

convert d dim. to b dim. $b < d$

why?

- visualization (usually 2 or 3 features)
- to get less noisy and irrelevant features
- to reduce the size (to be easier in computing)
(to avoid out of memory)

- It's NOT Feature Selection

IT is Transformation (new features)

	x		x
s_1	5	1	4
s_2	5	3	3
s_3	5	5	4
s_4	5	7	3

① remove low variance

② projection

Projection is

$$x \cdot u = u \cdot x = u^T x = \sum_{i=1}^d u_i x_i$$

features must be normalized

How to transform:-

① PCA (linear)

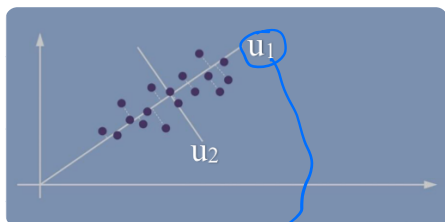
come up with new dimensions
that have high variance

to project any point to any axis \Rightarrow dot product

variance of x
in direction u

$$\text{total } \sigma = u^T \Sigma u$$

natural coordinate



Vector with maximum
variance

in python
SVD

$$\begin{bmatrix} 1, 0 \end{bmatrix} \begin{bmatrix} 1 & 0.85 \\ 0.85 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$D \times d \sum = \begin{bmatrix} \text{eigenvalues} \\ \text{eigenvectors} \end{bmatrix}$$

dot product with your
data to get the new
features

$$\begin{bmatrix} \lambda_1 & \lambda_2 & \lambda_3 & \lambda_4 \end{bmatrix}$$

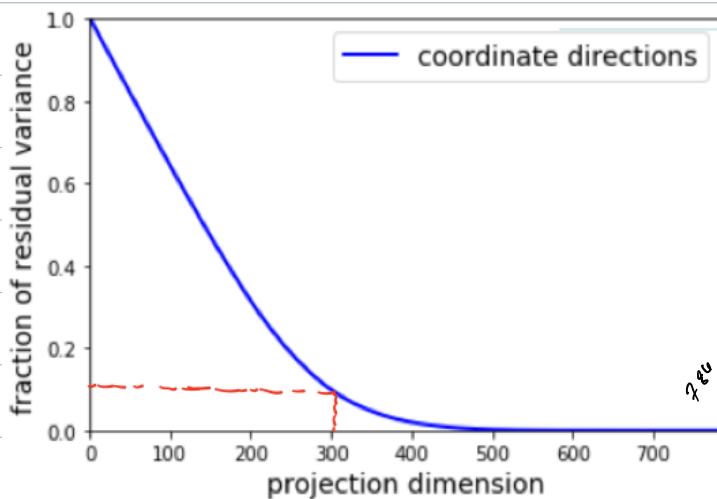
the axis itself
the best axis (highest value)

if you don't have
exact σ of axis
you fix % of

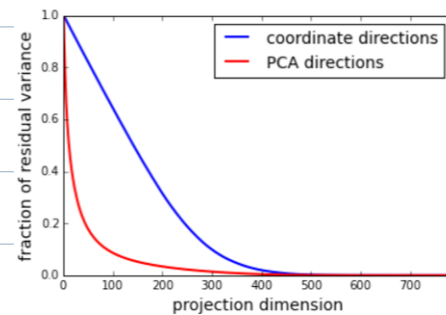
σ then get $\lambda_1, \dots, \lambda_n$
that give you % you want

how σ I will get if I select 3

$$\sigma = \frac{\lambda_1 + \lambda_2 + \lambda_3}{\sum \lambda}$$



if we keep 300 features
we will lose just 10% of
the variance

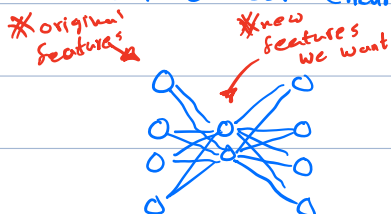


● B: If you use PCA to project d -dimensional points down to j principal coordinates, and then you run PCA again to project those j -dimensional coordinates down to k principal coordinates, with $d > j > k$, you always get the same result as if you had just used PCA to project the d -dimensional points directly down to k principal coordinates.

True

② T-Sne (non-linear)

auto-encoder (neural network)



loss = squared loss

$$-y \log \hat{y} - (1-y) \log (1-\hat{y})$$

$$-y \log \hat{y}$$

$$c = 1 - \max p_i$$

$$\text{Gini} = 1 - \sum (p_i)^2$$

$$\text{Entropy} = -\sum p_i \log_2 p_i$$

$$IG(D_p) = IG(D_p) - IG(D_l) \frac{N_l}{N_p} - IG(D_R) \frac{N_R}{N_p}$$

$$J = \text{error rate} + \lambda (\# \text{ of leaves})$$

$$e = \frac{\sum \text{of mistakes}}{\sum \text{of All}} \propto e^{-x}$$

$$x = \frac{1}{2} \ln \left(\frac{1-e}{e} \right)$$