

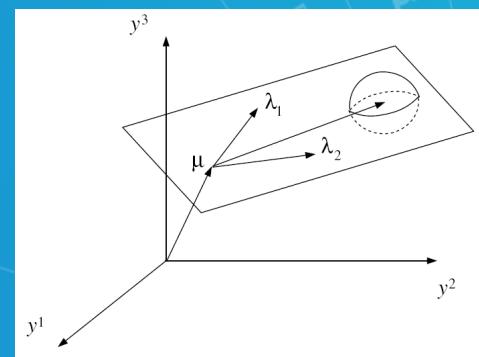
10701: Introduction to Machine Learning

Factor Analysis and Independent Component Analysis

Eric Xing

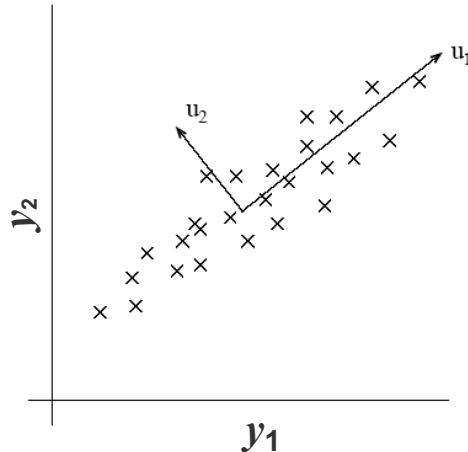
Lecture 16, October 28, 2020

Reading: 12.3-4, C.B book



Recap of PCA

- Popular dimensionality reduction technique
- Project data onto directions of greatest variation



$$\begin{aligned} u &= \arg \max \frac{1}{m} \sum_{i=1}^m (\vec{y}_i^T u)^2 \\ &= \arg \max u^T \left(\frac{1}{m} \sum_{i=1}^m \vec{y}_i \vec{y}_i^T \right) u \\ &= \arg \max (u^T \text{Cov}(y) u) \end{aligned}$$

$$\begin{aligned} \vec{x}_i &= \begin{bmatrix} u_1^T \vec{y}_i \\ u_2^T \vec{y}_i \\ \vdots \\ u_q^T \vec{y}_i \end{bmatrix} = U_q^T \vec{y}_i \in \mathbb{R}^q \\ \vec{y}_i &= U \vec{x}_i \end{aligned}$$

- Consequence:
 - x_i are uncorrelated such that the covariance matrix
 - Truncation error

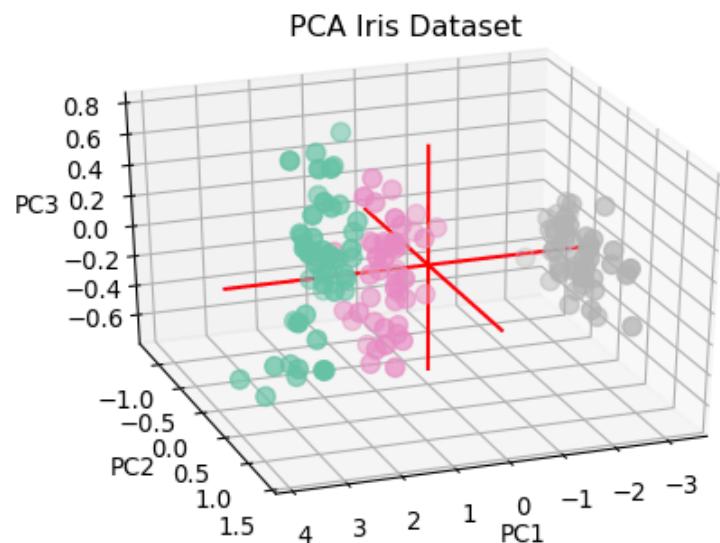
$$\Sigma_y = \sum_{k=1}^K \gamma_k (u_k u_k^T) \approx \sum_{k=1}^q \gamma_k (u_k u_k^T) = \Sigma_x$$

$$\frac{1}{m} \sum_{i=1}^m \vec{x}_i \vec{x}_i^T \quad \text{is} \quad \begin{bmatrix} \gamma_1 & & \\ & \ddots & \\ & & \gamma_q \end{bmatrix}$$



Recap of PCA

- ❑ Popular dimensionality reduction technique
- ❑ Project data onto directions of greatest variation



Useful tool for visualising patterns and clusters within the data set, but ...

Need centering

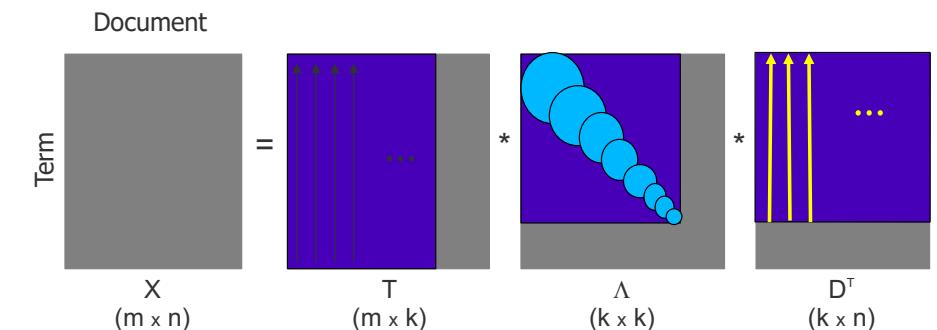
Does not explicitly model data noise



PCA , SVD, and Eigenvalue Decomposition

- The first root is called the principal eigenvalue which has an associated orthonormal ($u^T u = 1$) *eigenvector* u
- Subsequent roots are ordered such that $\lambda_1 > \lambda_2 > \dots > \lambda_M$ with $\text{rank}(D)$ non-zero values.
- Eigenvectors form an orthonormal basis i.e. $u_i^T u_j = \delta_{ij}$
- The eigenvalue decomposition of $XX^T = U\Sigma U^T$,
where $U = [u_1, u_2, \dots, u_M]$ and $\Sigma = \text{diag}[\lambda_1, \lambda_2, \dots, \lambda_M]$
- Similarly the eigenvalue decomposition of $X^T X = V\Sigma V^T$

- The SVD is closely related to the above $X=U \Sigma^{1/2} V^T$
- The left eigenvectors U , right eigenvectors V ,
- singular values = square root of eigenvalues.



