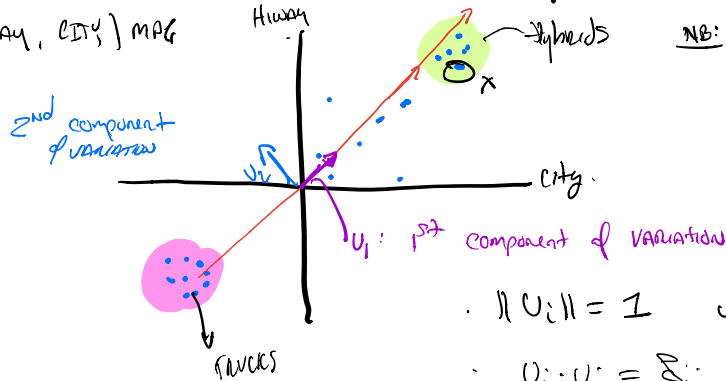# PCA & ICA

## PCA: Principal Component Analysis

GIVEN (HIWAY, CITY) MPG



NB: WE CENTERED DATA

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

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

- $\|U_i\| = 1$    unit length
- $U_i \cdot U_j = \delta_{ij}$    (orthonormal)

- $U_1$ — "How good is the mpg"

- $U_2$ — "variation in city/hiway from 'good'"

$$X = \alpha_1 U_1 + \alpha_2 U_2$$

$$X^{(i)} = \alpha_1^{(i)} U_1 + \alpha_2^{(i)} U_2$$

TODAY   HOW WE find directions

Think about dimension  1000s → 10

## Preprocess

GIVEN $x^{(i)} \dots x^{(n)} \in \mathbb{R}^d$

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

2. MAY NEED TO RESCALE components.. "FEET PER GALLON"
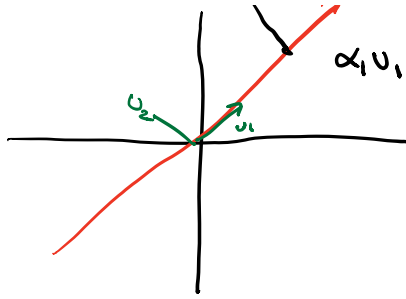   "MPG"
   "WHITEN" compute sample variance ...

   $$x_j^{(i)} \mapsto \frac{(x^{(i)} - \mu)_j}{\sigma_j}$$

WE will ASSUME DATA IS preprocessed.

## PCA AS OPTIMIZATION

$U_1$  $\{t U_1 : t \in \mathbb{R}\}$ line corresponds to $U_1$

How do you find closest point to the line generated by $u_1$?

$$\alpha_1 = \underset{\alpha}{\text{argmin}} \; \| x - \alpha u_1 \|^2$$

$$= \underset{\alpha}{\text{argmin}} \; \underbrace{\| x \|^2 + \alpha^2 \| u_1 \|^2 - 2\alpha (u_1 \cdot x)}_{g(\alpha)}$$

$$\nabla_\alpha g(\alpha) = 2\alpha - 2(u_1 \cdot x) = 0 \implies \alpha = u_1 \cdot x$$

<u>Generalize</u>: $u_1 \cdots u_k \in \mathbb{R}^d$ AND $x \in \mathbb{R}^d$ $\quad (u_i \cdot u_j = \delta_{i,j})$

$$\underset{\alpha_1 \cdots \alpha_d}{\text{Argmin}} \; \| x - \sum_{j=1}^{k} \alpha_j u_j \|^2 \cdots$$

Hence $\boxed{\alpha_j = u_j \cdot x}$

$$\| x - \sum_{j=1}^{k} \alpha_j u_j \|^2 \Leftarrow \text{RESIDUAL}$$

WE CAN find PCA by either:

1. MAXIMIZE Projected subspace

2. MINIMIZE Residuals

$$\underset{\substack{u \in \mathbb{R}^d \\ \| u \| = 1}}{\text{MAX}} \; \frac{1}{n} \sum_{i=1}^{n} (x^{(i)} \cdot u)^2$$

Solve this Optimization problem we need some facts.

