

# Course 495: Advanced Statistical Machine Learning/Pattern Recognition

---

- Goal (Lecture): To present Probabilistic Principal Component Analysis (PPCA) using both Maximum Likelihood (ML) and Expectation Maximization (EM).
- Goal (Tutorials): To provide the students the necessary mathematical tools for deeply understanding PPCA.

# Materials

---

- Pattern Recognition & Machine Learning by C. Bishop Chapter 12
- **PPCA:** Tipping, Michael E., and Christopher M. Bishop. "Probabilistic principal component analysis." *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 61.3 (1999): 611-622
- **PPCA:** Tipping, Michael E., and Christopher M. Bishop. "Mixtures of probabilistic principal component analyzers." *Neural computation* 11.2 (1999): 443-482.

# An overview

---

PCA: Maximize the global variance

LDA: Minimize the class variance while maximizing the mean variance

LPP: Minimize the local variance

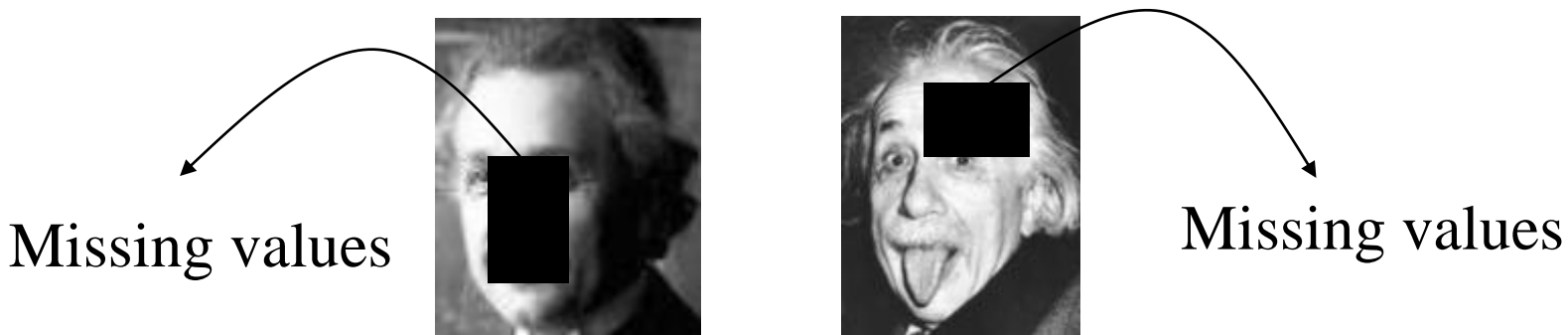
ICA: Maximize independence by maximizing non-Gaussianity

All are deterministic!!

