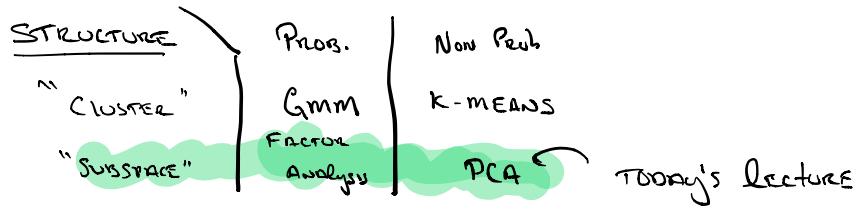


Factor Analysis \neq PCA Subspace structure.



Factor Analysis Given $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$
 "n < d" \rightarrow Data points much smaller than dims

Idea: Posit some structure, use to reduce dimensions
 (Sparse Subspace)

Technical Challenge Estimating Gaussian means when needed

Challenge $\Rightarrow P(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d}} \exp \left\{ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right\}$

\downarrow det = 0! \downarrow undefined

Look at cases (Assumptions) so that Σ is full rank but lower params.

RECALL MLE for Gaussian is equivalent to

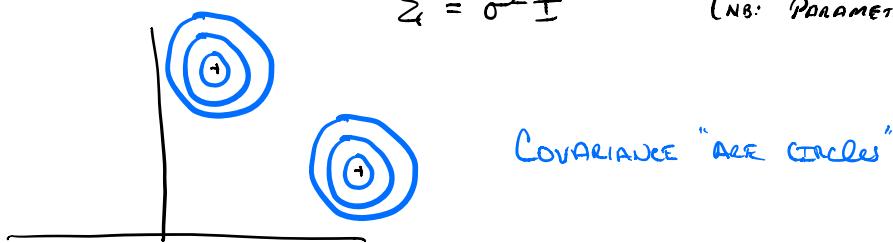
$$\min_{\mu, \Sigma} \sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)^T \Sigma^{-1} (\mathbf{x}^{(i)} - \mu) + \log |\Sigma|$$

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

Building Block 1

Suppose INDEPENDENT AND IDENTICAL COVARIANCE

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



WHAT IS MLE FOR Σ ?

$$|\Sigma| = 2\lambda$$

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

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

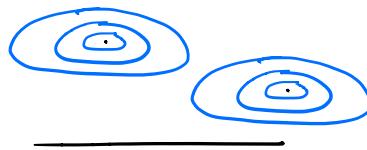
$$\Rightarrow \frac{1}{z} = -z^{-2} C + \frac{d}{z} = 0 \Rightarrow z = \frac{C}{nd}$$

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

"SUBTRACT MEAN AND SQUARE ALL ENTRIES."

Building Block 2

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



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

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

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

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

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

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

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

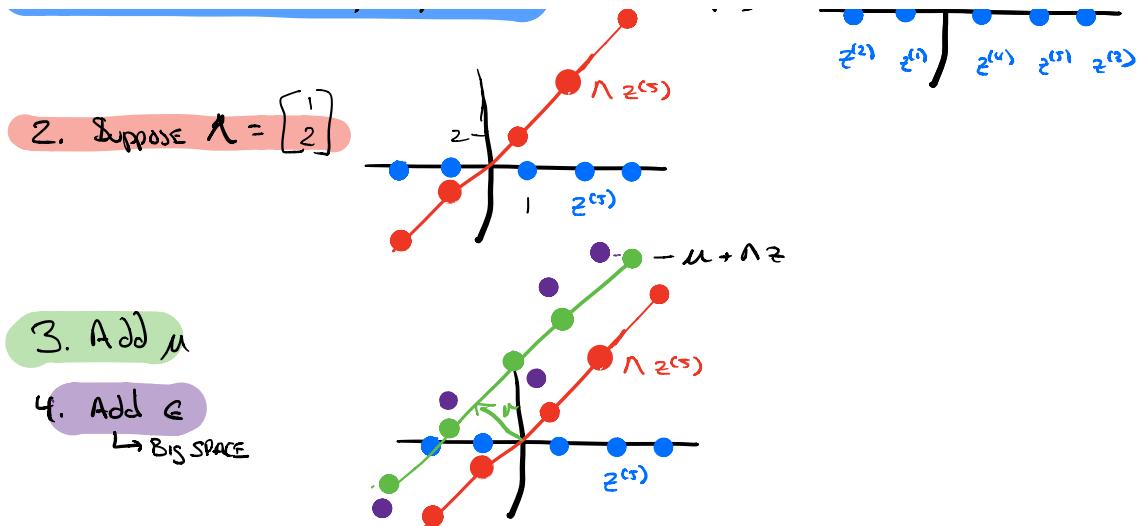
$$\epsilon \sim N(0, \Phi) \quad \text{Noisy}$$

