

Dimensionality Reduction

why? \rightarrow to visualize high dimension data.

how? \rightarrow PCA & t-SNE.

* By default, a vector is column vector.

$x_i \in \mathbb{R}^d$ — column vector $d \times 1$.

* Representing dataset

$$D = \left\{ x_i, y_i \right\}_{i=1}^n \quad \text{--- data pts.} \quad x_i \in \mathbb{R}^d, \quad y_i \in \{setosa, virginica, versicolor\}$$

D is usually represented in 2 matrices — X & Y.

X — rows are data pts & cols are features,
Y — rows are data pts. usually 1 col. only.

Data pre-processing: Column normalization

why? — so that the data becomes nicely structured so that data modelling algs will perform well.

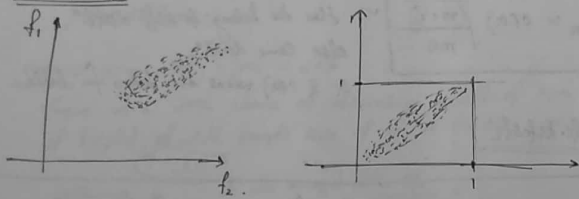
normalization: — for every feature, we get

$$a_i' = \frac{a_i - a_{\min}}{a_{\max} - a_{\min}}$$

$$a_i' \in [0, 1]$$

we become independent of units \rightarrow getting rid of scales.
 All in the same scale of $[0,1]$ — helps many algos become faster.

Geometrical view.



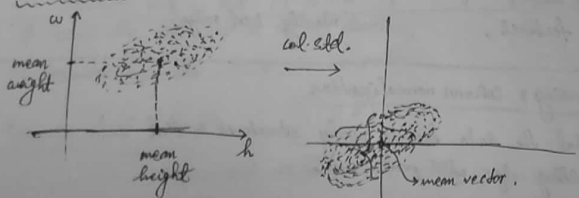
- * without 'spreading out' b/w the pts we are confining them to a small space.
- * In high dimension space, we squish data to hyper-cubes.

Column Standardization — more often used than col. normⁿ.

for every feature, $a_i' = \frac{a_i - \mu}{\sigma}$ μ = mean of feature
 σ = std devⁿ of feature.

So, now for each feature after standardization, mean = 0 & std-dev = 1.
 a_i can come from any distribⁿ, not necessarily gaussian.

Geometrical view



- ① moving the mean vector to origin
- ② squishing/expanding st. std dev for any feature = 1

Co-variance matrix

$$X = \begin{bmatrix} f_1 & f_2 & \dots & f_d \\ x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \text{def: } \Sigma = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} \quad \begin{matrix} n \times d \\ d \times d \end{matrix}$$

square matrix.
covariance matrix

$$\Sigma_{ij} = \text{cov}(f_i, f_j)$$

$i: 1 \rightarrow d$
 $j: 1 \rightarrow d$

$$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

Hence, the diagonal elements will be variance of features

$$\begin{aligned} \text{cov}(f_i, f_i) &= \text{var}(f_i) \\ \text{cov}(f_i, f_j) &= \text{cov}(f_j, f_i) \end{aligned}$$

& the matrix will be symmetric.

* Let X matrix is col. standardized.

$$\begin{aligned} \text{cov}(f_1, f_2) &= \frac{1}{n} \sum_{i=1}^n (x_{i1} - \mu_1)(x_{i2} - \mu_2) \\ &= \frac{1}{n} \sum_{i=1}^n x_{i1} * x_{i2} \rightarrow \text{multiplying the features & features values of corresponding data pts. \& taking summatⁿ.} \\ &= (f_1^T \cdot f_2) * \frac{1}{n} \end{aligned}$$

$$\Sigma = \frac{1}{n} (X_{d \times n}^T X_{n \times d}) \quad \text{can be found using simple dotⁿ. given } X \text{ is col. standardized.}$$

MNIST dataset (refer colab: github/20).

60k — training pts
 10k — test pts.

Data of handwritten digit grayscale photos of size 28 x 28 pixels.

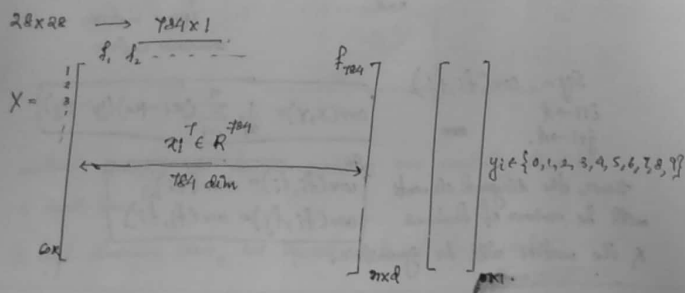
$$D = \{x_i, y_i\}_{i=1}^{60k}$$

$$x_i: \begin{matrix} \square \\ \updownarrow 28 \\ \square \end{matrix} \quad y_i \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

objective: - classify which of the class the image belongs to.

