# PCA & ICA

+ RECAP PCA & SOLVE IT
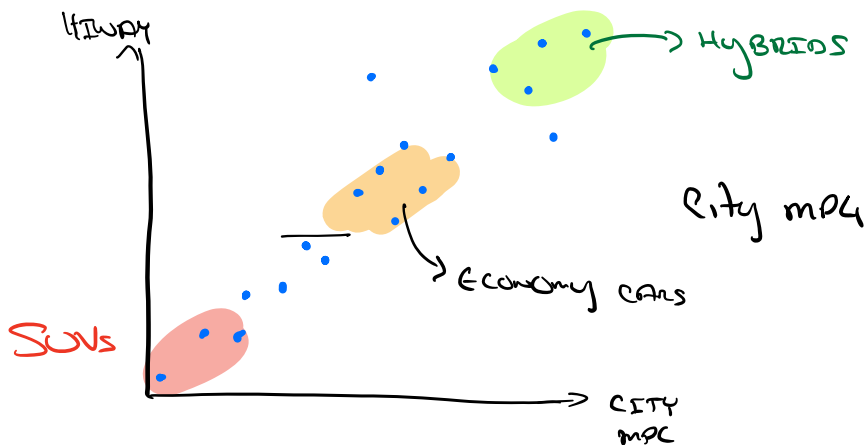
+ ICA & "the Cocktail Party"     "LEARN UPTO Symmetry"

UP NEXT: Self-Supervised MACHINE LEARNING!

# PCA: Principal Component Analysis

STRUCTURE

|  | Prob. | Non Prob |
|---|---|---|
| "CLUSTER" | GMM | K-MEANS |
| "SUBSPACE" | FACTOR ANALYSIS | PCA ← THIS SECTION |

Ex: GIVEN PAIRS (HIWAY mpg, City mpg) of SOME CARS

HIWAY

→ HYBRIDS

City mpg

→ ECONOMY CARS

SUVs

CITY mpg

Question: "GOOD mpg"

① CENTER DATA

$$\mu = \frac{1}{x} \sum_i x^{(i)}$$

$$x^{(i)} \mapsto x^{(i)} - \mu$$

HIWAY

$u_1$

Component of Principal VARIATION $(U_1)$

x

CITY

$2^{nd}$ Component of VARIATION $(U_2)$

$u_2$

Now $\|u_1\| = \|u_2\| = 1$ by convention.

• $U_1$ IS "HOW good IS mpg"

• $U_2$ IS "difference between HIWAY & city" (roughly)

WE CAN WRITE $x = \alpha_1 U_1 + \alpha_2 U_2$

↳ WE may just keep this component

"Explains more variation"

How we find these directions, and some caveats
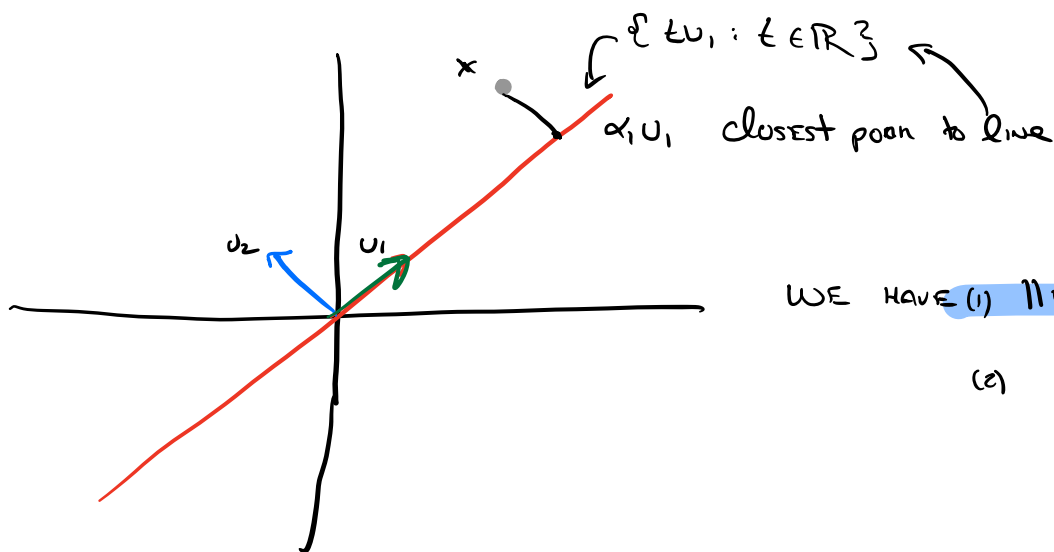- think about 1000s of dims ⟶ 10s of dims
- A dimensionality reduction method