Probabilistic PCA

- PCA can be cast as a probabilistic model

$$y_n = \Lambda x_n + \mu + \varepsilon_n \quad \varepsilon_n \sim \mathbf{N}(0, \sigma^2 I)$$

with q -dimensional latent variables $x_n \sim \mathbf{N}(0, I)$

- The resulting data distribution is

$$y_n \sim \mathbf{N}(\mu, \Lambda \Lambda^T + \sigma^2 I)$$

- Maximum likelihood solution is equivalent to PCA

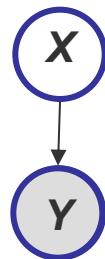
$$\mu^{ML} = \frac{1}{N} \sum_n y_n \quad \Lambda^{ML} = U_q (\Gamma_q - \sigma^2 I)^{1/2}$$

Diagonal Γ_q contains the top q sample covariance eigen-values and U_q contains associated eigenvectors



Factor analysis

- An unsupervised linear regression model

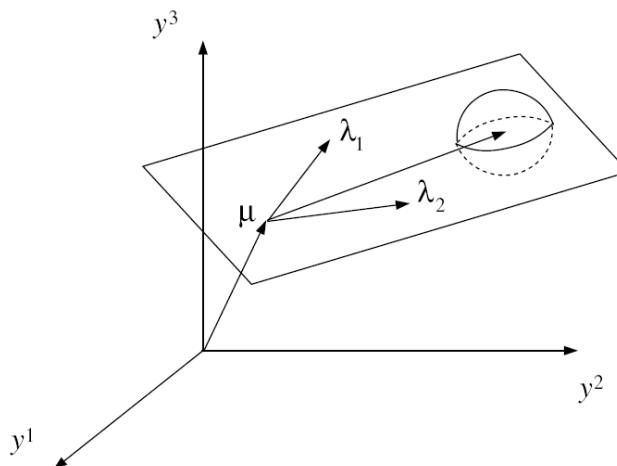


$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, I)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{x}, \boldsymbol{\Psi})$$

where $\boldsymbol{\Lambda}$ is called a **factor loading matrix**, and $\boldsymbol{\Psi}$ is diagonal.

- Geometric interpretation



- To generate data, first generate a point within the manifold then add noise.
Coordinates of point are components of latent variable.



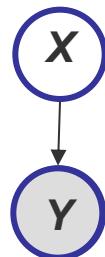
Relationship between PCA and FA

- Probabilistic PCA is equivalent to factor analysis with equal noise for every dimension, i.e., $\varepsilon_n \sim \text{isotropic Gaussian } \mathcal{N}(0, \sigma^2 I)$
- In factor analysis $\varepsilon_n \sim \mathcal{N}(0, \Psi)$ for a diagonal covariance matrix Ψ
- An iterative algorithm (eg. EM) is required to find parameters if precisions are not known in advance



Factor analysis

- An unsupervised linear regression model

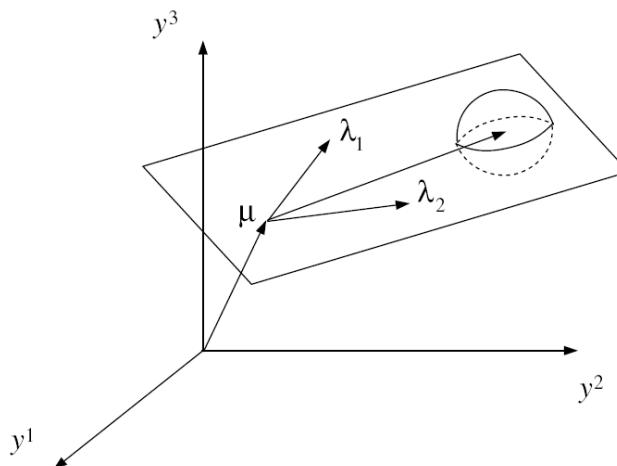


$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, I)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{x}, \boldsymbol{\Psi})$$

where $\boldsymbol{\Lambda}$ is called a factor loading matrix, and $\boldsymbol{\Psi}$ is diagonal.

- Geometric interpretation

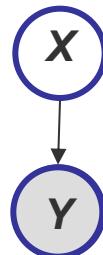


- To generate data, first generate a point within the manifold then add noise.
Coordinates of point are components of latent variable.



Marginal data distribution

- A marginal Gaussian (e.g., $p(\mathbf{x})$) times a conditional Gaussian (e.g., $p(\mathbf{y}|\mathbf{x})$) is a joint Gaussian
- Any marginal (e.g., $p(\mathbf{y})$) of a joint Gaussian (e.g., $p(\mathbf{x}, \mathbf{y})$) is also a Gaussian
- Since the marginal is Gaussian, we can determine it by just computing its mean and variance. (Assume noise uncorrelated with data.)



$$\begin{aligned}E[\mathbf{Y}] &= E[\mu + \Lambda \mathbf{X} + \mathbf{W}] \quad \text{where } \mathbf{W} \sim \mathcal{N}(\mathbf{0}, \Psi) \\&= \mu + \Lambda E[\mathbf{X}] + E[\mathbf{W}] \\&= \mu + \mathbf{0} + \mathbf{0} = \mu \\Var[\mathbf{Y}] &= E[(\mathbf{Y} - \mu)(\mathbf{Y} - \mu)^T] \\&= E[(\mu + \Lambda \mathbf{X} + \mathbf{W} - \mu)(\mu + \Lambda \mathbf{X} + \mathbf{W} - \mu)^T] \\&= E[(\Lambda \mathbf{X} + \mathbf{W})(\Lambda \mathbf{X} + \mathbf{W})^T] \\&= \Lambda E[\mathbf{X}\mathbf{X}^T] \Lambda^T + E[\mathbf{W}\mathbf{W}^T] \\&= \Lambda \Lambda^T + \Psi\end{aligned}$$



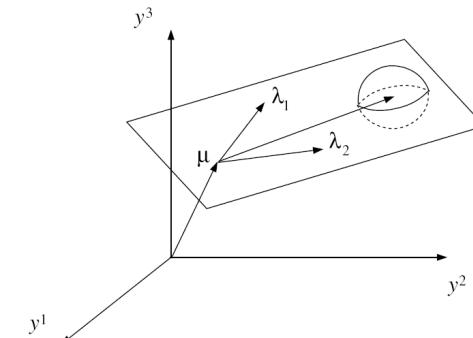
FA = Constrained-Covariance Gaussian

- Marginal density for factor analysis (\mathbf{y} is p -dim, \mathbf{x} is k -dim):

$$p(\mathbf{y} | \theta) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \Lambda \Lambda^T + \boldsymbol{\Psi})$$

- So the effective covariance is the low-rank outer product of two long skinny matrices plus a diagonal matrix:

$$\text{Cov}[\mathbf{y}] = \begin{matrix} \boxed{\Lambda} \end{matrix} \begin{matrix} \boxed{\Lambda^T} \\ + \end{matrix} \begin{matrix} \boxed{\Psi} \end{matrix}$$



- In other words, factor analysis is just a constrained Gaussian model. (If were not diagonal then we could model any Gaussian and it would be pointless.)