* how to go from image to vector?

→ grayscale images - 2D matrix. simply flattening by concatenating 1 row after other will do.



we'll visualize this in 2D

Principal Component Analysis (PCA)

dimensionality reduction d -dim to d' -dim. $d' < d$.

why → visualize data.

Geometric Interpretation

we'll get intuitions from examples reducing

d -dim → 1 dim & generalize to

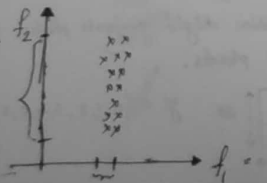
d -dim → d' -dim $d' < d$.

f_1 : blackness of hair

f_2 : height

spread of f_2 high

spread on f_1 very minimal.

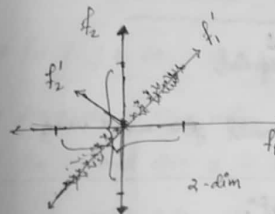


If we have to convert this data to 1-D, skipping the blackness feature is better option since ~~the~~ spread is a measure of data.

* So, we should preserve the dirⁿ with maximal spread.

Ex: X = 2-dim dataset.

say its al standardized. $\begin{cases} \text{mean}\{f_1\} = \text{mean}\{f_2\} = 0 \\ \text{var}\{f_1\} = \text{var}\{f_2\} = 1. \end{cases}$



spread on both f_1 & f_2 is significant.
but if we consider f_1' & f_2'
spread(f_2') \ll spread(f_1').
 $\{f_1' \& f_2'\}$ are \perp to each other.

So, basically we are rotating the given axis pair f_1 - f_2 by an angle θ such that the variance of f_1' is maximum.

& then drop f_2'

Variance maximizing optimization

Let u_1 be unit vector in dirⁿ of f_1' .

So, $\|u_1\| = 1$.

$$x_i' = \text{proj}_{u_1} x_i = \frac{u_1^T \cdot x_i}{\|u_1\|^2} = u_1^T \cdot x_i$$

$$\text{So, } \bar{x}_i' = u_1^T \cdot \bar{x}_i$$

mean of original pts.

mean of projected pts

we have to find u_1 s.t. $\text{var}\{\text{proj}_{u_1} x_i\}_{i=1}^n$ is maximum.

$$\text{var} \{ u_1^T x_i \}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n (u_1^T x_i - u_1^T \bar{x})^2 \quad \text{since col std. } [0, 0, 0, \dots, 0]$$

$$\text{var} \{ u_1^T x_i \}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n (u_1^T x_i)^2$$

$$\max_{u_1} \frac{1}{n} \sum_{i=1}^n (u_1^T x_i)^2 \quad \text{s.t. } u_1^T u_1 = 1 = \|u_1\|^2$$

objective of optimizatⁿ problem

constraints

If no constraint, then $u_1 = [a, 0]$ would give max^m value.

Alternative formulation:- distance minimizatⁿ.

Smaller formulatⁿ:- find u_1 which maximizes projected variance.

Alternative:-

for each x_i , we have d_i : distance of x_i from u_1

$$\min_{u_1} \sum_{i=1}^n d_i^2$$

$$\min_{u_1} \sum_{i=1}^n [(x_i^T x_i) - (u_1^T x_i)^2]$$

$$\text{s.t. } u_1^T u_1 = 1$$

$$\text{proj of } x_i \text{ on } u_1 = u_1^T x_i$$

$$d_i^2 = \|x_i\|^2 - (u_1^T x_i)^2 = (x_i^T x_i) - (u_1^T x_i)^2$$

Both the optimizatⁿ problems are the same.

Solution to our optimization problem:-

S_{cov} = co-variance matrix. will give d-pairs of (eigen value, eigen vector)

$$\begin{matrix} \lambda_1, \lambda_2, \lambda_3 & \dots & \lambda_d \\ \downarrow & & \downarrow \\ v_1 & v_2 & v_3 & \dots & v_d \end{matrix} \quad (\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d)$$

defⁿ of E-vectors: $\lambda_i v_i = S v_i$ - all pair of λ_i, v_i will satisfy this relationship.

* for v_i, v_j $i \neq j$, $v_i \perp v_j \Rightarrow v_i^T v_j = v_i \cdot v_j = 0$

* Every pair of eigen vectors are perpendicular.

* Our required dirⁿ u_1 is basically v_1 (vector corresponding to largest eigen val).

