

Factor Analysis  $\not\equiv$  PCA Continue unsupervised learning

+ Factor Analysis (many more dimensions than points  $d \gg n$ ) EM

+ PCA (old standby)

Linear Algebra

→ UP NEXT ICA "cocktail Party" (end of classical unsupervised)

## Factor Analysis

MANY fewer parts than dimensions "n < d"

cf: GMMs n ≈ d lots of neurons, few sources.

### How does this happen?

PLACE SENSORS ALL OVER CAMPUS, RECORD @ 100s of locations

$$\Rightarrow d \approx 100s$$

But Only record for 30 days (n < d)

WANT TO FIT A DENSITY but seems hopeless.

**KEY IDEA:** ASSUME THERE IS SOME LATENT r.v. THAT

IS NOT TOO COMPLEX AND EXPLAINS BEHAVIOR.

1<sup>st</sup> Let's see Problems w/ GMMs... Even 1 GAUSSIAN

$$\mu = \frac{1}{n} \sum_{i=1}^n x^{(i)} \rightarrow \text{this is OK}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

RANK( $\Sigma$ ) ≤ n < d - NOT FULL RANK.

Problem IN GAUSSIAN likelihood

$$P(x; \mu, \Sigma) = \frac{1}{2\pi |\Sigma|^{1/2}} \exp \left\{ -(x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

| IS NOT DEFINED.

$$\hookrightarrow |\Sigma| = 0$$

WE will fix these issues by examining THREE models

that are Simpler. Spoiler: we'll combine these in the end!

RECALL MLE FOR GAUSSIAN

$$\underset{\mu, \Sigma}{\text{MAX}} \sum_{i=1}^n \log \frac{1}{2\pi |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

Equivalent

$$\underset{\mu, \Sigma}{\text{MIN}} \sum_{i=1}^n (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) + \log |\Sigma|$$

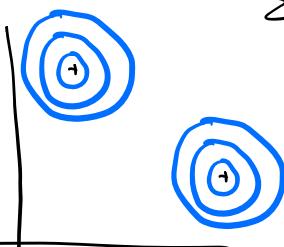
If  $\Sigma$  is full rank,  $\nabla_{\mu} = \sum_{i=1}^n \Sigma^{-1} (\mathbf{x} - \mu) = 0 \Rightarrow \mu = \frac{1}{n} \sum_i \mathbf{x}^{(i)}$

We'll use this as plugin below.

### Building Block 1

Suppose INDEPENDENT AND IDENTICAL COVARIANCE

$$\Sigma = \sigma^2 I \quad (\text{NB: PARAMETER } \sigma^2)$$



COVARIANCE "ARE CIRCLES"

WHAT IS MLE FOR  $\Sigma$ ?

$$|\Sigma| = 2\sigma$$

$$\underset{\sigma^2}{\text{MIN}} \sigma^{-2} \underbrace{\sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)^T (\mathbf{x}^{(i)} - \mu)}_C + d \log \sigma^2$$

$$\text{let } z = \sigma^2 \quad \underset{z}{\text{MIN}} \frac{1}{z} C + d \log z$$

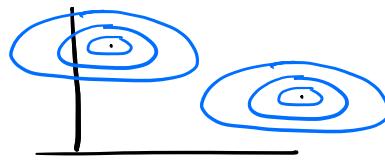
$$\Rightarrow \frac{1}{z} = -\frac{1}{z^2} C + \frac{nd}{z} = 0 \Rightarrow z = \frac{C}{nd}$$

$$\therefore \sigma^2 = \frac{1}{nd} \sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)^T (\mathbf{x}^{(i)} - \mu)$$

"SUBTRACT MEAN AND SQUARE ALL ENTRIES."

## Building Block 2

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \end{bmatrix}$$



Axis Aligned ellipse

SET  $z_i = \sigma_i^2$  (SAME IDEA AS ABOVE)

$$\min_{z_1 \dots z_d} \sum_{i=1}^n \sum_{j=1}^d z_j^{-1} (x^{(i)} - \mu_j)^2 + \log z_j$$