Review: A primer to multivariate Gaussian

- Multivariate Gaussian density:

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

- A joint Gaussian:

$$p\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} | \mu, \Sigma\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \middle| \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

- How to write down $p(\mathbf{x}_1)$, $p(\mathbf{x}_1|\mathbf{x}_2)$ or $p(\mathbf{x}_2|\mathbf{x}_1)$ using the block elements in μ and Σ ?
- Formulas to remember:

$$p(\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_2 | \mathbf{m}_2^m, \mathbf{V}_2^m)$$

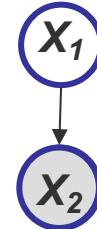
$$\mathbf{m}_2^m = \mu_2$$

$$\mathbf{V}_2^m = \Sigma_{22}$$

$$p(\mathbf{x}_1|\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 | \mathbf{m}_{1|2}, \mathbf{V}_{1|2})$$

$$\mathbf{m}_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{x}_2 - \mu_2)$$

$$\mathbf{V}_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$



Review: Some matrix algebra

- Trace and derivatives
 - Cyclical permutations

$$\text{tr}[A] \stackrel{\text{def}}{=} \sum_i a_{ii}$$

- Derivatives

$$\text{tr}[ABC] = \text{tr}[CAB] = \text{tr}[BCA]$$

$$\frac{\partial}{\partial A} \text{tr}[BA] = B^T$$

$$\frac{\partial}{\partial A} \text{tr}[x^T Ax] = \frac{\partial}{\partial A} \text{tr}[xx^T A] = xx^T$$

- Determinants and derivatives

$$\frac{\partial}{\partial A} \log|A| = A^{-T}$$



FA joint distribution

- Model

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, I)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu} + \Lambda \mathbf{x}, \Psi)$$

- Covariance between \mathbf{x} and \mathbf{y}

$$\begin{aligned} \text{Cov}[\mathbf{X}, \mathbf{Y}] &= E[(\mathbf{X} - \mathbf{0})(\mathbf{Y} - \boldsymbol{\mu})^T] = E[\mathbf{X}(\boldsymbol{\mu} + \Lambda \mathbf{X} + \mathbf{W} - \boldsymbol{\mu})^T] \\ &= E[\mathbf{X}\mathbf{X}^T \Lambda^T + \mathbf{X}\mathbf{W}^T] \\ &= \Lambda^T \end{aligned}$$

- Hence the joint distribution of \mathbf{x} and \mathbf{y} :

$$p\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \mid \begin{bmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda\Lambda^T + \Psi \end{bmatrix}\right)$$

- Assume noise is uncorrelated with data or latent variables.



Inference in Factor Analysis

- Apply the Gaussian conditioning formulas to the joint distribution we derived above, where

$$\Sigma_{11} = I$$

$$\Sigma_{12} = \Sigma_{12}^T = \Lambda^T$$

$$\Sigma_{22} = (\Lambda\Lambda^T + \Psi)$$

we can now derive the posterior of the latent variable \mathbf{x} given observation \mathbf{y} , $p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x} | \mathbf{m}_{1|2}, \mathbf{V}_{1|2})$, where

$$\begin{aligned}\mathbf{m}_{1|2} &= \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{y} - \mu_2) & \mathbf{V}_{1|2} &= \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \\ &= \Lambda^T(\Lambda\Lambda^T + \Psi)^{-1}(\mathbf{y} - \mu) & &= I - \Lambda^T(\Lambda\Lambda^T + \Psi)^{-1}\Lambda\end{aligned}$$

Applying the matrix inversion lemma $(E - FH^{-1}G)^{-1} = E^{-1} + E^{-1}F(H - GE^{-1}F)^{-1}GE^{-1}$

$$\Rightarrow \quad \mathbf{V}_{1|2} = (I + \Lambda^T\Psi^{-1}\Lambda)^{-1} \quad \mathbf{m}_{1|2} = \mathbf{V}_{1|2}\Lambda^T\Psi^{-1}(\mathbf{y} - \mu)$$

- Here we only need to invert a matrix of size $|\mathbf{x}|'|\mathbf{x}|$, instead of $|\mathbf{y}|'|\mathbf{y}|$.



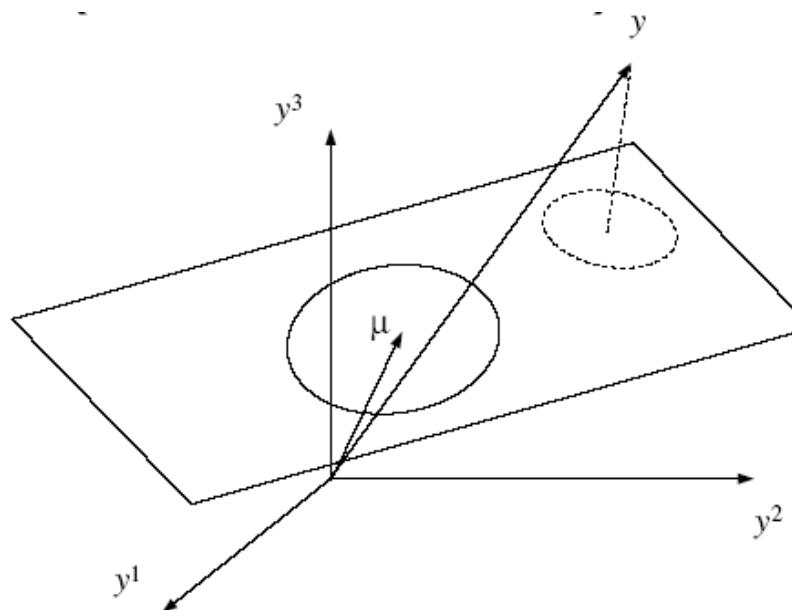
Geometric interpretation: inference is linear projection

- The posterior is:

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}; \mathbf{m}_{1|2}, \mathbf{V}_{1|2})$$

$$\mathbf{V}_{1|2} = (I + \Lambda^T \Psi^{-1} \Lambda)^{-1} \quad \mathbf{m}_{1|2} = \mathbf{V}_{1|2} \Lambda^T \Psi^{-1} (\mathbf{y} - \mu)$$

- Posterior covariance does not depend on observed data \mathbf{y} !
- Computing the posterior mean is just a linear operation:



EM for Factor Analysis

- Incomplete data log likelihood function (marginal density of y)

$$\begin{aligned}\ell(\theta, D) &= -\frac{N}{2} \log |\Lambda \Lambda^T + \Psi| - \frac{1}{2} \sum_n (\mathbf{y}_n - \boldsymbol{\mu})^T (\Lambda \Lambda^T + \Psi)^{-1} (\mathbf{y}_n - \boldsymbol{\mu}) \\ &= -\frac{N}{2} \log |\Lambda \Lambda^T + \Psi| - \frac{1}{2} \text{tr}[(\Lambda \Lambda^T + \Psi)^{-1} \mathbf{S}], \quad \text{where } \mathbf{S} = \sum_n (\mathbf{y}_n - \boldsymbol{\mu})(\mathbf{y}_n - \boldsymbol{\mu})^T\end{aligned}$$