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

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

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





DATA WE would OBSERVE ARE Purple Dots

So small latent space produces data in high dim space.

TECHNICAL TOOLS : Block Gaussians

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

$$\hat{\Sigma} = \begin{bmatrix} \frac{d_1}{\sum_{11}} & \frac{d_2}{\sum_{21}} \\ \frac{\sum_{11}}{d_1} & \frac{\sum_{21}}{d_2} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} \quad \hat{\Sigma}_{ij} \in \mathbb{R}^{d_i \times d_j} \quad i, j \in \{1, 2\}$$

Notation is widely used and helpful.

$$\text{FACT 1: } P(x_1) = \sum_{x_2} P(x_1, x_2) \quad \text{MARGINALIZATION}$$

For Gaussians, $p(x_i) = N(\mu_{ii}, \Sigma_{ii})$ (not zeroising)

$$\underline{\text{Fact 2:}} \quad P(x_1 | x_2) \sim N(\mu_{1|2}, \Sigma_{1|2}) \quad \text{conditioning}$$

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

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

Proofs outline (apply to ABD)

Summary: Marginalization $\not\equiv$ 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 Σ ?

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

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

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

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

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

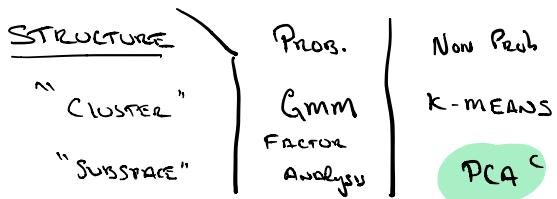
E-STEP : $Q_i(z) = \tilde{P}(z^{(i)} | x^{(i)}; \theta)$ - USE CONDITIONAL!

M-STEP : WE HAVE CLOSED FORMS!

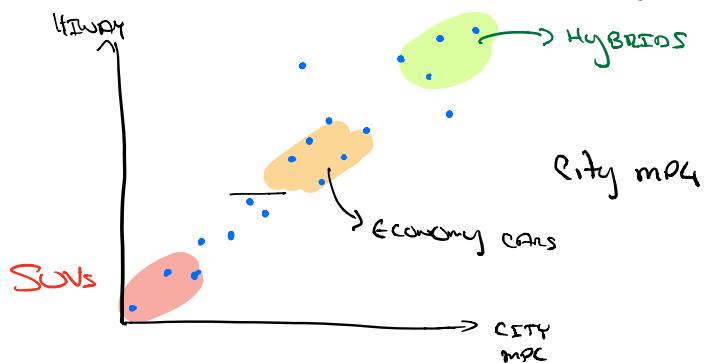
Summary of Factor Analysis

- 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 (hiway mpg, city mpg) of some cars