this is  $d$  problems for each 1 dimension

$$\Rightarrow \sum_{i=1}^n z_j^{-1} (x^{(i)} - \mu_j)^2 + \log z_j$$

$$\Rightarrow \sigma_j^2 = \frac{1}{n} \sum_i (x_j^{(i)} - \mu_j)^2$$

## Our FACTOR model

### PARAMETERS

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

$$\Lambda \in \mathbb{R}^{d \times s}$$

$$\Phi \in \mathbb{R}^{d \times d} \text{ - DIAGONAL MATRIX}$$

### MODEL

$$P(x, z) = P(x|z) P(z) \quad z \text{ IS LATENT}$$

$$z \sim N(\mathbf{0}, \mathbf{I}) \in \mathbb{R}^s \text{ for } s < d \text{ "small dim"}$$

$$x = \underbrace{\mu}_{\substack{\text{MEAN} \\ \text{IN} \\ \text{the space}}} + \underbrace{\Lambda z}_{\substack{\text{MAPS FROM SMALL LATENT SPACE TO LARGE SPACE}}} + \epsilon \quad \text{or} \quad x \sim N(\mu + \Lambda z, \Phi)$$

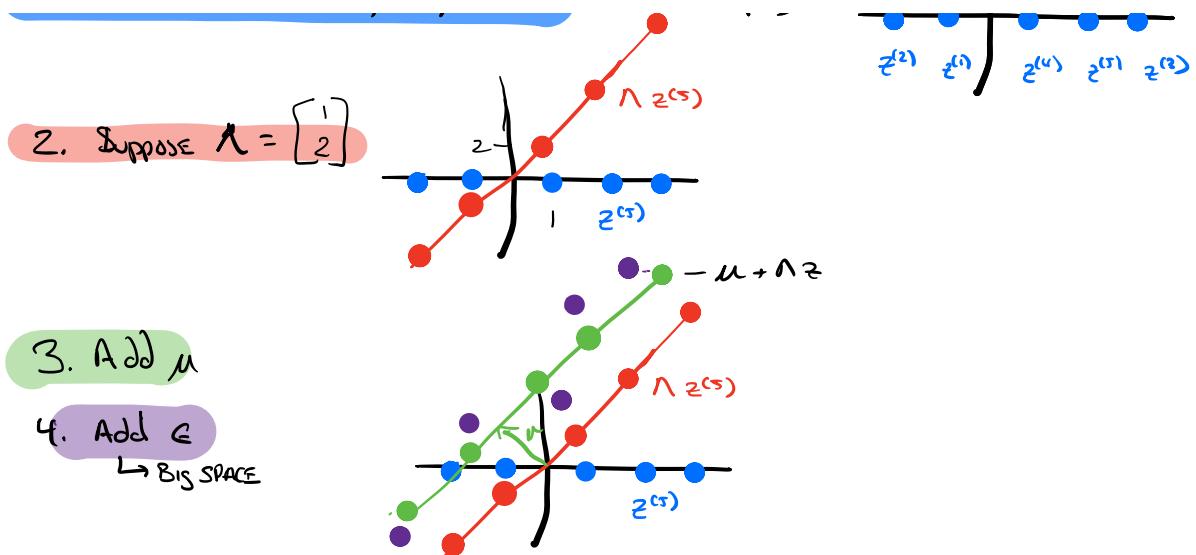
$\epsilon \sim N(\mathbf{0}, \Phi)$  Noisy

$$\text{Ex: } d=2, s=1, n=5$$

$$x = \underbrace{\mu}_{\text{mean}} + \underbrace{\Lambda z}_{\text{latent variable}} + \epsilon$$

1. GENERATE  $z^{(1)}, \dots, z^{(s)}$  from  $N(\mathbf{0}, \mathbf{I})$





DATA WE WOULD OBSERVE ARE Purple DOTS

SO SMALL LATENT SPACE PRODUCES DATA IN HIGH DIM SPACE.

TECHNICAL TOOLS: Block Gaussians

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad x_1 \in \mathbb{R}^{d_1}, x_2 \in \mathbb{R}^{d_2}$$

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

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}^{1 \times 2} \quad \Sigma_{ij} \in \mathbb{R}^{d_i \times d_j} \quad i, j \in \{1, 2\}$$

NOTATION IS widely used AND helpful.