- Estimating $\boldsymbol{\mu}$ is trivial: $\hat{\boldsymbol{\mu}}^{ML} = \frac{1}{N} \sum_n \mathbf{y}_n$
- Parameters Λ and Ψ are coupled nonlinearly in log-likelihood
- Complete log likelihood

$$\begin{aligned}\ell_c(\theta, D) &= \sum_n \log p(\mathbf{x}_n, \mathbf{y}_n) = \sum_n \log p(\mathbf{x}_n) + \log p(\mathbf{y}_n | \mathbf{x}_n) \\ &= -\frac{N}{2} \log |I| - \frac{1}{2} \sum_n \mathbf{x}_n^T \mathbf{x}_n - \frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n (\mathbf{y}_n - \Lambda \mathbf{x}_n)^T \Psi^{-1} (\mathbf{y}_n - \Lambda \mathbf{x}_n) \\ &= -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr}[\mathbf{x}_n \mathbf{x}_n^T] - \frac{N}{2} \text{tr}[\mathbf{S} \Psi^{-1}], \quad \text{where } \mathbf{S} = \frac{1}{N} \sum_n (\mathbf{y}_n - \Lambda \mathbf{x}_n)(\mathbf{y}_n - \Lambda \mathbf{x}_n)^T\end{aligned}$$



E-step for Factor Analysis

- Compute

$$\langle \ell_e(\theta, D) \rangle_{p(x|y)}$$

$$\langle \ell_e(\theta, D) \rangle = -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr} [\langle X_n X_n^T \rangle] - \frac{N}{2} \text{tr} [\langle S \rangle \Psi^{-1}]$$

$$\langle S \rangle = \frac{1}{N} \sum_n (Y_n Y_n^T - Y_n \langle X_n^T \rangle \Lambda^T - \Lambda \langle X_n^T \rangle Y_n^T + \Lambda \langle X_n X_n^T \rangle \Lambda^T)$$

$$\langle X_n \rangle = E[X_n | Y_n]$$

$$\langle X_n X_n^T \rangle = \text{Var}[X_n | Y_n] + E[X_n | Y_n] E[X_n | Y_n]^T$$

- Recall that we have derived:

$$V_{1|2} = (I + \Lambda^T \Psi^{-1} \Lambda)^{-1} \quad m_{1|2} = V_{1|2} \Lambda^T \Psi^{-1} (\mathbf{y} - \mu)$$

$$\Rightarrow \langle X_n \rangle = m_{x_n | y_n} = V_{1|2} \Lambda^T \Psi^{-1} (Y_n - \mu) \quad \text{and} \quad \langle X_n X_n^T \rangle = V_{1|2} + m_{x_n | y_n} m_{x_n | y_n}^T$$



M-step for Factor Analysis

- Take the derivates of the expected complete log likelihood wrt. parameters.
 - Using the trace and determinant derivative rules:

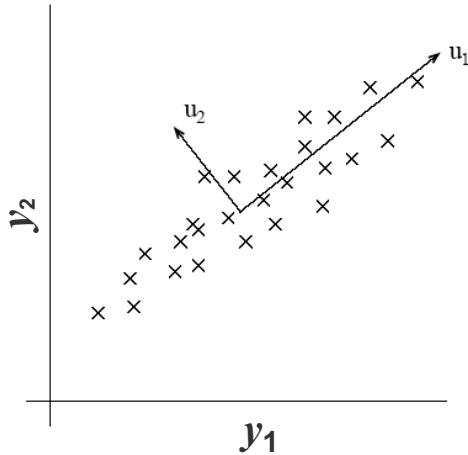
$$\begin{aligned}\frac{\partial}{\partial \Psi^{-1}} \langle \ell_c \rangle &= \frac{\partial}{\partial \Psi^{-1}} \left(-\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr} [\langle X_n X_n^T \rangle] - \frac{N}{2} \text{tr} [\langle S \rangle \Psi^{-1}] \right) \\ &= \frac{N}{2} \Psi - \frac{N}{2} \langle S \rangle \quad \Rightarrow \quad \Psi^{t+1} = \langle S \rangle\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \Lambda} \langle \ell_c \rangle &= \frac{\partial}{\partial \Lambda} \left(-\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr} [\langle X_n X_n^T \rangle] - \frac{N}{2} \text{tr} [\langle S \rangle \Psi^{-1}] \right) = -\frac{N}{2} \Psi^{-1} \frac{\partial}{\partial \Lambda} \langle S \rangle \\ &= -\frac{N}{2} \Psi^{-1} \frac{\partial}{\partial \Lambda} \left(\frac{1}{N} \sum_n (Y_n Y_n^T - Y_n \langle X_n^T \rangle \Lambda^T - \Lambda \langle X_n^T \rangle Y_n^T + \Lambda \langle X_n X_n^T \rangle \Lambda^T) \right) \\ &= \Psi^{-1} \sum_n Y_n \langle X_n^T \rangle - \Psi^{-1} \Lambda \sum_n \langle X_n X_n^T \rangle \quad \Rightarrow \quad \Lambda^{t+1} = \left(\sum_n Y_n \langle X_n^T \rangle \right) \left(\sum_n \langle X_n X_n^T \rangle \right)^{-1}\end{aligned}$$



Comparison of PCA and FA

□ PCA

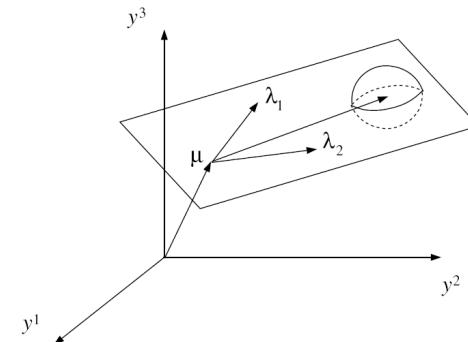


$$y_n = Ux_n$$

$$u = \arg \max(u^T \text{Cov}(y)u)$$

$$\vec{x}_i = \begin{bmatrix} u_1^T \vec{y}_i \\ u_2^T \vec{y}_i \\ \vdots \\ u_k^T \vec{y}_i \end{bmatrix} = U_q^T \vec{y}_i \in \mathbb{R}^q$$