Geometric interpretation of λ_i & v_i

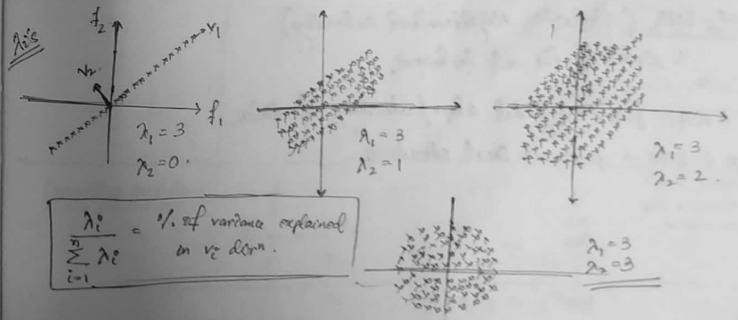
Say $x_i \in \mathbb{R}^{10}$ then dimensions = 10.

$$\begin{matrix} \lambda_1 \geq \lambda_2 \geq \lambda_3 & \dots & \lambda_{10} \\ \downarrow & & \downarrow \\ v_1 & v_2 & v_3 & \dots & v_{10} \end{matrix}$$

↑ directⁿ with max^m variance

↑ directⁿ with 2nd most variance

↑ directⁿ with least variance.



$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ \vdots & \vdots & & \vdots \\ x_1 & x_2 & \dots & x_n \end{bmatrix}_{n \times 10} \xrightarrow[\text{(PCA)}]{\text{dim}} X' = \begin{bmatrix} v_1 & v_2 \\ \vdots & \vdots \\ v_1 & v_2 \end{bmatrix}_{n \times 2}$$

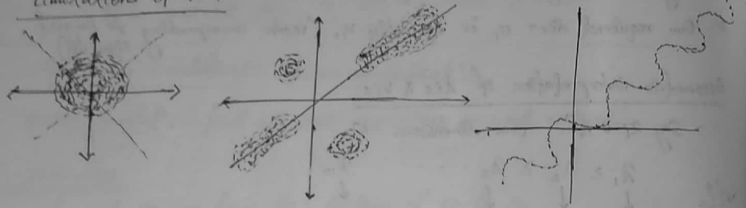
$$X' = [x_1^T v_1, x_1^T v_2]$$

$$v_1 \& v_2 \text{ --- top 2 eigen vectors}$$

* Preserve 99% of variance.

Then let $\lambda_1 + \lambda_2 + \dots + \lambda_k = 0.99$ So, k-dimensions used.

Limitations of PCA



Suppose we have to reduce from 184 dim to 200 dim.

$$X_{184 \times 784} \times V_{784 \times 200} = X'_{184 \times 200}$$

$$V = \begin{bmatrix} v_1 & v_2 & \dots & v_{200} \\ \vdots & \vdots & & \vdots \end{bmatrix}$$

t-SNE (Stochastic Neighbourhood Embedding)

↳ state of the art technique.

* PCA :- preserves global shape/structure of data.

* t-SNE → preserves local structure.

Neighbourhood, Embedding

d-dim space (high).

$$\|x_i - x_j\|^2 = \text{dist}^2$$



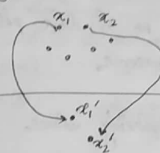
$N(x_i) = \{x_j \mid x_i, x_j \text{ are geometrically close}\}$

$N(x_1) = \{x_2, x_3, x_4, x_5, x_6, x_7, x_8\}$
doesn't contain $x_9, x_{10}, x_{11}, \dots$

Embedding (mapping)

d-dim

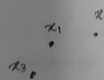
2-dim



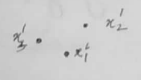
for every point in the high dimensional space if we are finding corresponding pt in low dimensional space, such a thing is called embedding.

* t-SNE preserves distance of pts in neighbourhood.

d-dim



2-dim



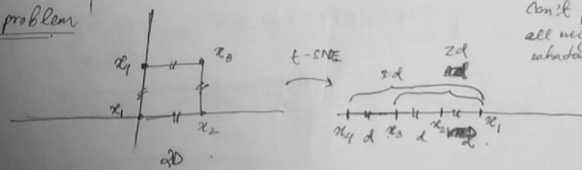
$$d(x_1, x_2) \approx d(x'_1, x'_2)$$

$$d(x_1, x_3) \approx d(x'_1, x'_3)$$

$$d(x_1, x_5) \approx d(x'_1, x'_5)$$

$$d(x_1, x_4) \neq d(x'_1, x'_4)$$

Crowding problem



Can't preserve all neighbours whatever we do

So, sometimes it becomes impossible to preserve distance of all the neighbours.

So, use t -distribution to resolve this. do its best to minimize errors.

Analysing t -SNE (distill.pub).

- * Iterative algo. — goes through all data pts on iterations.
- * Probabilistic algo — can give diff o/p on multiple runs.
- * perplexity can be loosely thought of no. of data pts in neighbourhood of each pt we are trying to preserve.
- * Always run t -SNE with multiple perplexity values.
- * always perplexity $<$ #data pts.
- * Always check whether shape has stabilized by running it for more no of times.
- * cluster size (more tightly packed, more diffused) doesn't say anything.
- * distance b/w clusters don't mean anything.