# Advantages of PPCA over PCA

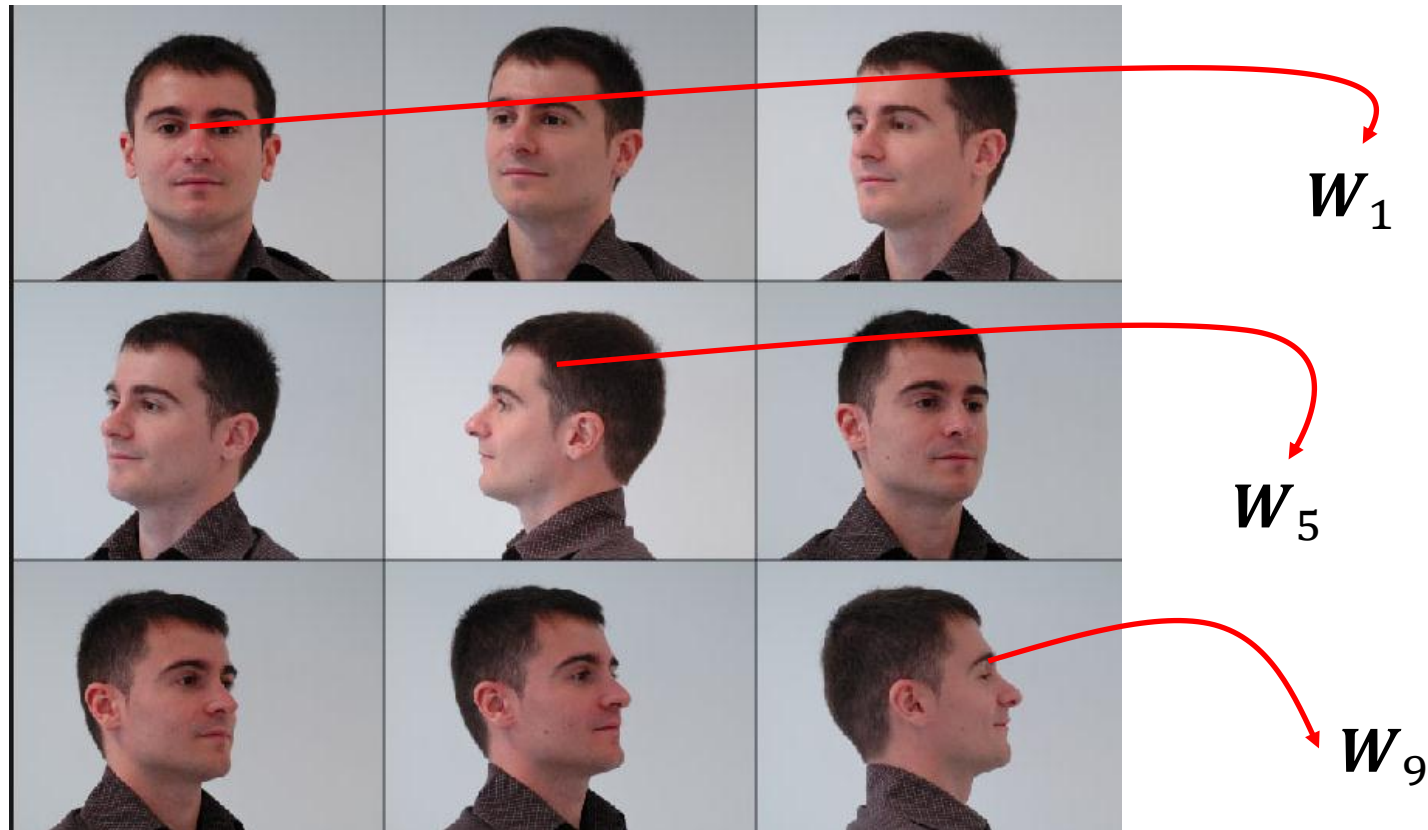
---

- An EM algorithm for PCA that is computationally efficient in situations where only a few leading eigenvectors are required and that avoids having to evaluate the data covariance matrix as an intermediate step.
- A combination of a probabilistic model and EM allows us to deal with missing values in the dataset.



# Advantages of PPCA over PCA

- Mixtures of a probabilistic PCA models can be formulated in a principled way and trained using the EM algorithm.



# Advantages of PPCA over PCA

---

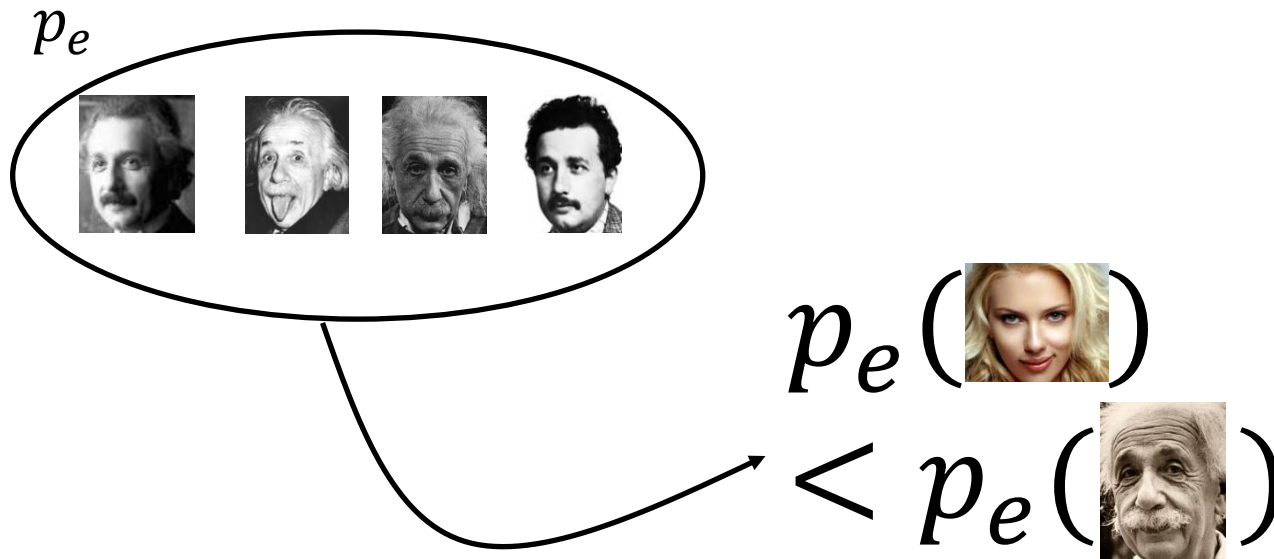
- Probabilistic PCA forms the basis for a Bayesian treatment of PCA in which the dimensionality of the principal subspace can be found automatically from the data.