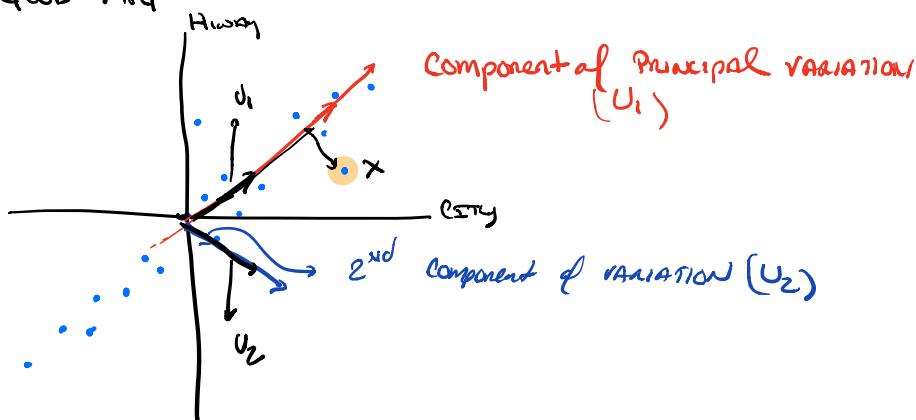


Question: "Good mpg"

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 hiway & 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

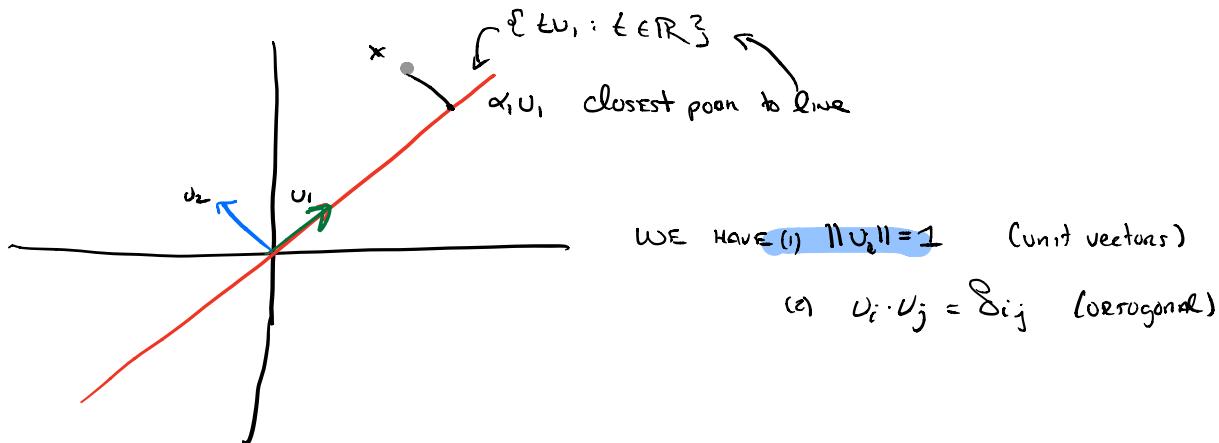
Preprocessing

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 PG

WE will assume data is preprocessed

PCA AS OPTIMIZATION



How do you find closest point to the line?

$$\begin{aligned}\alpha_i &= \underset{\alpha}{\operatorname{argmin}} \|x - \alpha u_i\|^2 \\ &= \underset{\alpha}{\operatorname{argmin}} \|x\|^2 + \alpha^2 \|u_i\|^2 - 2\alpha(u_i \cdot x)\end{aligned}$$

Differentiate w.r.t α

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

Generalize: $u_1 \dots u_k \in \mathbb{R}^d$ AND $x \in \mathbb{R}^d$ use $u_i \cdot u_j = \delta_{ij}$

$$\underset{\alpha_1, \dots, \alpha_d}{\operatorname{Argmin}} \|x - \sum_{i=1}^k \alpha_i u_i\|^2 = \underset{\alpha}{\operatorname{argmin}} \|x\|^2 + \sum_{i=1}^k \alpha_i^2 \|u_i\|^2 - 2\alpha_i \langle u_i, x \rangle$$

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

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

LET A be symmetric & square, then

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

- $U U^T = I$ (orthonormal)
- Λ 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 $[u_1 \dots u_n] = U$

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

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!

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

THE projection onto \mathbf{U}

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

Covariance of data
(WE SUBTRACTED MEAN)

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

WHAT IF WE WANT MORE DIMENSIONS? WE KEEP \mathbf{U}_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"

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

$j=1$

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

Ranking 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.
- cf Factor Analysis