FACT 1:  $P(x_1) = \int_{x_2} P(x_1, x_2)$  MARGINALIZATION

FOR GAUSSIANS,  $P(x_1) = N(\mu_{11}, \Sigma_{11})$  (Not surprising)

FACT 2:  $P(x_1 | x_2) \sim N(\mu_{1|2}, \Sigma_{1|2})$  CONDITIONING

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)$$

$$\hat{\Sigma}_{12} = \hat{\Sigma}_{11} - \hat{\Sigma}_{12} \hat{\Sigma}_{22}^{-1} \hat{\Sigma}_{21} \quad (\text{matrix inversion lemma})$$

Proofs outline (appy to Add)

Summary: MARGINALIZATION  $\neq$  CONDITIONING GAUSSIAN  $\Rightarrow$   
ANOTHER GAUSSIAN (CLOSED)  
WE HAVE formula for PARAMETERS.

Back to Factor Analysis

$$x = \mu + \Lambda z + \epsilon$$

$$\begin{pmatrix} z \\ x \end{pmatrix} \sim N\left(\begin{pmatrix} 0 \\ \mu \end{pmatrix}, \Sigma\right) \quad \text{SINCE } \mathbb{E}[z] = 0$$

$$\mathbb{E}[x] = \mu$$

WHAT IS  $\Sigma$ ?

$$\hat{\Sigma}_{11} = \mathbb{E}[zz^T] = I$$

$$\begin{aligned} \hat{\Sigma}_{12} &= \mathbb{E}[z(x-\mu)^T] = \mathbb{E}[zz^T\Lambda^T] + \cancel{\mathbb{E}[z\epsilon^T]} \\ &= \Lambda^T \end{aligned}$$

$$\hat{\Sigma}_{21} = \hat{\Sigma}_{12}^T$$

$$\begin{aligned} \hat{\Sigma}_{22} &= \mathbb{E}[(x-\mu)(x-\mu)^T] \\ &= \mathbb{E}[(\Lambda z + \epsilon)(\Lambda z + \epsilon)^T] \\ &= \mathbb{E}[\Lambda z z^T \Lambda^T] + \mathbb{E}[\epsilon \epsilon^T] \\ &= \Lambda \Lambda^T + \Phi \end{aligned}$$

$$\hat{\Sigma} = \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Phi \end{bmatrix}$$

E-STEP :  $Q_i(z) = P(z^{(i)} | x^{(i)}; \theta)$  - USE CONDITIONAL!

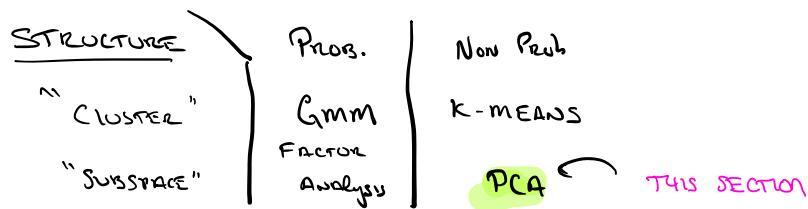
M-STEP : WE HAVE CLOSED FORMS!

### Summary:

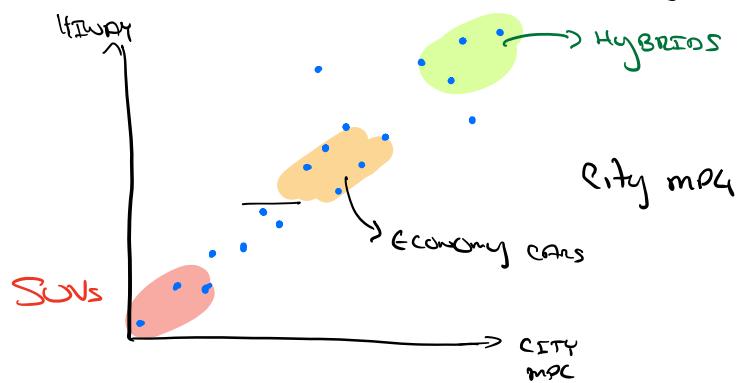
- WE SAW THAT EM CAPTURES GMM
- WE LEARNED ABOUT FACTOR Analysis (Latent low dim. STRUCTURE)
- WE SAW HOW TO ESTIMATE PARAMETERS OF FA USING EM.