□ FA



$$y_n = \Lambda x_n + \mu + \varepsilon_n$$

$$\varepsilon_n \sim \mathcal{N}(0, \Psi)$$

$$\langle \boldsymbol{X}_n \rangle = \mathbf{m}_{x_n|y_n} = \mathbf{V}_{1|2} \Lambda^T \Psi^{-1} (\boldsymbol{y}_n - \mu)$$

$$\text{and } \langle \boldsymbol{X}_n \boldsymbol{X}_n^T \rangle = \mathbf{V}_{1|2} + \mathbf{m}_{x_n|y_n} \mathbf{m}_{x_n|y_n}^T$$

$$\Lambda^{t+1} = \left(\sum_n \boldsymbol{y}_n \langle \boldsymbol{X}_n^T \rangle \right) \left(\sum_n \langle \boldsymbol{X}_n \boldsymbol{X}_n^T \rangle \right)^{-1}$$



Comparison of PCA and FA

□ PCA

$$u = \arg \max (u^T Cov(y) u)$$

$$\vec{x}_i = \begin{bmatrix} u_1^T \vec{y}_i \\ u_2^T \vec{y}_i \\ \vdots \\ u_k^T \vec{y}_i \end{bmatrix} = U_q^T \vec{y}_i \in \mathbf{R}^q$$

- SVD on a $K'K$ matrix
- Covariant under rotation: Ay
- Principle axis can be found incrementally

□ FA

$$\langle X_n \rangle = \mathbf{m}_{x_n|y_n} = \mathbf{V}_{1|2} \Lambda^T \Psi^{-1} (y_n - \mu)$$

$$\text{and } \langle X_n X_n^T \rangle = \mathbf{V}_{1|2} + \mathbf{m}_{x_n|y_n} \mathbf{m}_{x_n|y_n}^T$$

$$\Lambda^{t+1} = \left(\sum_n y_n \langle X_n^T \rangle \right) \left(\sum_n \langle X_n X_n^T \rangle \right)^{-1}$$

$$\Psi^{t+1} = \langle \mathbf{S} \rangle$$

- Invert a $q'q$ matrix
- Covariant under rescaling: $diag(\alpha)y$
- Neither of the factors found by a two-factor model is necessarily the same as that found by a single factor model, and ...



Example:

Original data matrix \longrightarrow Correlation matrix \longrightarrow Factor matrix

		Variables									
		v_1	v_2	.	.	.	v_{10}				
Obs	O_1										
	O_2										
	.										
	O_3										
	.										
	O_4										
	.										
	O_5										
	.										
	O_n										

Observational data

		Variables									
		v_1	v_2	.	.	.	v_{10}				
Variables	v_1										
	v_2										
	.										
	v_3										
	.										
	v_4										
	.										
	v_5										
	.										
	v_{10}										

Correlation coefficients



- Decisions
- Factoring method
 - # of factors to retain
 - Factor rotation

		Factors		
		F_1	F_2	F_3
Factors	Spd	.87	.07	.14
	Long	.77	.27	.13
	High	.51	.34	.34
	110m	.63	.32	.05
	400m	.74	.06	.38
	Discuss	.22	.79	.06
	Shot	.31	.82	.10
	Javelin	.02	.70	.15
	Pole	.24	.40	.50
	1500m	.02	.12	.89

Factor loadings



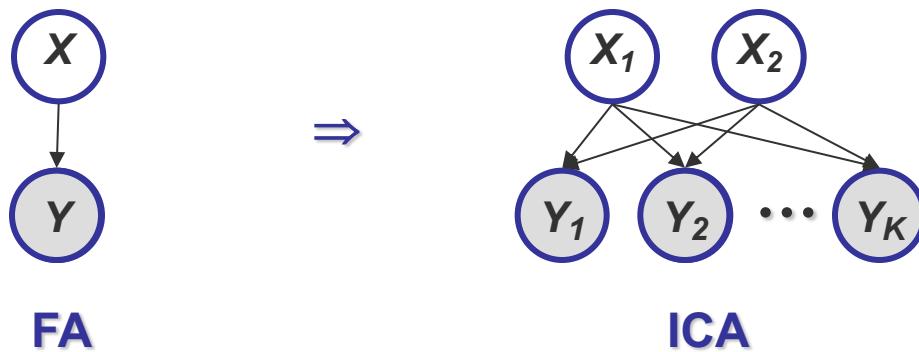
Model Invariance and Identifiability

- There is *degeneracy* in the FA model.
- Since Λ only appears as outer product $\Lambda\Lambda^T$, the model is invariant to rotation and axis flips of the latent space.
- We can replace Λ with ΛQ for any orthonormal matrix Q and the model remains the same: $(\Lambda Q)(\Lambda Q)^T = \Lambda(QQ^T)\Lambda^T = \Lambda\Lambda^T$.
- This means that there is no “one best” setting of the parameters. An infinite number of parameters all give the ML score!
- Such models are called *un-identifiable* since two people both fitting ML parameters to the identical data will not be guaranteed to identify the same parameters.

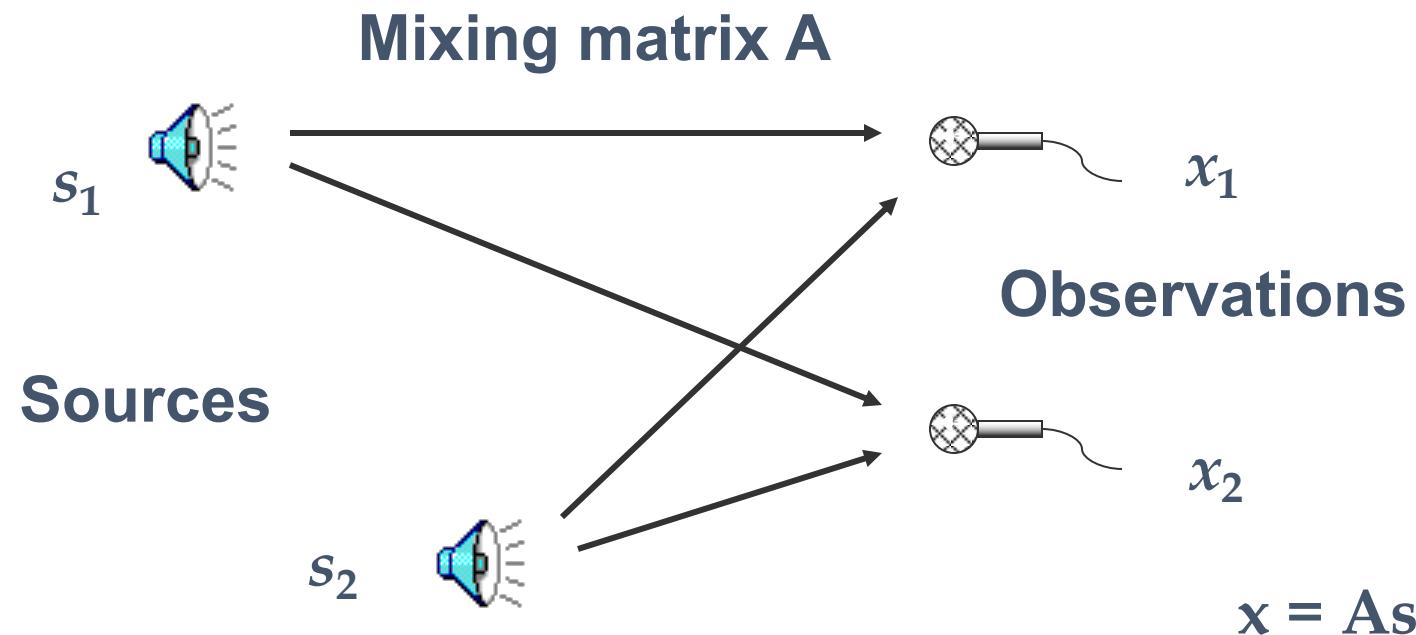


Independent Components Analysis (ICA)

