCA Method /Algorithm

- Mean Normalization $\bar{y} = 0, \sigma^2 = 1$
- Feature Scaling: to bring dimensions to similar range.
- Covariance Matrix $C = X^T X = (N \times N) \text{ dim}$

Eigen Analysis.

Covmat: Σ be an $N \times N$ matrix

λ be eigen value of Σ if $p \neq 0$ vector

$$\Sigma p = \lambda p. \quad \lambda \text{ is scalar}$$

$$\left[\begin{array}{c} \\ \\ \end{array} \right]_{N \times N} P_{N \times 1} = \lambda P_{N \times 1}$$

Solve for Eigen values

$$\det(\Sigma - \lambda I) = 0$$

$$\sum \vec{P} = \lambda \vec{P}$$

$$\lambda_1 \rightarrow p_1$$

$$\lambda_2 \rightarrow p_2$$

\vdots

$$\lambda_q \rightarrow p_q$$

Descending order

$$\lambda_1 > \lambda_2$$

$$\dots \lambda_q$$

if we want only k dimension

$$\lambda_1 > \lambda_2 > \dots > \lambda_k$$

we will capture max variance.

disjoint cap

not a good fit at 3rd

conflict of fit is a bad sign idea

make it a high order

fit to first

term

extra term with value