# PCA: Principal Component Analysis



Ex: GIVEN PAIRS (Highway mpg, City mpg) of some cars

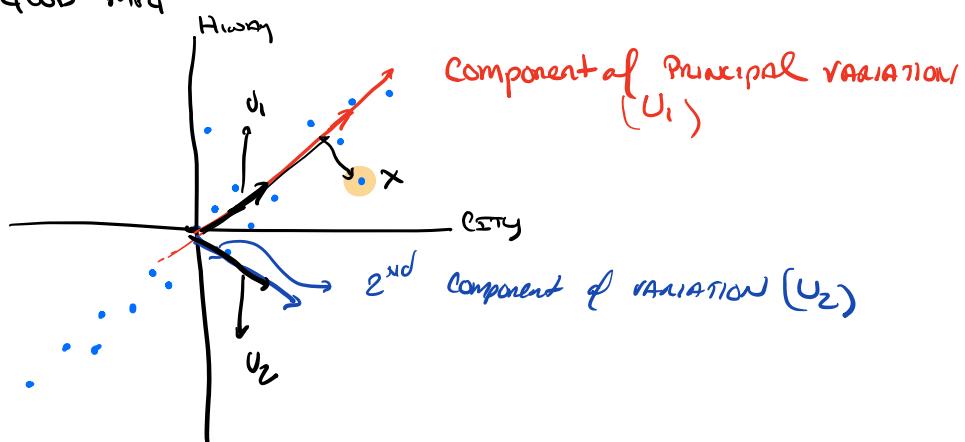


Question: "Good mpg"

① CENTER DATA

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

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



Now  $\|U_1\| = \|U_2\| = 1$  by convention.

- $U_1$  is "How good is mpg"
- $U_2$  is "difference between highway & city" (Roughly)

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

WE may just keep this component

“Explains more variation”

Today: How we find these directions, and some caveats

- think about 1000s of dims  $\rightarrow$  10s of dims
- A dimensionality reduction method

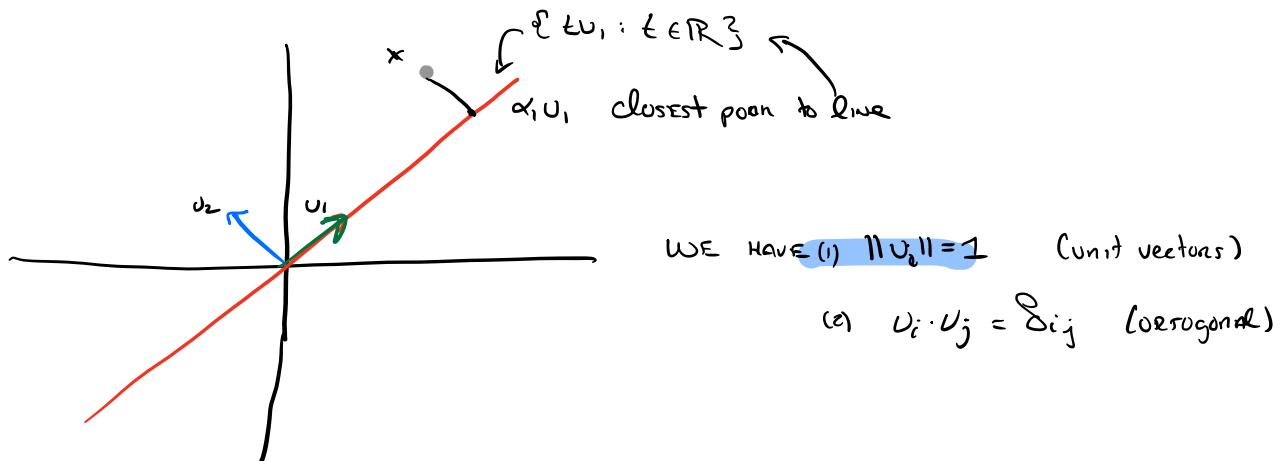
## Preprocessing

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”  
?  $m_p g$

WE will assume data is preprocessed

## PCA AS OPTIMIZATION



How do you find closest point to the line?

$$\alpha_i = \underset{\alpha}{\operatorname{argmin}} \|x - \alpha u_i\|^2$$

$$= \underset{\alpha}{\operatorname{argmin}} \|x\|^2 + \alpha^2 \|u_i\|^2 - 2\alpha (u_i \cdot x)$$

Differentiate w.r.t  $\alpha$

$$2(\alpha - u_i \cdot x) = 0 \Rightarrow \alpha = u_i \cdot x$$

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

$$\underset{\alpha \in \mathbb{R}^d}{\operatorname{Argmin}} \|x - \sum_{i=1}^k \alpha_i U_i\|^2 = \underset{\alpha}{\operatorname{argmin}} \|x\|^2 + \sum_{i=1}^k \alpha_i^2 \|U_i\|^2 - 2\alpha_i \langle U_i, x \rangle$$