- ❑ ICA is similar to FA, except it assumes the latent source has non-Gaussian density.
- ❑ Hence ICA can extract higher order moments (not just second order).
- ❑ It is commonly used to solve blind source separation (cocktail party problem).



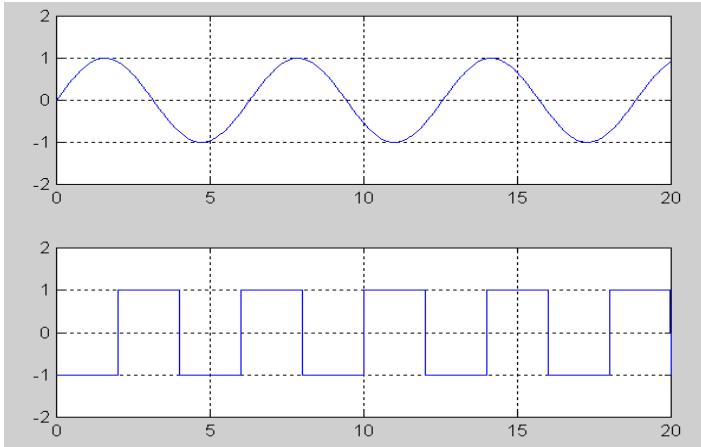
The simple “Cocktail Party” Problem



n sources, $m=n$ observations



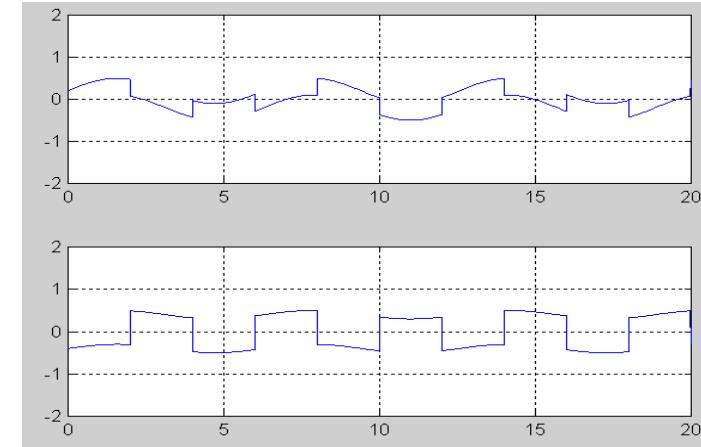
Motivation



Two Independent Sources

$$x_1(t) = a_{11}s_1(t) + a_{12}s_2(t)$$

$$x_2(t) = a_{21}s_1(t) + a_{22}s_2(t)$$

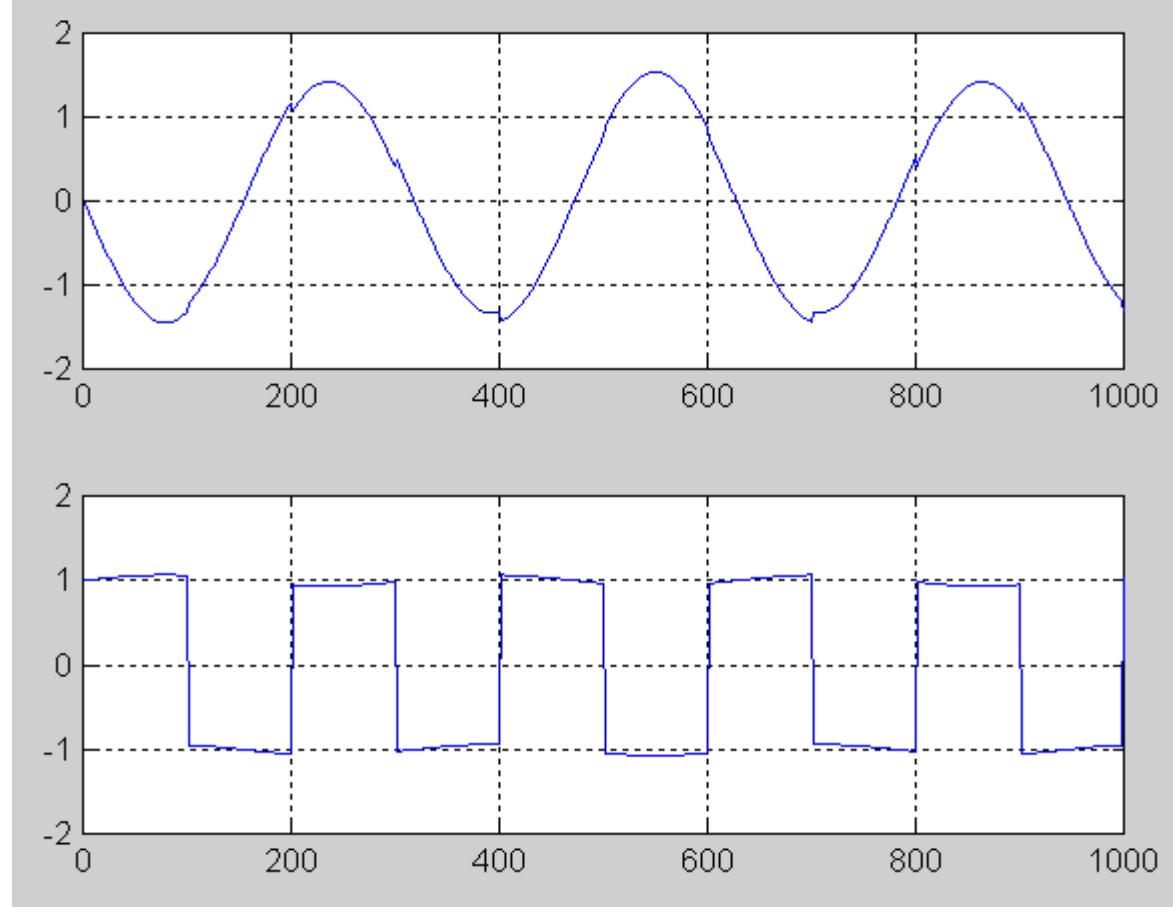


Mixture at two Mics

a_{ij} ... Depend on the distances of the microphones from the speakers



Motivation



Get the Independent Signals out of the Mixture



Blind Source Separation

- Suppose that there are k unknown independent sources

$$\mathbf{s}(t) = [s_1(t), \dots, s_k(t)]^T \quad \text{with} \quad E[\mathbf{s}(t)] = \mathbf{1}$$

