# Lecture 7: Principal Component Analysis

Q: How do we find salient directions in data?

$$X_1, \dots, X_n \in \mathbb{R}^d \quad, \quad d \text{ large.}$$

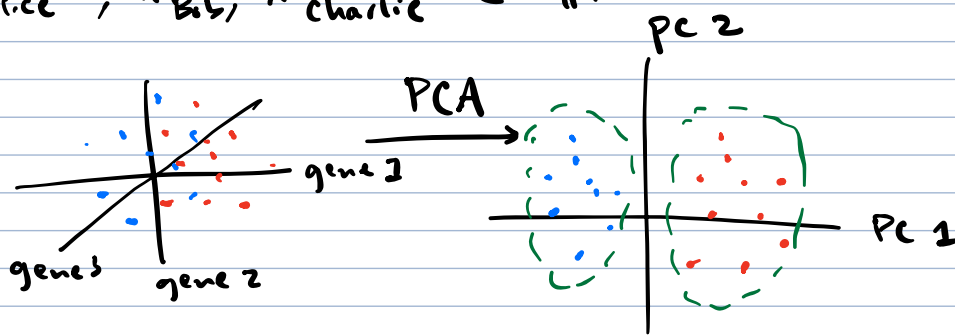Q: How do we find the "best" low dimensional representation of data?

Many applications:

- data visualization
- discovery
- data clustering.
- ...

e.x. genetics data

|  | gene 1 | gene 2 | .... | gene d |
|---|---|---|---|---|
| Alice | 0 | 0 | 1 | 0 ... 0 |
| Bob | 1 | 0 | 0 | 1 ... 0 |
| Charlie | 0 | 0 | 1 | 0 ... 0 |
| ⋮ |  |  |  |  |

1 if Bob has a mutation at position 1.

$$X_{Alice}, X_{Bob}, X_{Charlie} \in \mathbb{R}^d$$



Toy example:

4 people rate foods from 1-10.

|  | kale | taco bell | sashimi | pop tarts |
|---|---|---|---|---|
| Alice | 10 | 1 | 2 | 7 |
| Bob | 7 | 2 | 1 | 10 |
| Carol | 2 | 9 | 7 | 3 |
| Dave | 3 | 6 | 10 | 2 |

Q: How do we visualize this data?

step 1: center the data

$$\mu = (5.5, 4.5, 5, 5.5)$$

step 2: find 2 good directions $v_1, v_2$ for the data s.t.

$$x - \mu \approx a_1 v_1 + a_2 v_2 \quad \forall \ x \in \{ \text{Alice, Bob, Carol, Dave} \}.$$
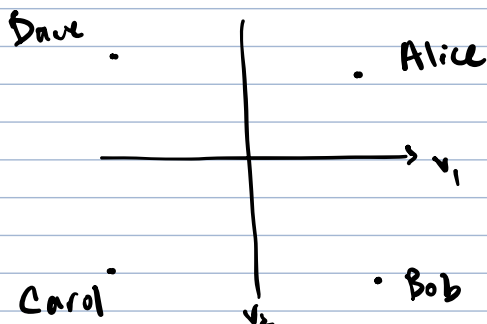
turns out, if you take

$$v_1 = (3, -3, -3, 3)$$
$$v_2 = (1, -1, 1, -1),$$

then

$$
\begin{aligned}
\text{Alice} - \mu &\approx v_1 + v_2 &&\rightarrow (1, 1) \\
\text{Bob} - \mu &\approx v_1 - v_2 &&\rightarrow (1, -1) \\
\text{Carol} - \mu &\approx -v_1 - v_2 &&\rightarrow (-1, -1) \\
\text{Dave} - \mu &\approx -v_1 + v_2 &&\rightarrow (-1, 1)
\end{aligned}
$$

e.g. $v_1 + v_2 = (4, -4, -2, -2)$

$\text{Alice} - \mu = (4.5, -3.5, -3, -1.5)$   pretty close!



big $v_1 \rightarrow$ like kale, pop-tarts, dislike TB + sashimi

big $v_2 \rightarrow$ like kale, sushi, dislike TB + pop-tarts

we can use this to infer more properties of their food prefs!

# Principal Component Analysis

Given $X_1, \dots, X_n \in \mathbb{R}^d$

Typically, de-mean them. Let $\mu = \frac{1}{n} \sum_{i=1}^{n} X_i$, set $X_i' = X_i - \mu$, so that new mean is $(0, \dots, 0)$. So to slightly simplify notation, let's just work with de-meaned data, ie. let's assume $\mu = 0$.

Goal: Find a __subspace__ $V \subseteq \mathbb{R}^d$, $\dim(V) = k$

parameter.

so that $X_i \approx \text{proj}_V(X_i)$.
Want $k << d$ (often constant).

