# Second task of unsupervised learning
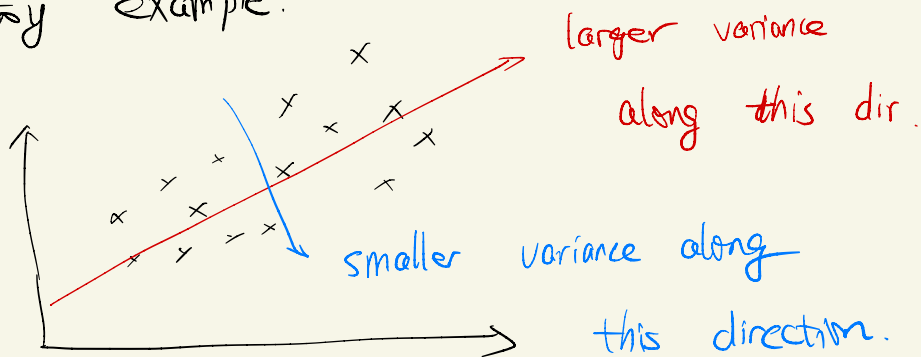
## dimensionality reduction

goal: reduce high dimensional data to low-D.

classical method: PCA: principal component analysis

- Given data $x_1, x_2, \ldots x_n \in \mathbb{R}^d$.

- Would like to have a reduced-dimension representation $y_1, y_2, \ldots y_n \in \mathbb{R}^l$ with $l << d$. such that important info is kept.

---

A toy example.



larger variance along this dir.

smaller variance along this direction.

$\longrightarrow$ project onto direction with larger variance.

# First interpretation of PCA

$\longrightarrow$ maximize the variance of the reduced data

assume data is centered i.e.

$$\sum_{i=1}^{n} x_i = 0.$$

Want to find a <u>direction</u> $u$.

$$u \in \mathbb{R}^d \qquad \|u\|_2 = 1.$$

s.t. the projections

$$\left\{ u^\top x_i \right\}_{i=1}^{n} \qquad \text{have max variance.}$$

$\Longleftarrow$ $\max\limits_{u:\, \|u\|_2 = 1} \sum\limits_{i=1}^{n} (x_i^\top u)^2.$

$\Longleftarrow$ $\max\limits_{u:\, \|u\|_2 = 1} \underbrace{u^\top \sum\limits_{i=1}^{n} x_i x_i^\top\, u}$

$\underset{n}{\triangleq} X X^\top,$ where

covariance matrix

$$X = d \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_n \\ | & | & & | \end{bmatrix}$$

This gives us the first PC. (principal component)

$$u_1 = \underset{u : \|u\|_2 = 1}{\arg\max} \sum_{i=1}^{\Lambda} (x_i^T u)^2$$

How to obtain the second PC ??!

$$u_2 = \underset{\substack{u : \|u\|_2 = 1 \\ u^T u_1 = 0}}{\arg\max} \sum_{i=1}^{\Lambda} (x_i^T u)^2$$

$\hookrightarrow$ the 2nd eigenvector

of $XX^T$.

$\vdots$    can extend all the way to $\underline{k\ PCA}$.

---

Algorithm for PCA.

· center the data set.

· compute eigen-decomposition of $XX^T = U\Sigma U^T$.

· return $U = [\ u_1\ u_2\ \cdots\ u_k\ u_{k+1}\ \cdots\ u_d\ ]$.

$\underbrace{\phantom{u_1\ u_2\ \cdots\ u_k}}$
top-$k$ eigenvectors.

What's the low-D representation of x ???

