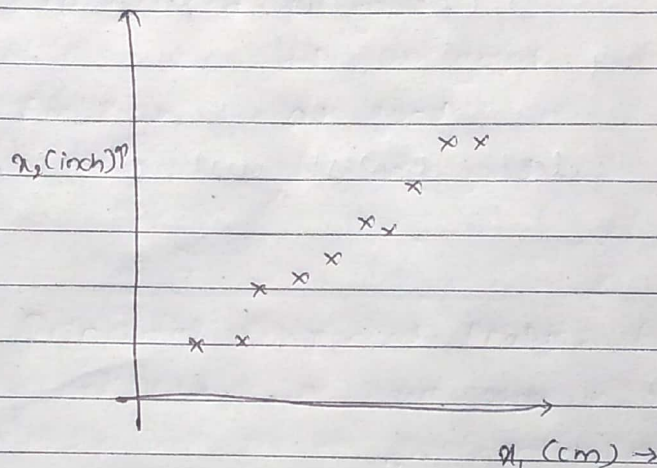


Dimensionality reduction

Dimensionality reduction helps by not only to reduce the memory usage but also the learning speed is improved.

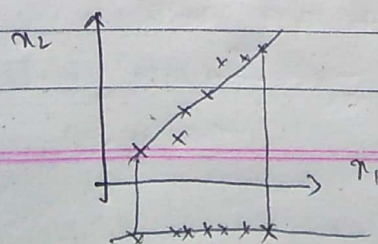
Let's say there are 2 features length in cms & length in inches of same object. both rounded to nearest value.

We can reduce this to one feature. In a bigger problem we have loads of features hence it is difficult to keep track.



for example let's say in plane driving example: if Pilot's skill, Pilot's enjoyment are the features they are highly correlated and can be clubbed together to form a new feature 'aptitude'.

So in above example we draw a st. line that fits the curve & drop the corresponding pts on new feature.



Thus each feature will be indicated just by one number rather than 2 hence we halve the reduce memory requirement to half

Similarly data reduction can be from 3D to 2D. or in other words we reduce k dimensional data to c dimensional ($c < k$).

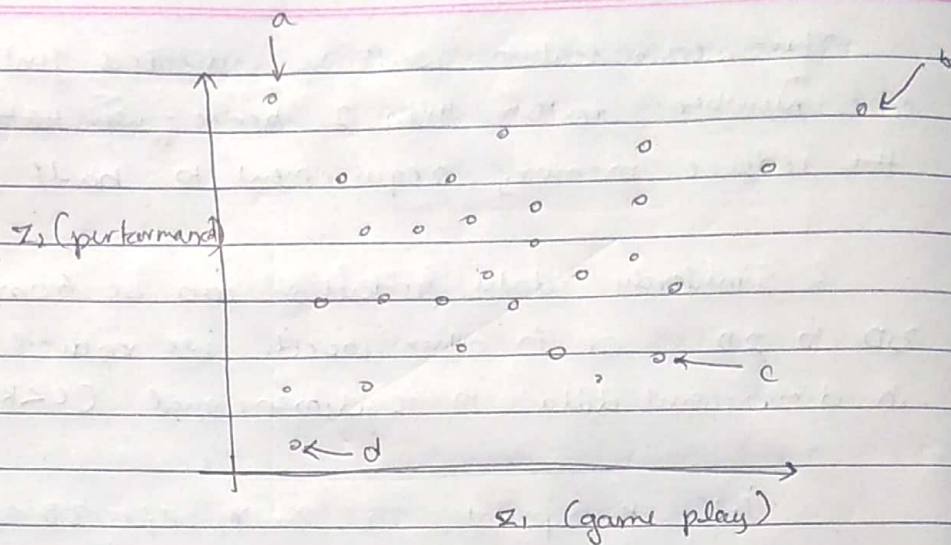
Let's the points in x_1, x_2, x_3 3D space lie on a single plane then this can be projected on a plane z_1, z_2 .

Data visualization

Let's say we have huge no. of features & a huge number of dataset. This can not be plotted. Let's say there are 50 features, we reduce it to 2 which summarizes the data. Thus we can plot the data & visualize the trend.

Let's take an example of a TT player
Matches played (x_1), Strokes per match (x_2), Smashes (x_3), successful serve (x_4), back hand shots (x_5), time of play (x_6), Spinning ball shot (x_7).