## Preprocessing

Given $x^{(1)} \dots x^{(n)} \in \mathbb{R}^d$

1. Center the data $x^{(i)} \mapsto x^{(i)} - \mu$ in which $\mu = \frac{1}{n} \sum_i x^{(i)}$

2. May need to rescale components e.g. "feet per gallon"
$$\frac{?}{} \text{m pg}$$

We will assume data is preprocessed

## PCA as Optimization



$\{tu_1 : t \in \mathbb{R}\}$

$\alpha_1 u_1$ closest point to line

We have (1) $\|u_i\| = 1$  (unit vectors)

(2) $u_i \cdot u_j = \delta_{ij}$  (orthogonal)

How do you find closest point to the line?
$$\alpha_1 = \arg\min_\alpha \|x - \alpha u_1\|^2$$
$$= \arg\min_\alpha \|x\|^2 + \alpha^2 \|u\|^2 - 2\alpha (u_1 \cdot x)$$

differentiate wrt $\alpha$        $2(\alpha - u_i x) = 0 \implies \alpha = u_i \cdot x$

Generalize: $U_1 \cdots U_k \in \mathbb{R}^d$ AND $x \in \mathbb{R}^d$, USE $U_i \cdot U_j = \delta_{ij}$

$$\underset{\alpha_1 \cdots \alpha_d}{\text{Argmin}} \|x - \sum_{i=1}^{k} \alpha_i U_i\|^2 = \underset{\alpha}{\text{argmin}} \|x\|^2 + \sum_{i=1}^{k} \alpha_i^2 \|U\|^2 - 2\alpha_i \langle U_i \cdot x \rangle$$

Hence $\alpha_i = U_i \cdot x$

WE call $\|x - \sum_{i=1}^{k} \alpha_i U_i\|^2$ the $\underline{\text{RESIDUAL}}$

WE CAN find PCA by either

In class ① MAXIMIZE Projected Subspace

② MINIMIZE Residual

$$\underset{\substack{U \in \mathbb{R}^d \\ \|U\| = 1}}{\text{MAX}} \frac{1}{n} \sum_{i=1}^{n} (U \cdot x^{(i)})^2 \qquad \text{WE NEED some facts}$$

to solve this

LET A be symmetric & square, then

$$A = U \Lambda U^T \quad \text{IN which}$$

- $UU^T = I$ (ORTHONORMAL)

- $\Lambda$ is diagonal

$\Lambda_{ii} = \lambda_i$ AND $\lambda_1 \geq \cdots \geq \lambda_n$ by convention

eigenvalues

$\underline{\text{Recall}}$: If $x = \sum_{i=1}^{n} \alpha_i U_i$ where $[U_1 \cdots U_n] = U$

STANDARD BASIS vector

$$Ax = U \Lambda U^T x = U \Lambda \sum_{i=1}^{n} \alpha_i e_i \qquad (U_i \cdot U_j = \delta_{ij})$$

$$= U \sum_{i=1}^{n} \lambda_i \alpha_i e_i \qquad \text{diagonal } \Lambda$$

$$= \sum \lambda_i \alpha_i U_i$$

If $x = c U_i$ then $x$ is AN eigenvector, AND $Ax = \lambda_i x$

$$\max_{x: \|x\|^2 = 1} x^T A x = \max_{\alpha: \|\alpha\|^2 = 1} \sum_{i=1}^{n} \alpha_i^2 \lambda_i$$

Hence, we set $\alpha_i = 1$, the principal eigenvalue

Which $x$ attains it? If $\lambda_1 = \lambda_2$?

