# PCA: Principal Component Analysis

| STRUCTURE | PROB. | NON PROB |
|-----------|-------|----------|
| "CLUSTER" | GMM | K-MEANS |
| "SUBSPACE" | FACTOR ANALYSIS | PCA → TODAY'S lecture |

**Ex:** GIVEN PAIRS (HIWAY MPG, CITY MPG) of SOME CARS
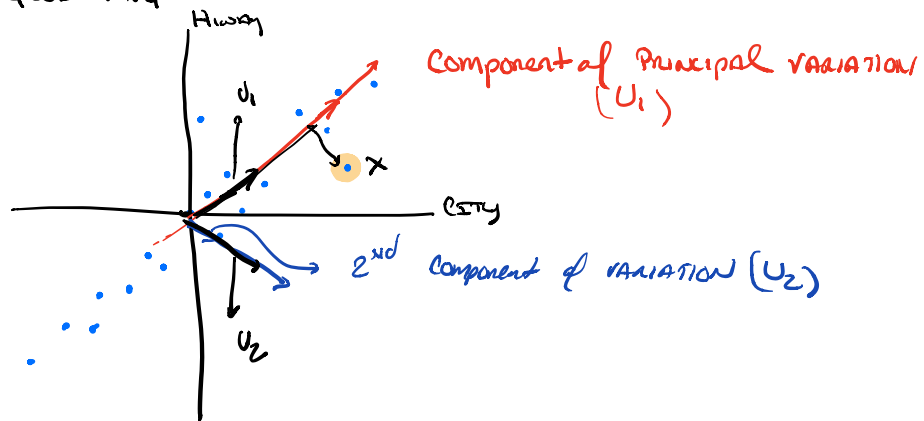


HIWAY

→ HYBRIDS

City mpg

ECONOMY CARS

SUVs

CITY MPC

**Question:** "GOOD MPG"

① CENTER DATA

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

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



HWAY

$U_1$

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

**TODAY:** How we find these directions, and some caveats

- think about 1000s of dims ⟶ 10s of dims
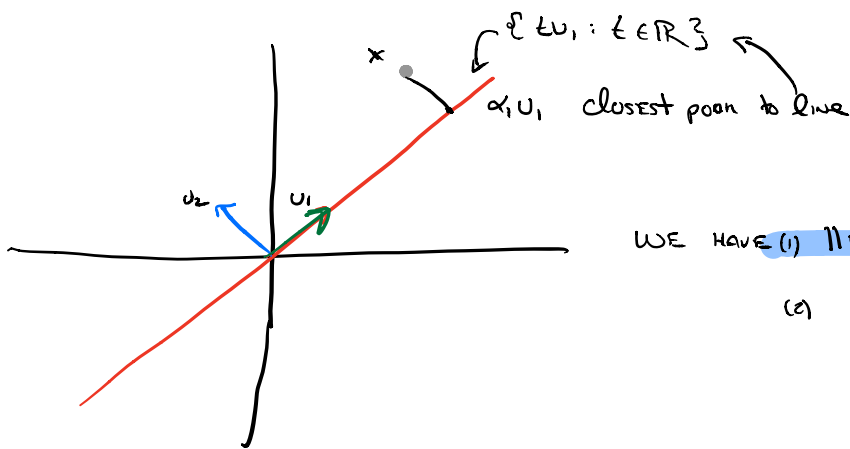- A dimensionality reduction method

## Preprocessing

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

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

2. May need to rescale components e.g. "feet per gallon" $\updownarrow$ MPG

We will assume data is preprocessed

## PCA as optimization



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

$\alpha_1 u_1$ 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?

$$\alpha_1 = \underset{\alpha}{\arg\min} \| x - \alpha u_1 \|^2$$

$$= \underset{\alpha}{\arg\min} \| x \|^2 + \alpha^2 \| u_1 \|^2 - 2\alpha (u_1 \cdot x)$$

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

$\underline{\text{Generalize}}$: $\quad u_1 \cdots u_k \in \mathbb{R}^d \quad$ AND $\quad x \in \mathbb{R}^d \quad$, USE $u_i \cdot u_j = \delta_{ij}$

$$\underset{\alpha_1 \cdots \alpha_d}{\text{Argmin}} \; \|x - \sum_{i=1}^{k} \alpha_i v_i\|^2 \;\; \overset{\bullet}{=} \;\; \underset{\alpha}{\text{argmin}} \; \|x\|^2 + \sum_{i=1}^{n} \alpha_i^2 \|u_i\|^2 - 2\alpha_i \langle u_i \cdot x \rangle$$

$\quad\quad$ Hence $\quad \alpha_i = u_i \cdot x$

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

WE CAN $\text{find}$ PCA $\text{by}$ either

In class $\quad$ ① $\quad$ MAXIMIZE $\quad$ Projected Subspace

$\quad\quad$ ② $\quad$ MINIMIZE $\quad$ Residual

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

LET A be symmetric & SQUARE, then

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

$\quad\quad \cdot \; UU^T = I \quad$ (ORTHONORMAL)

$\quad\quad \cdot \; \Lambda \;$ is diagonal

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

$\quad\quad\quad\quad\quad$ eigenvalues

$\underline{\text{Recall}}$: $\quad$ If $\; x = \sum_{i=1}^{n} \alpha_i v_i \;$ where $\; [u_1 \cdots u_n] = U$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ STANDARD BASIS vector

$$Ax = U\Lambda U^T x = U\Lambda \sum_{i=1}^{n} \alpha_i e_i \quad (u_i \cdot u_j = \delta_{ij})$$
$$= U \sum_{i=1}^{n} \lambda_i \alpha_i e_i \quad\quad \text{diagonal } \Lambda$$
$$= \sum_i \lambda_i \alpha_i v_i$$

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!