Let A be symmetric & square then

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

$\cdot \ U U^T = U^T U = I$ (ORTHONORMAL BASIS)

$\cdot \ \Lambda$ is diagonal matrix

$\Lambda_{ii} = \lambda_i$ AND $\lambda_1 \geq \lambda_2 \cdots \geq \lambda_D$ by convention

$\hookrightarrow$ eigenvalues

Recall If $x = \sum_{i=1}^{d} \alpha_i U_i$ where $[U_1 \cdots U_d] = U$

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

$$= U \sum_{i=1}^{d} \alpha_i \lambda_i e_i \qquad \text{if } i = j \Rightarrow \delta_{ij} = 1$$

$$= \sum_{i=1}^{d} \alpha_i \lambda_i U_i \qquad \text{else } 0$$

fix $i$, and let $c \in \mathbb{R} \neq 0$

$x = c U_i$     eigenvectors     $Ax = \lambda_i x$

$$\max_{\substack{U \in \mathbb{R}^d \\ \|U\| = 1}} \frac{1}{n} \sum_{i=1}^{n} (x^{(i)} \cdot U)^2 \quad \Rightarrow \max_{\substack{U \in \mathbb{R}^d \\ \|U\| = 1}} U^T A U \qquad A = \frac{1}{n} \sum_{i=1}^{n} x^{(i)} x^{(i)T}$$

$$\Rightarrow \max_{\substack{\alpha \in \mathbb{R}^d \\ \|\alpha\|^2 = 1}} \sum_{i=1}^{d} \alpha_i^2 \lambda_i$$

WHAT SHOULD WE PICK TO MAXIMIZE? $\alpha_1 = 1 \quad \alpha_2 = \alpha_3 = \cdots = \alpha_d = 0$

IS IT UNIQUE? $\lambda_1 = \lambda_2$ WHAT HAPPENS? (PCA "loses insrability")

$\lambda_1 > \lambda_2 \Rightarrow$ UNIQUE

$U_1$ IS the PRINCIPAL eigenvector

WHAT IF WE WANT THE TOP-K such vectors?

$U_1 \cdots U_K$ BECAUSE $\lambda_1 \geq \cdots \geq \lambda_K$

$$x^{(i)} \mapsto \sum_{i=1}^{K} (x^{(i)} \cdot U_i) U_i$$

$J =$

$\quad\quad\quad \hookrightarrow$ keep these Knumbers.

# How do we pick k?

ONE Approach   "Amount of Explained Variance"

$$\frac{\sum\limits_{i=1}^{k} \lambda_i}{\sum\limits_{j=1}^{d} \lambda_j} \geq 0.9 \qquad\qquad \lambda_i \geq 0$$

<u>Lurking Instability</u>   $\lambda_{k-1} = \lambda_k = \lambda_{k+1}$   Are you were looky for top k?
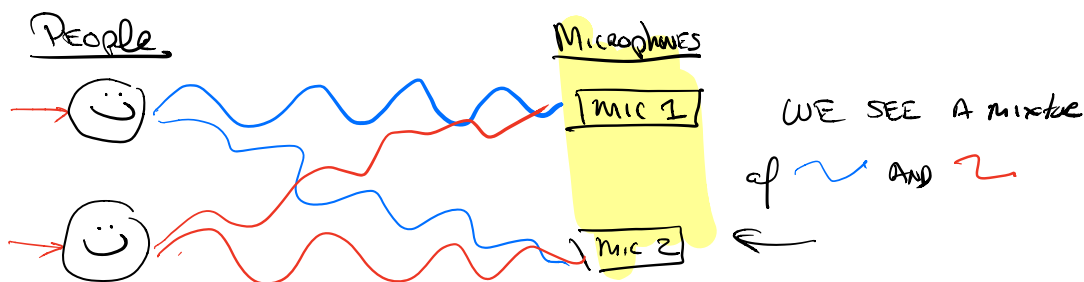
Pick any 2 of these! $\Rightarrow$ <u>different bases</u>

<u>PCA</u> · Dimensionality Reduction technique

· MAIN IDEA IS Project on Subspace

· Nice theory! Contrast w/ EM

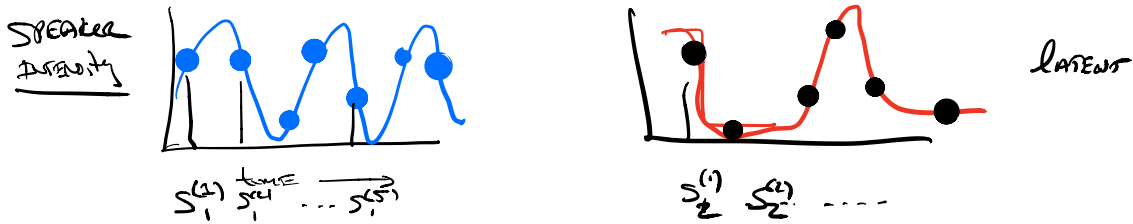<u>ICA</u>   Independent Component Analysis

· High-level Story

· Key facts AND likelihood

· model

<u>People</u>                                    <u>Microphones</u>



MIC 1              WE SEE A mixture

MIC 2              of ~ AND 2

SPEAKER SEPARATION.