## Now, back to PCA!

$$\max_{U: \|U\|^2 = 1} \frac{1}{n} \sum_{i=1}^{n} (U \cdot x^{(i)})^2 \quad \longrightarrow \quad \text{THE projection onto } U$$

$$U \in \mathbb{R}^d$$

$$= \frac{1}{n} \sum_{i=1}^{n} U^T x^{(i)} (x^{(i)})^T U = U^T \left( \frac{1}{n} \sum_i x^{(i)} (x^{(i)})^T \right) U$$

$\longrightarrow$ Covariance of data (we subtracted mean)

$\therefore$ $U$ is principal Eigenvector

What if we want more dimensions? We keep top-$k$

## How do we represent data?

$$x^{(i)} \longmapsto \sum_{j=1}^{k} (x^{(i)} \cdot U_j) U_j$$

$\longrightarrow$ we keep these $k$ scalars (the $\alpha_k$ above)

A map from $\mathbb{R}^d \to \mathbb{R}^k$

## How do we choose $k$?

One approach "Amount of Explained Variance"

$$\frac{\sum_{i=r_k}^{k} \lambda_i}{\sum \lambda_i} \geq 0.9 \qquad \left( \text{Aside } tr(A) = \sum_i A_{ii} = \sum_i \lambda_i \right)$$

$j=1$

NB: ONLY MAKES SENSE IF $\lambda_j \geq 0$. HENCE COVARIANCE IS IMPORTANT

<u>**Lurking Instability**</u>: SUPPOSE $\lambda_k = \lambda_{k+1}$ ... WHAT HAPPENS?

REP IS <u>UNSTABLE</u> HERE

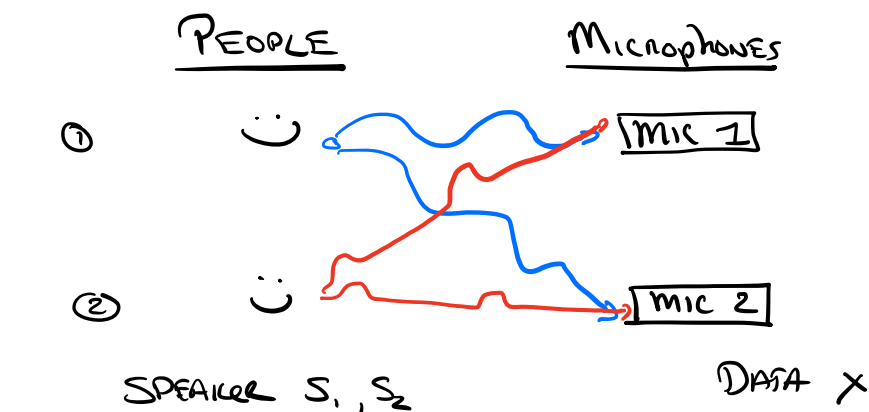# <u>Recap of PCA</u>

· Dimensionality Reduction technique (e.g. Visualization)

· MAIN IDEA IS TO project on a subspace, nice theory.

# <u>ICA</u>  INDEPENDENT Component Analysis
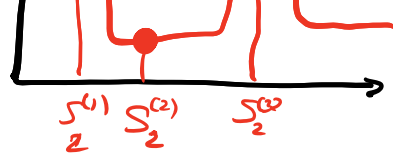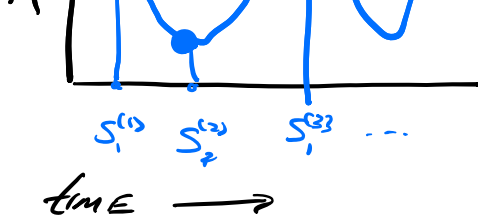
· high-level story

· Key facts & likelihood

· model

# <u>Cocktail Party Problem</u>  (IN HW!)

<u>PEOPLE</u>        <u>Microphones</u>

① ☺    |mic 1|

② ☺    |mic 2|

NB: WE SEE A MIXTURE of ~ & ~ AT EACH MIC

SPEAKER $S_1$, $S_2$        DATA $X$

INTENSITY        $S_1$        $S_2$

$S_1^{(1)} \quad S_2^{(2)} \quad S^{(3)} \cdots$

$S_2^{(1)} \quad S_2^{(2)} \quad S_2^{(3)}$

time $\longrightarrow$

$S_j^{(t)}$ IS __INTENSITY__ AT __TIME__ $t$ from __SPEAKER__ $j$

WE DO __NOT__ Observe $S^{(t)}$ only $x^{(t)}$ — the microphones

__EX model__  $X_j^{(t)} = a_{j1} S_1^{(t)} + a_{j2} S_2^{(t)}$

"MICROPHONE $j$ SEES A MIXTURE OF $S_1^{(t)}$ & $S_2^{(t)}$ "