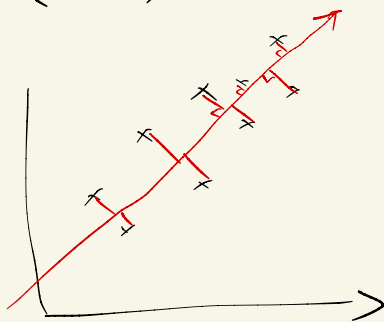$$x \rightarrow z = [\ u_1^T x \ , u_2^T x, \ \cdots \ u_k^T x\ ] \in \mathbb{R}^k.$$

___

second interpretation: minimizing reconstruction error.

$$x \rightarrow (u^T x)\ u \qquad \text{for any } \|u\|_2 = 1.$$

would like to have small reconstruction error.

$$\| x - (u^T x)\ u) \|_2^2 \quad \text{is} \quad \text{small}.$$

$$\rightarrow \qquad \underset{u:\ \|u\|_2 = 1}{\arg\min} \ \sum_{i=1}^{n} \| x_i - (u^T x_i)\ u \|_2^2$$



Claim: this yields exactly the same

direction as $u_1$.

more generally, we have.

$$\min_{u_1, u_2, \dots u_k} \sum_{i=1}^{n} \left\| x_i - \sum_{j=1}^{k} (u_j^\top x_i) u_j \right\|_2^2$$
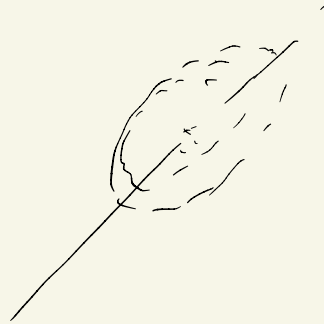
$$s.t. \quad u_i^\top u_i = 1 \quad \boxed{u_i^\top u_j = 0 \quad i \neq j}$$

$$\forall i.$$

---

# How to choose k ??!

- choose k s.t. the PCs have low reconstruction error.

- cross validation on downstream tasks.

---

PCA likes Gaussian data.

dislikes other structured data.



$\longrightarrow$ kernel PCA.

Advanced: random projections
and Johnson-Lindenstrauss lemma.

Given $x_1, \ldots x_n \in \mathbb{R}^d$ $\rightarrow$ hard to store and manipulate.

construct a mapping $\pi: \mathbb{R}^d \rightarrow \mathbb{R}^k$  $k \ll d$.

s.t, all distances are nearly preserved, i.e.

$$\|x_i - x_j\|_2 \approx \|\pi(x_i) - \pi(x_j)\|_2 .$$

$\underbrace{\phantom{\|x_i - x_j\|_2}}_{\mathbb{R}^d} \qquad \underbrace{\phantom{\|\pi(x_i) - \pi(x_j)\|_2}}_{\mathbb{R}^k}$

example: $\underline{k-\text{means clustering}}$.

More precisely: we want to achieve.

$$1 - \varepsilon \leq \frac{\|\pi(x_i) - \pi(x_j)\|_2}{\|x_i - x_j\|_2} \leq 1 - \varepsilon \qquad *$$

for some small $\varepsilon$ say $\varepsilon = 0.001$.

Lemma (Johnson-Lindenstrauss 1984).

As long as $k > \dfrac{4 \log n}{\dfrac{\varepsilon^2}{2} - \dfrac{\varepsilon^3}{3}}$

for any set of data points in $\mathbb{R}^d$.

there exists a map s.t. (*) is true

# Remarkable property:

$k$ is dimension independent.

only depends $log$ on $\underline{\underline{n}}$.

In fact: you can achieve $*$ simply by random projection:

let $W \in \mathbb{R}^{k \times d}$ be Gaussian random matrix

define $\pi(x) = \dfrac{Wx}{\sqrt{m}}$. this "almost always" works.

why ??! fix some $x$.

① $\mathbb{E} \| \pi(x) \|_2^2 = \mathbb{E} \dfrac{\| Wx \|_2^2}{m} = \| x \|_2^2$.

② $\| \pi(x) \|_2^2$ concentrates well around $\| x \|_2^2$.

③ we only need this to hold for $n^2$ pairs.

Extensions

_____. other distances.