Priors on  $\mathbf{W}$ : Automatic Relevance Determination  
(find the relevant components)

$$p(\mathbf{W}|\mathbf{a}) = \prod_{i=1}^d \left(\frac{a_i}{2\pi}\right)^{\frac{d}{2}} \exp\left\{-\frac{1}{2} a_i \mathbf{w}_i^T \mathbf{w}_i\right\}$$

# Advantages of PPCA over PCA

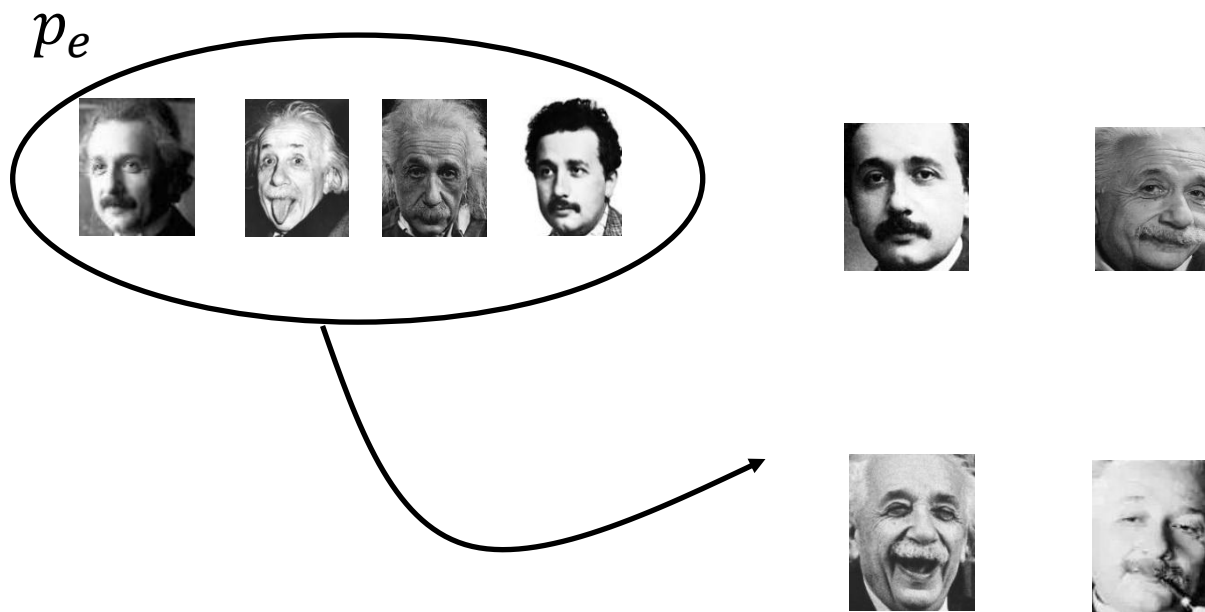
- Probabilistic PCA can be used to model class-conditional densities and hence be applied to classification problems.