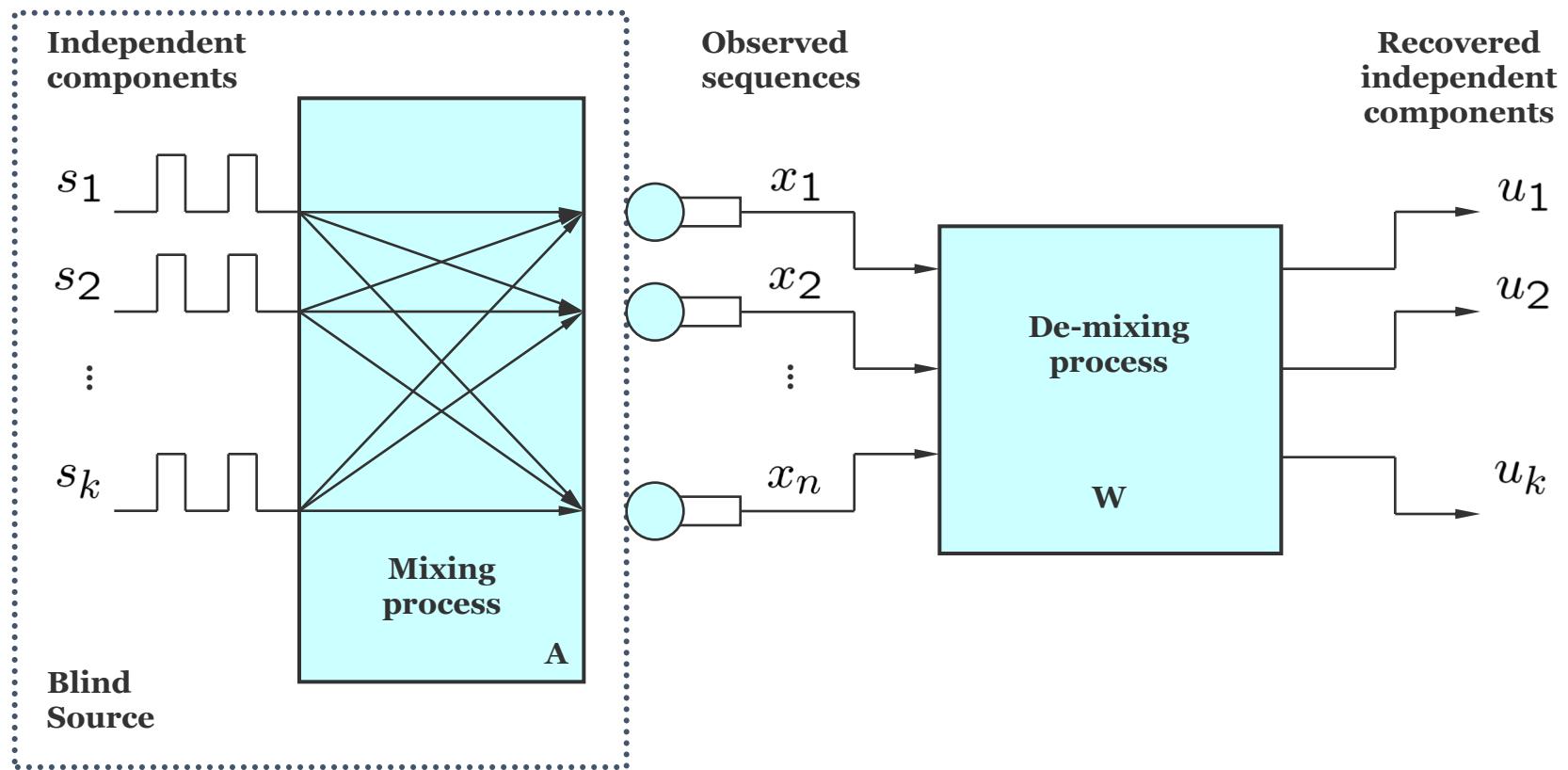
- A data vector $\mathbf{x}(t)$ is observed at each time point t , such that

$$\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t)$$

where \mathbf{A} is a $n \times k$ full rank scalar matrix

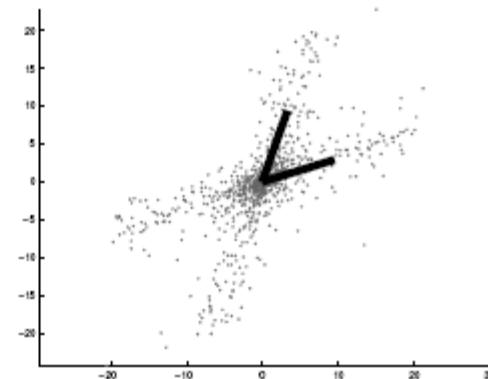
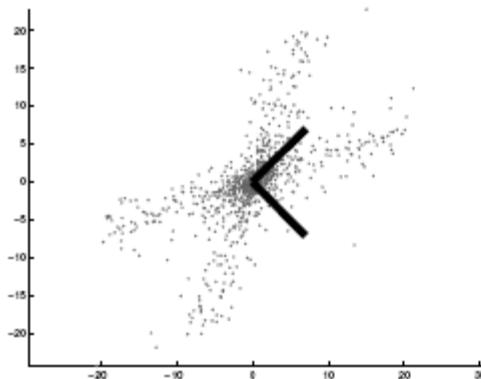


Blind source separation



ICA versus PCA (and FA)

- ❑ Similarity
 - ❑ Feature extraction
 - ❑ Dimension reduction
- ❑ Difference
 - ❑ PCA uses up to second order moment of the data to produce uncorrelated components
 - ❑ ICA strives to generate components as independent as possible



Probability 101: Pearson's correlation

- ❑ Normalized covariance

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}}$$

- ❑ Captures **linear** dependency
 - ❑ Linear regression from X to Y gives $\beta = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$
- ❑ Important properties:
 - ❑ $X \perp\!\!\!\perp Y$ implies $\rho(X, Y) = 0$ (Why?)
 - ❑ $\rho(X, Y) = 0$ does **not** imply $X \perp\!\!\!\perp Y$ (Counterexamples?)
- ❑ Q1: Is there any measure that implies independence?
- ❑ Q2: What kind of dependency should they consider?