SPEAKER $S_1^{(i)}$, $S_2^{(i)}$                    DATA $X_1^{(i)}$, $X_2^{(i)}$

SPEAKER
INTENSITY

latent

$S_1^{(1)} \xrightarrow{\text{time}} S_1^{(5)}$    $S_2^{(1)} \; S_2^{(4)} \; \dots$

$S_j^{(t)}$ IS the INTENSITY AT TIME $t$ of SPEAKER $j$

WE DO **NOT** OBSERVE $S_j^{(t)}$ only $X_j^{(t)}$ — the MICROPHONES

model    $X_j^{(t)} = a_{j1} S_1^{(t)} + a_{j2} S_2^{(t)}$ ← Hidden

"MICROPHONE $j$ SEES A MIXTURE of the SPEAKER INTENSITY"

OBSERVED $\longleftarrow$         $\longrightarrow$ LATENT

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

for simplicity, $\underline{d}$ is the number of SPEAKERS & MICROPHONES

**GIVEN:** $X^{(1)}, \dots, X^{(n)} \in \mathbb{R}^d$

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

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 $W = \begin{bmatrix} \omega_1^T \\ \vdots \\ \omega_d^T \end{bmatrix}$ so that $S_j^{(t)} = \omega_j \cdot X^{(t)}$

· WE CENTER the DATA $X^{(i)} \mapsto X^{(i)} - \mu \quad \mu = \frac{1}{n} \sum X^{(i)}$

· $A$ does not VARY w/ time, $A$ is full RANK

· THERE ARE SOME INHERENT Ambiguity:

   · WE CAN'T DETERMINE SPEAKER ID

   · CAN'T DETERMINE INTENSITY (ABSOLUTE)

$$x^{(i)} = As^{(i)}$$
$$= (cA)(c^{-1}s^{(i)}) \quad \text{for any } c \neq 0.$$

· <u>Surprising</u>  Speakers cannot <u>be Gaussian</u>

$$x^{(i)} = As^{(i)} \qquad s^{(i)} \sim N(0, I)$$

$$\Rightarrow \quad x^{(i)} \sim N(0, AA^T) \qquad UU^T = I$$
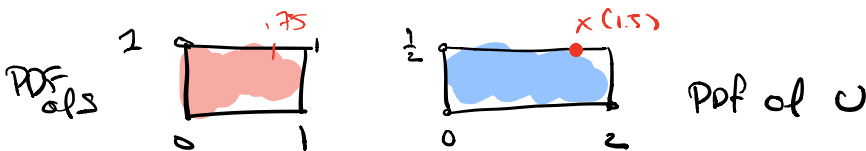
$$AUU^T A^T = AA^T$$

Nevertheless, we can recover something meaningful

<u>Algorithm</u>: Just MLE w/ Grad Descent.

<u>Detour</u> : Density under linear transform

Ex: S ~ uniform $[0,1]$    $U = 2S$    What's the PDF of $U$?

tempted  $P_U(x) = P_S(\frac{x}{2})$  (not right)



PDF of S

PDF of U

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

The key issue is the NORMALIZATION constant

A square & invertible ,    $U = As$    $s \sim$ PDF of $P_S$

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

$$= P_S(Wx) \, |\det(W)|$$

CHANGE of VARIABLES formula



volume  $vol(B)$    $vol(AB) = |\det(A)| \, vol(B)$

FROM HERE ICA IS MLE!

latent $\longrightarrow$ $P(s) = \prod_{j=1}^{d} P_S(s_j)$ " SOURCES INDEPENDENT

AND HAVE distribution "

OBSERVED $\longrightarrow$ $P(x) = \prod_{j=1}^{d} P_S(w_j \cdot x) \, |det(\omega)|$

Key technical track    NOT ROTATIONALLY Symmetric

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

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

· log det k

· USE GD & goolie done

RECAP: · SAW PCA. WORKHORSE dimensidy REDUCTION

· ICA - Key IDEAS. Handly symmetry.