More concretely:

$V = \text{span}\{v_1, \dots, v_k\}$. We can choose $v_1, \dots, v_k$ __orthonormal__

$$\langle v_i, v_j \rangle = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases} \quad \leftarrow \|v_i\|_2^2 = 1.$$

PCA objective:

$$\underset{\substack{v_1, \dots, v_k \in \mathbb{R}^d \\ \text{orthonormal}}}{\arg\max} \sum_{i=1}^{n} \sum_{j=1}^{k} \langle X_i, v_j \rangle^2 \quad \leftarrow \|\text{proj}_V(X_i)\|_2^2$$
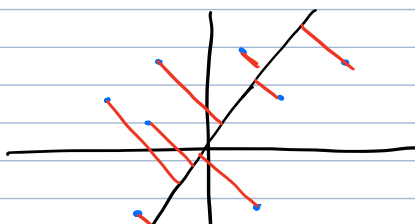
orthonormal basis.

Interpretation: For $V = \text{span}(\{v_1, \dots, v_k\})$, the projection of $X \in \mathbb{R}^d$ onto $V$ is

$$\text{proj}_V(X) = \sum_{j=1}^{k} \langle X, v_j \rangle \cdot v_j,$$

$$\|\text{proj}_V(X)\|_2^2 = \langle \text{proj}_V(X), \text{proj}_V(X) \rangle$$
$$= \sum_{j, \ell} \langle X, v_j \rangle \langle X, v_\ell \rangle \langle v_j, v_\ell \rangle = \sum_{j=1}^{k} \langle X, v_j \rangle^2$$

PCA: What is the $k$-dimensional subspace that explains the most variance in the dataset?

e.g. $k = 1$.

Given principal components $v_1, \dots, v_k$, we can approximate data points with their projection:

$$X_i \approx \text{proj}_V(X_i) = \sum c_{ij} v_j$$

Since $v_j$ are orthonormal basis, we can rewrite the projection in this basis as a $k$-dimensional vector

$$\text{proj}_V(X_i) = \begin{pmatrix} c_{i1} \\ c_{i2} \\ \vdots \\ c_{ik} \end{pmatrix}$$

← component of $X_i$ along the 2nd PC.

Some structural facts:

The PCs are not always unique!

However, if they are unique, then the solutions are nested!

$k=1 \longrightarrow \text{span}(\{v_1\})$

$k=2 \longrightarrow \text{span}(\{v_1, v_2\})$

$k=3 \longrightarrow \text{span}(\{v_1, v_2, v_3\})$

we'll see why next lecture!

Next lecture: there are efficient algorithms for PCA using connections to singular value decomposition.

## Relationship to Johnson Lindenstrauss

JL also gives a low-dimensional representation of data.

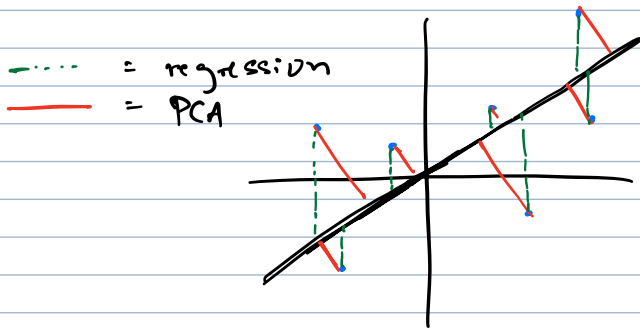| PCA | JL |
|---|---|
| - doesn't preserve distances | - preserves distances |
| - data dependent | - data oblivious |
| - PCs are meaningful | - JL directions are random → not meaningful. |

## Relationship with linear regression

Regression is a way to explain one dependent variable using data.

$$(\underbrace{x_1}_{\in \mathbb{R}^b}, \underbrace{y_1}_{\in \mathbb{R}}), \cdots, (x_n, y_n)$$
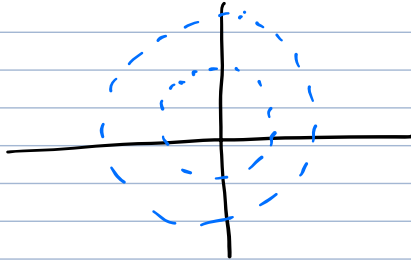
$$\underbrace{\qquad\qquad\qquad}_{\in \mathbb{R}^{d+1}}$$

$$y \approx \langle \theta, x \rangle$$

even in 2D is a bit different.



$-----$ = regression
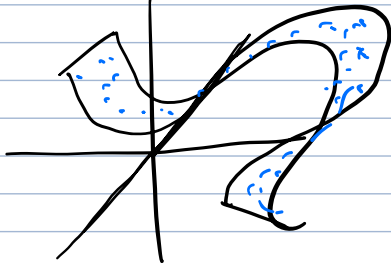$\textcolor{red}{-----}$ = PCA

---

Failure modes of PCA

Main issue: can only discover linear structure.



---

A natural idea: kernelize data!

Related concept: "manifold learning"

"nonlinear dimension reduction"



| Another visualization tool

t - SNE

not linear but tries to find low-d representation that "looks" like original data.