Hence  $\alpha_i = U_i \cdot x$

WE CALL  $\|x - \sum_{i=1}^k \alpha_i U_i\|^2$  THE Residual

WE CAN find PCA by either

in class ① MAXIMIZE Projected Subspace

② MINIMIZE Residual

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

WE NEED some facts  
to solve this

LET A BE symmetric & square, then

$$A = U \Lambda U^T$$
 in which

- $U U^T = I$  (orthonormal)
- $\Lambda$  is diagonal

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

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

$$\begin{aligned} Ax &= U \Lambda U^T x = U \Lambda \sum_{i=1}^n \alpha_i e_i && \text{STANDARD BASIS VECTOR} \\ &= U \sum_{i=1}^n \lambda_i \alpha_i e_i && \text{diagonal } \Lambda \\ &= \sum_{i=1}^n \lambda_i \alpha_i U_i \end{aligned}$$

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

$$\underset{\mathbf{x}: \|\mathbf{x}\|^2=1}{\text{MAX}} \quad \mathbf{x}^T \mathbf{A} \mathbf{x} = \underset{\alpha: \|\alpha\|^2=1}{\text{MAX}} \sum_{i=1}^n \alpha_i^2 \lambda_i$$

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

Which  $\mathbf{x}$  attains it? If  $\lambda_1 = \lambda_2$ ?

Now, back to PCA!

$$\underset{\mathbf{U}: \|\mathbf{U}\|^2=1}{\text{MAX}} \frac{1}{n} \sum_{i=1}^n (\mathbf{U}_i \cdot \mathbf{x})^2$$

THE PROJECTION onto  $\mathbf{U}$

$$\underset{\mathbf{U} \in \mathbb{R}^{d \times d}}{=} \frac{1}{n} \sum_{i=1}^n \mathbf{U}^T \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \mathbf{U} = \mathbf{U}^T \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \right) \mathbf{U}$$

COVARIANCE of DATA  
(WE SUBTRACTED MEAN)

$\therefore \mathbf{U}$  is principal Eigenvector

WHAT IF WE WANT MORE DIMENSIONS? WE KEEP TOP-1

How do we represent DATA?

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

WE KEEP THESE  $k$  SCALARS

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

How do we choose  $k$ ?

ONE APPROXIMATE "Amount of Explained Variance"

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

1-1

NB: Only makes sense if  $\lambda_j \geq 0$ . Hence covariance is important

Lurking Instability: Suppose  $\lambda_k = \lambda_{k+1} \dots$  what happens?