THE Projection onto $U$

$$\max_{\substack{U:\,\|U\|^2=1 \\ U \in \mathbb{R}^d}} \frac{1}{n} \sum_{i=1}^{n} (U_i \cdot x^{(i)})^2$$

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

→ COVARIANCE of DATA (WE SUBTRACTED MEAN)

∴ $U$ IS principal Eigenvector

WHAT if WE WANT MORE dimensions? WE keep top-$k$

How do WE REPRESENT DATA?

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

→ WE KEEP these $k$ scalars

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

How do WE CHOOSE $k$?

ONE APPROACH "Amount of Explained VARIANCE"

$$\frac{\sum_{i=1}^{k} \lambda_i}{\sum_{i=1}^{n} \lambda_i} \geq 0.9 \qquad \left( \text{ASIDE } tr(A) = \sum_i A_{ii} = \sum_i \lambda_i \right)$$

$j=1$

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

_Lurking_ INSTABILITY: SUPPOSE $\lambda_k = \lambda_{k+1} \ldots$ WHAT HAPPENS?
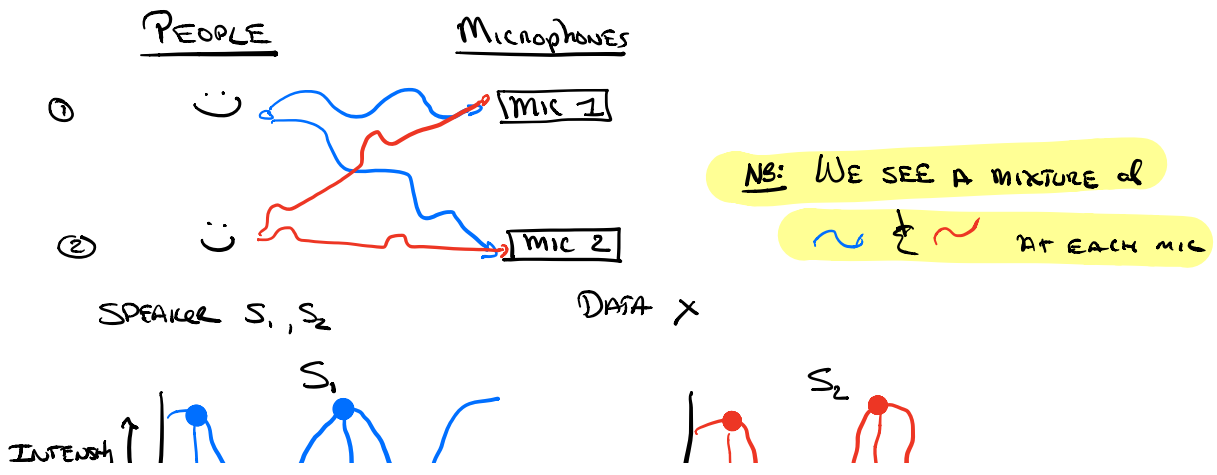
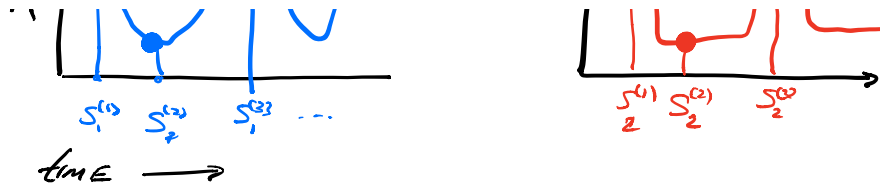REP IS UNSTABLE HERE

## RECAP of PCA

· Dimensionality REDUCTION technique (eg. Visualization)

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

## ICA   INDEPENDENT Component Analysis

· high-level story
· Key facts & likelihood
· model

## Cocktail Party Problem   (IN HW!)

PEOPLE          Microphones



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

① ☺          MIC 1

② ☺          MIC 2

SPEAKER $S_1, S_2$          DATA $X$

INTENSITY ↑          $S_1$          $S_2$

time $\longrightarrow$

$S_j^{(t)}$ IS UNDERLINED INTENSITY AT TIME $t$ from SPEAKER $j$

WE DO NOT OBSERVE $S^{(t)}$ only $x^{(t)}$ — the microphones

ex model  $\quad X_j^{(t)} = a_{j1} S_1^{(t)} + a_{j2} S_2^{(t)}$

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

OR

$X^{(t)} = A \, S^{(t)}$

OBSERVED $\quad\quad$ LATENT $\quad\quad$ LATENT

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

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

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

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

WE call $A$ the MIXING MATRIX AND $W = A^{-1}$ UNMIXING MATRIX

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

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

  $$(cA)(c^{-1}s^{(t)}) = As^{(t)} \quad \text{for ANY } c \neq 0$$

- <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!
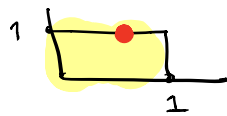
<u>Algorithm</u> : Just MLE, SOLVED By GRAD DESCENT

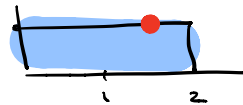<u>DETOUR</u> : Density under linear transform (Key Confusion)

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

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

PDF of S



PDF of U

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

THE key ISSUE IS the **NORMALIZATION constant**

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

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

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

**CHANGE of VAR** formula for **INTEGRALS**

From HERE ICA IS MLE:

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

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

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

Key technical bit: USE non-ROTATONALLY INVARIANT distribution

SET $P_S(x) \propto g'(x)$ 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."