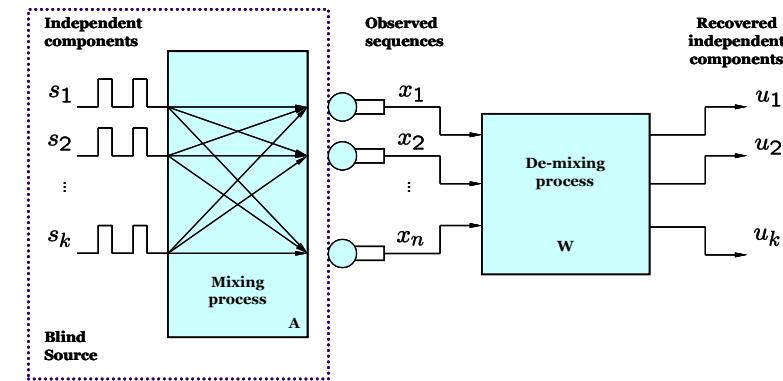


Problem formulation

- The goal of ICA is to find a linear mapping \mathbf{W} such that the unmixed sequences \mathbf{u}

$$\mathbf{u}(t) = \mathbf{Wx}(t) = \mathbf{WA}s(t)$$

are maximally statistically independent



- Find some

$$\mathbf{V} = \mathbf{WA} = \mathbf{PC}$$

where \mathbf{C} is a diagonal matrix and \mathbf{P} is a permutation matrix.



Principle of ICA: Nongaussianity

- The fundamental restriction in ICA is that the independent components must be nongaussian for ICA to be possible.
- This is because gaussianity is invariant under orthogonal transformation and hence make the matrix \mathbf{A} not identifiable for gaussian independent components.



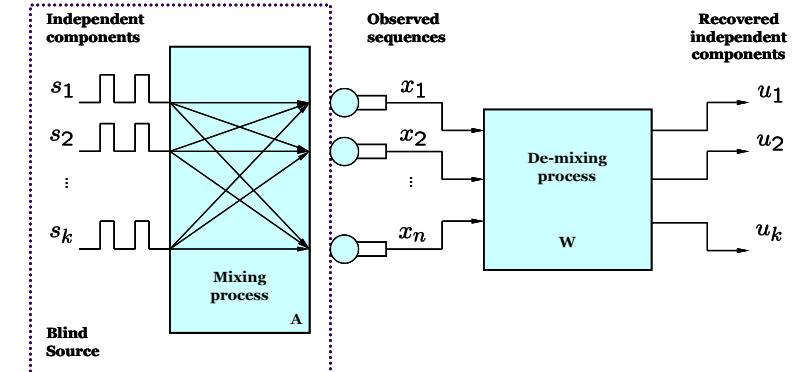
Measures of nongaussianity (1)

- Kurtosis
 - $\text{Kurt}(y) = E \{y^4\} - 3 (E \{y^2\})^2$
 - Kurtosis can be very sensitive to outliers, when its value has to be estimate from a measured sample.
- Mutual information
- Negative Entropy



FastICA — Preprocessing

- ❑ Centering:
 - ❑ Make the x -s mean 0 variables
- ❑ Whitening
 - ❑ Transform the observed vector x linearly so that it has unit variance:



$$E \left\{ \tilde{x} \tilde{x}^T \right\} = I$$

- ❑ One can show that:

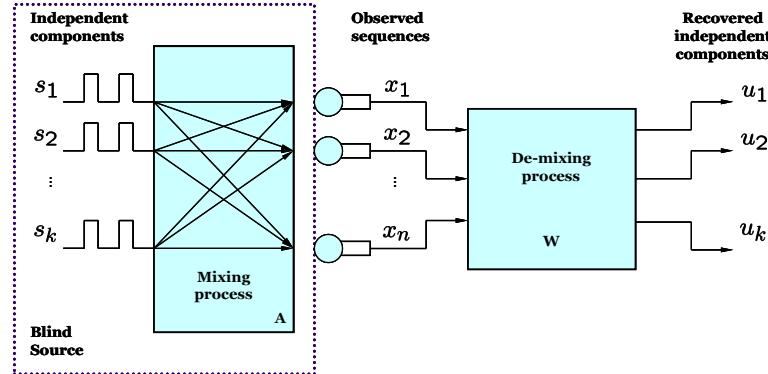
$$\tilde{x} = E D^{-1/2} E^T x = \tilde{A} s$$

where

$$E \left\{ x x^T \right\} = E D E^T$$



FastICA algorithm



- A fixed-point iteration scheme, that maximizes a measure of non-Gaussianity of the rotated components. Non-Gaussianity serves as a proxy for statistical independence
- To measure non-Gaussianity, FastICA relies on a nonquadratic nonlinear function $f(u)$, its first derivative $g(u)$, and its second derivative $g'(u)$.

$$f(u) = \log \cosh(u),$$

$$f(u) = -e^{-u^2/2}$$

Algorithm FastICA

Input: C Number of desired components

Input: $\mathbf{X} \in \mathbb{R}^{N \times M}$ Prewhitened matrix, where each column represents an N -dimensional sample, where $C \leq N$

Output: $\mathbf{W} \in \mathbb{R}^{N \times C}$ Un-mixing matrix where each column projects \mathbf{X} onto independent component.

Output: $\mathbf{S} \in \mathbb{R}^{C \times M}$ Independent components matrix, with M columns representing a sample with C dimensions.

for p **in** 1 to C :

$\mathbf{w}_p \leftarrow$ Random vector of length N

while \mathbf{w}_p changes

$$\mathbf{w}_p \leftarrow \frac{1}{M} \mathbf{X} g(\mathbf{w}_p^T \mathbf{X})^T - \frac{1}{M} g'(\mathbf{w}_p^T \mathbf{X}) \mathbf{1}_M \mathbf{w}_p$$

$$\mathbf{w}_p \leftarrow \mathbf{w}_p - \sum_{j=1}^{p-1} (\mathbf{w}_p^T \mathbf{w}_j) \mathbf{w}_j$$

$$\mathbf{w}_p \leftarrow \frac{\mathbf{w}_p}{\|\mathbf{w}_p\|}$$

output $\mathbf{W} \leftarrow [\mathbf{w}_1, \dots, \mathbf{w}_C]$

output $\mathbf{S} \leftarrow \mathbf{W}^T \mathbf{X}$



Summary

- ❑ There has been a wide discussion about the application of Independence Component Analysis (ICA) in Signal Processing, Neural Computation and Finance.
- ❑ First introduced as a novel tool to separate blind sources in a mixed signal.
- ❑ The Basic idea of ICA is to reconstruct from observation sequences the hypothesized independent original sequences.