# Advantages of PPCA over PCA

---

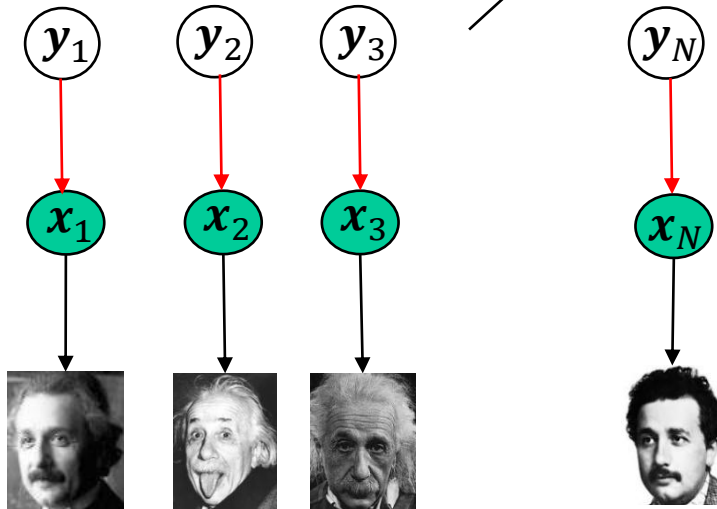
- The probabilistic PCA model can be run generatively to provide samples from the distribution.





# Probabilistic Principal Component Analysis

General Concept:



Share a common linear structure

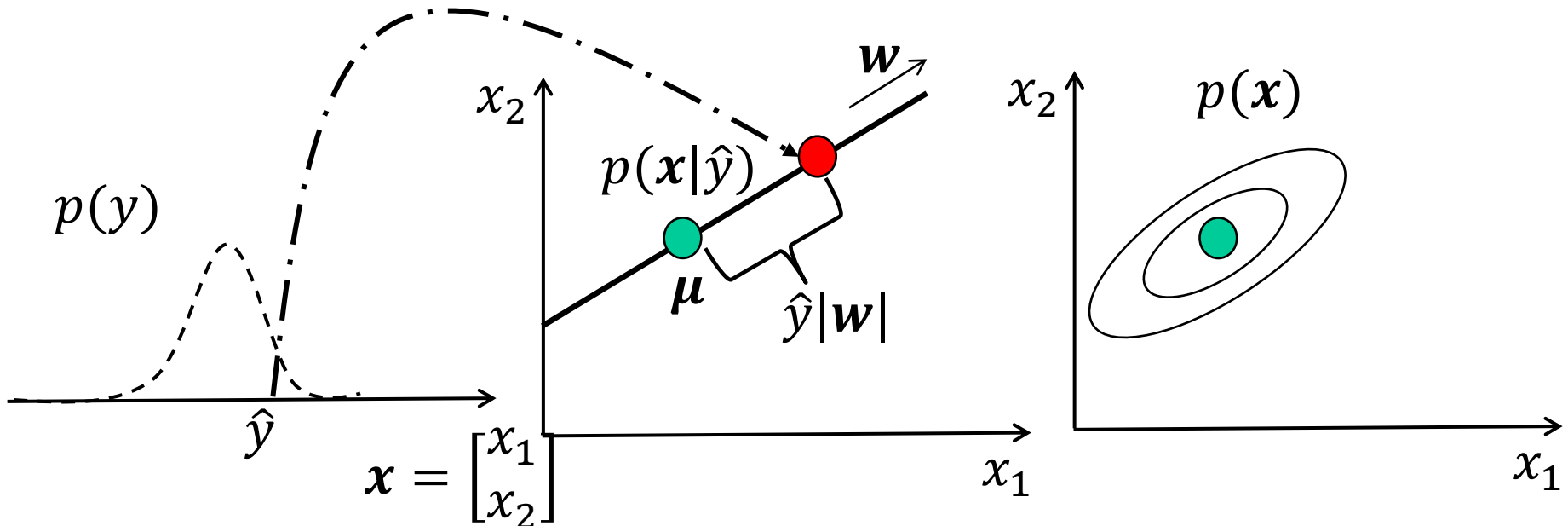
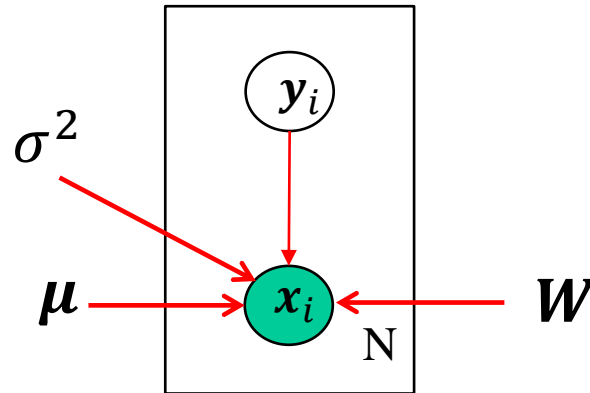
$$\begin{aligned}x &= \mathbf{W}y + \mu + e \\ e &\sim N(e | \mathbf{0}, \sigma^2 I) \\ y &\sim N(y | \mathbf{0}, I)\end{aligned}$$