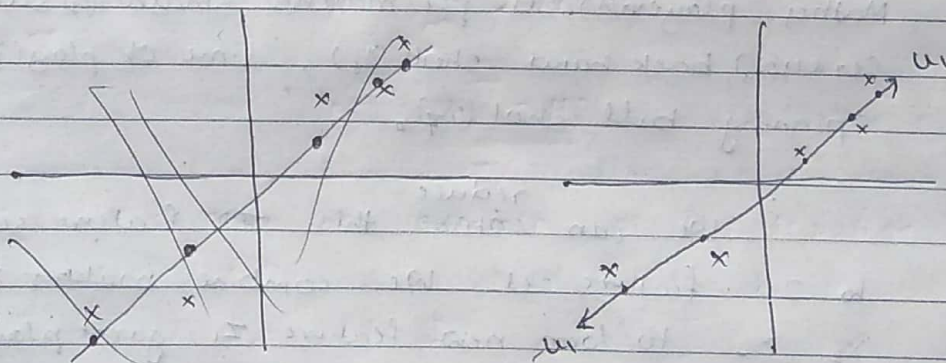
We can ^{reduce} combine this 7D feature set to 2D feature set. We combine matches x_1 , x_2 , x_6 to form new feature z_1 game play & x_3, x_4, x_5, x_7 to form overall performance z_2 .



- interpret to
we may have kinds of players
- a- Experienced & good
 - b- " but not good
 - c- Amateur but good
 - d- " & not good

Principal Component Analysis (PCA)

let's say we have some 2D data & want to reduce it to 1D. hence we need find a best fit line

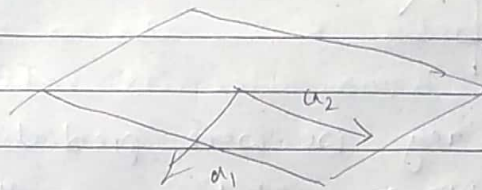


What PCA does is find a lower dimensional line, which best fits the data i.e. has least (or vector)

projection error

The direction of line u_1 doesn't matter

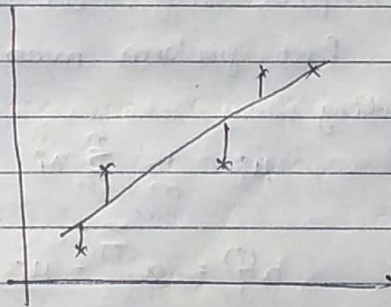
When we are reducing 3D to 2D we should find 2 vectors



Similarly when we reduce to k dimensions we need k vectors

PCA is not linear regression

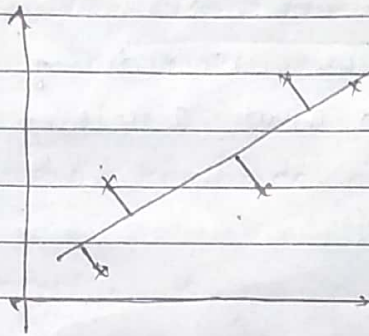
In linear regression we predict some y for given input feature x_1 . So as to minimize the squared distance betⁿ predicted & actual values



we do find squared distance that is vertical or \perp to x_1 axis

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

While PCA tries to minimize the squared distance drawn from to fitting line.



In lin. reg we were predicting some distinct variable y . That is not there in PCA.

In PCA all variables are treated symmetrically there is no special variable as in the case of lin. reg, the variable y .

PCA algorithm.

Step 1) → Data preprocessing: on a given training set we should first perform mean normalization or feature scaling.

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$x_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{s_j}$$

s_j may be (max-min) or std deviation

Step 2) Compute covariance matrix

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \underbrace{\left(x^{(i)} \right) \left(x^{(i)} \right)^T}_{n \times n}$$

$n \times 1 \quad 1 \times n$

Now compute 'eigenvectors' of matrix Σ
by $U, \Sigma, V = \text{Singular Value Decomposition}(\Sigma)$

Σ will be $n \times n$ matrix

We need matrix U for which is $n \times n$ matrix

$$U = \begin{bmatrix} | & | & | & & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & & | \end{bmatrix}$$

To reduce n dimension to k dimension
we should choose first k columns of U matrix
i.e. from $u^{(1)}$ to $u^{(k)}$