REP IS UNSTABLE HERE

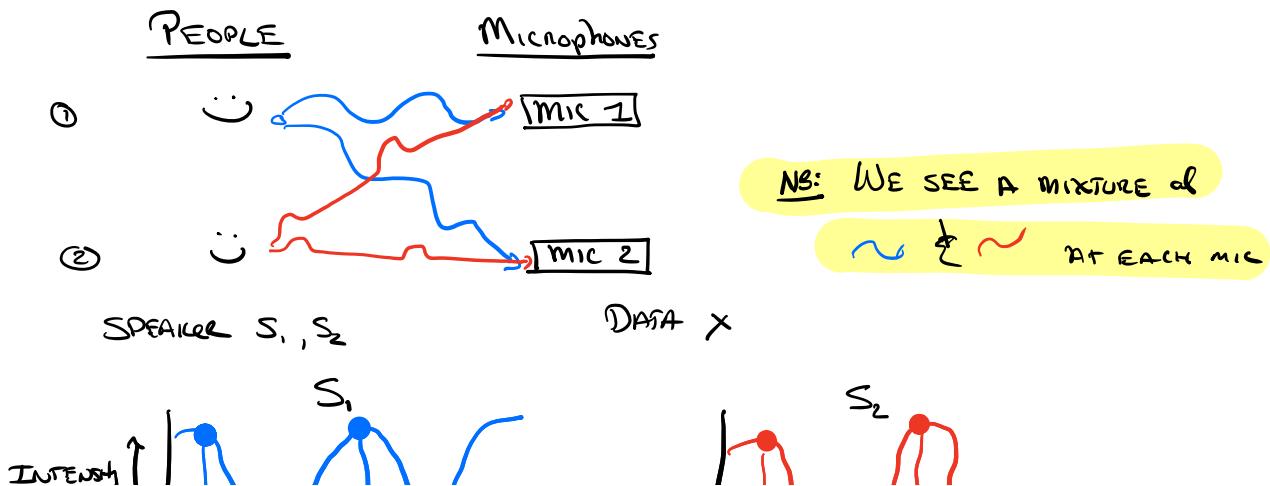
## RECAP of PCA

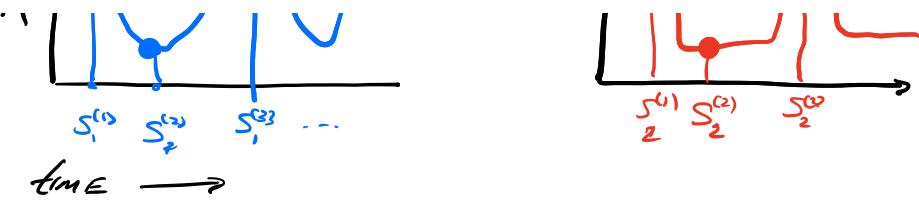
- Dimensionality Reduction technique (e.g. Visualization)
    - MAIN IDEA IS TO PROJECT ON A SUBSPACE, NICE theory.

# ICA INDEPENDENT Component Analysis

- high-level story
  - Key facts  $\not\equiv$  likelihood
  - model

## Cocktail Party Problem (in hw!)



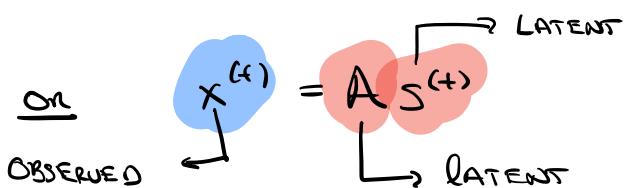


$S_j^{(t)}$  IS INTENSITY AT TIME  $t$  FROM SPEAKER  $j$

WE DO NOT OBSERVE  $S^{(t)}$  ONLY  $x^{(t)}$  - THE MICROPHONES

ex model  $x_j^{(t)} = \alpha_{j1} S_1^{(t)} + \alpha_{j2} S_2^{(t)}$

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



For simplicity, assume # of SPEAKERS = # of mics =  $d$

GIVEN:  $x^{(1)}, x^{(2)}, \dots, x^{(n)} \in \mathbb{R}^d$  d is # of microphones & speakers

DO: find  $S^{(1)}, \dots, S^{(n)} \in \mathbb{R}^d$   
AND  $A \in \mathbb{R}^{d \times d}$  st.  $x^{(t)} = AS^{(t)}$

WE CALL  $A$  THE **MIXING MATRIX** AND  $W = A^{-1}$  **UNMIXING MATRIX**

WRITE  $W = \begin{bmatrix} W_1 \\ \vdots \\ W_d \end{bmatrix}$  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  $\underline{\text{ID}}$  (cold swap 1 to 2)
  - CAN'T DETERMINE ABSOLUTE INTENSITY
$$(cA)(c^{-1}s^{(c)}) = As^{(+)}$$
 for any  $c \neq 0$

Surprisingly Speakers cannot be Gaussian  
 Suppose so  $x^{(t)} \sim N(\mu_t, A A^T)$  show if  $U^T U = I$  AU generates the SAME data.

Nevertheless, we can recover something meaningful!

Algorithm: Just MLE, solved by grad descent

Detour: Density under linear transform (Key Confusion)

Ex:  $S \sim \text{Uniform}[0, 1]$   $U = 2S$  what is PDF of  $U$ ?

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



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

THE key ISSUE is the NORMALIZATION constant

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

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

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

CHANGE of VAR  
formula for  
INTEGRALS

From here ICA is MLE:

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

"sources are independent,

AND HAVE SAME distribution"

$$P(x) = \prod_{j=1}^d P_s(w \cdot x) \cdot |\det(w)|$$

(use linear transform rule)

Now written in terms of  $w$  and  $A$ .

Key technical bit: Use non-rotational invariant distribution

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

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

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

RECAP:

- SAW PCA. Workhorse Dimensionality Reduction
- ICA. Key ideas for now. Introduce "up to symmetry".