__OR__    $\underset{\text{OBSERVED}}{X^{(t)}} = A \underset{\text{LATENT}}{S^{(t)}}$

$\nearrow$ LATENT

for simplicity, ASSUME # of SPEAKERS = # of mics $= d$

__GIVEN__:  $X^{(1)}, X^{(2)}, \ldots, X^{(n)} \in \mathbb{R}^d$   $d$ is # of microphones & SPEAKERS

__DO__:  find $S^{(1)}, \ldots, S^{(n)} \in \mathbb{R}^d$

AND $A \in \mathbb{R}^{d \times d}$  s.t.  $X^{(t)} = A S^{(t)}$

WE call $A$ the __MIXING MATRIX__ AND $W = A^{-1}$ __UNMIXING MATRIX__

WRITE $W = \begin{bmatrix} W_1^T \\ \vdots \\ W_d^T \end{bmatrix}$  so that  $S_j^{(t)} = W_j \cdot X^{(t)}$

__SOME CAVEATS__

· WE ASSUME $A$ does __NOT__ VARY w/ TIME AND IS __full RANK__

- THERE ARE INHERENT AMBIGUITY
  - WE CAN'T DETERMINE SPEAKER ID (cald SWAP 1 & 2)
  - CAN'T DETERMINE ABSOLUTE INTENSITY

$$(cA)(c^{-1}S^{(t)}) = AS^{(t)} \text{ for any } c \neq 0$$

- <u>SUPRISING</u> SPEAKERS <u>CANNOT</u> be GAUSSIAN

SUPPOSE SO $\quad x^{(t)} \sim N(\mu, AA^T) \quad$ then if $U^TU = I \quad AU$ GENERATES

the <u>SAME</u> DATA.

NEVERTHELESS, WE CAN RECOVER SOMETHING MEANINGFUL!
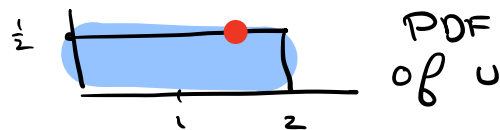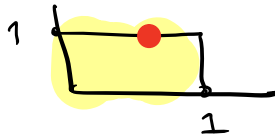
<u>Algorithm</u> : JUST MLE, SOLVED BY GRAD DESCENT

<u>DETOUR</u> : Density under linear transform (Key Confusion)

<u>Ex</u>: $S \sim$ Uniform $[0,1] \quad U = 2S \quad$ WHAT IS PDF of $U$?

TEMPTED TO WRITE $\quad P_U(\frac{x}{2}) = P_S(x)$

PDF of S



PDF of U

$$P_S(x) = \begin{cases} 1 & \text{if } x \in [0,1] \\ 0 & \text{o.w} \end{cases} \qquad P_U(x) = P_S(\frac{x}{2}) \cdot \frac{1}{2}$$

THE key ISSUE IS the NORMALIZATION CONSTANT

FOR INVERTIBLE MATRIX $A, \quad U = As$

$$P_U(x) = P_S(A^{-1}x) \, |\det(A^{-1})|$$

CHANGE OF VAR formula for INTEGRALS

$$= P_S(Wx) \cdot |\det(W)| \qquad (\frac{1}{\det(A)} = \det(A^{-1}))$$

FROM HERE ICA IS MLE:

$$P(s) = \prod_{j=1}^{d} P_S(s_j)$$

$$P(x) = \prod_{j=1}^{d} P_S(W \cdot x) \cdot |det(w)| \quad (\text{USE linear transform Abve})$$

Now WRITTEN IN TERMS of $x$ AND $A$.

Key technical bit : USE non-ROTATONLLY INVARIANT distribution

SET $P_S(k) \propto g'(x)$ for $g(x) = (1 + e^{-x})^{-1}$

Solve $\ell(w) = \sum_{t=1}^{n} \sum_{j=1}^{d} \log g'(w_j \cdot x^{(t)}) + \log |det(w)|$

- $\log |det(w)|$
- USE gD & you're done!

RECAP: · SAW PCA. WORKHORSE dimensionality Reduction
· ICA. Key IDEAS for HW. Introduce "up to symmetry"