$\therefore U_{\text{reduced}}$ will be $n \times k$ matrix

$$Z = \underbrace{(U_{\text{reduced}})^T}_{k \times n} \underbrace{X}_{n \times 1} = k \times 1$$

This exclude bias variable x_0

Choosing the value of k
 k is no. of principle components (z_1, z_2, \dots, z_k)

$$\text{Avg squared proj}^2 \text{ error} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$$

$$\text{Total variation} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

we typically choose K for which

$$\frac{\text{Avg. squared proj}^2 \text{ error}}{\text{Total variation}} \leq 0.01 \quad (17.7)$$

i.e. 99% variance is retained

Generally K is chosen such that 95-99% of variance is retained

Algorithm

Start with $K=1$
for $K=1$ to n
Compute $U^{(K)}, Z^{(1)}, Z^{(2)}, \dots, Z^{(K)}, \lambda_{\text{approx}}^{(1)}, \dots, \lambda_{\text{approx}}^{(K)}$

calculate

check if

$$\frac{1/m \sum_{i=1}^m \|x^{(i)} - z_{\text{approx}}^{(K)}\|^2}{1/m \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01, ?$$

This calculation can be done in a simpler way by using for given K

$$1 - \frac{\sum_{i=1}^K s_{ii}}{\sum_{i=1}^n s_{ii}} \leq 0.01, ?$$

Where s is matrix S in U, S, V obtained from $U, S, V = \text{svd}(\Sigma)$.

S is $n \times n$ matrix whose diagonal elements are non-zero rest all will be zero

Now the smallest K that satisfying the

condition is selected

Going back from KD (2D) to nD (3D)

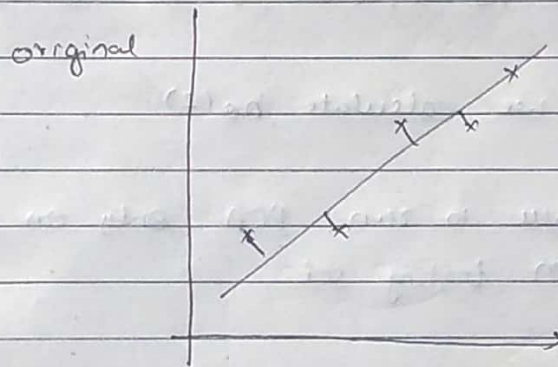
let's say we have

$$Z = (U_{\text{reduced}})^T X$$

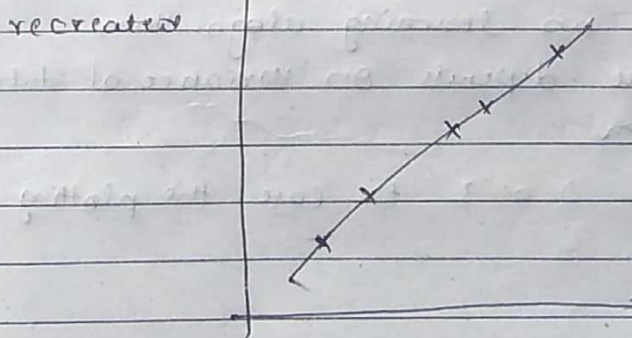
$$X_{\text{approx}} = U_{\text{reduced}} \times Z$$

$\underbrace{n \times k \quad k \times 1}_{n \times 1}$

As due to PCA projection error is too small X_{approx} will be almost coinciding X



projected



Supervised learning speedup

We have supervised learning problems with dataset $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})$
let's say $x^{(i)}$ are very high dimensional

→ We extract the inputs & decrease the labels
 $x^{(1)} \dots x^{(m)} \rightarrow 10,000$

→ supply PCA to get

$z^{(1)} \dots z^{(m)} \rightarrow 1000$

→ new training set will be
 $(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}) \dots (z^{(m)}, y^{(m)})$

→ Now we calculate $h_\theta(z)$

→ We have to run PCA only on training
& not on testing set

Applications

- Compression

- reduce memory needed to store data
- Speed up learning algorithm
(K value depends on variance of data)

- Visualization

(K is used 2 or 3 to ease the plotting)

→ PCA should not be used ^{as} for solution for overfitting

PCA might work in this regard but it is not a good way to select PCA to address overfitting instead use regularisation

→ PCA might throw out some piece of information

→ Sometimes PCA is used where it shouldn't have been used

hence run learning algorithm without PCA meaning use the original training set if it is not doing what we want then go for PCA