We want to find the parameters:

$$\theta = \{\mathbf{W}, \mu, \sigma^2\}$$

# Probabilistic Principal Component Analysis

Graphical Model:



# ML-PPCA

---

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{y}_1, \dots, \mathbf{y}_N | \theta) = \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) \prod_{i=1}^N p(\mathbf{y}_i)$$

$$p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) = \mathcal{N}(\mathbf{x}_i | \mathbf{W}\mathbf{y}_i + \boldsymbol{\mu}, \sigma^2)$$

$$p(\mathbf{y}_i) = \mathcal{N}(\mathbf{y}_i | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta) =$$

$$\int_{\mathbf{y}_i} p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{y}_1, \dots, \mathbf{y}_N | \theta) d\mathbf{y}_1 \dots d\mathbf{y}_N$$
$$= \int_{\mathbf{y}_i} \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) \prod_{i=1}^N p(\mathbf{y}_i) d\mathbf{y}_1 \dots d\mathbf{y}_N$$

# ML-PPCA

---

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta) &= \int_{\mathbf{y}_i} \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) \prod_{i=1}^N p(\mathbf{y}_i) d\mathbf{y}_1 \dots d\mathbf{y}_N \\ &= \prod_{i=1}^N \int_{\mathbf{y}_i} p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) d\mathbf{y}_i \end{aligned}$$

$$\begin{aligned} p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) &= \frac{1}{\sqrt{(2\pi)^F \sigma^F}} e^{-\frac{1}{\sigma^2} (\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{y}_i)^T (\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{y}_i)} \\ &\quad \frac{1}{\sqrt{(2\pi)^d}} e^{-\mathbf{y}_i^T \mathbf{y}_i} \end{aligned}$$

# ML-PPCA

---

Complete the square:

$$\begin{aligned} \frac{1}{\sigma^2} \left[ \underbrace{(\mathbf{x}_i - \boldsymbol{\mu})}_{\bar{\mathbf{x}}_i} - \mathbf{W}\mathbf{y}_i \right]^T \left[ \underbrace{(\mathbf{x}_i - \boldsymbol{\mu})}_{\bar{\mathbf{x}}_i} - \mathbf{W}\mathbf{y}_i \right] + \sigma^2 \mathbf{y}_i^T \mathbf{y}_i &= \\ &= (\bar{\mathbf{x}}_i - \mathbf{W}\mathbf{y}_i)^T (\bar{\mathbf{x}}_i - \mathbf{W}\mathbf{y}_i) + \sigma^2 \mathbf{y}_i^T \mathbf{y}_i \\ &= \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - 2\bar{\mathbf{x}}_i^T \mathbf{W}\mathbf{y}_i + \mathbf{y}_i^T \mathbf{W}^T \mathbf{W} \mathbf{y}_i + \sigma^2 \mathbf{y}_i^T \mathbf{y}_i \\ &= \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - 2\bar{\mathbf{x}}_i^T \mathbf{W}\mathbf{y}_i + \mathbf{y}_i^T \underbrace{(\sigma^2 \mathbf{I} + \mathbf{W}^T \mathbf{W})}_{\mathbf{M}} \mathbf{y}_i \end{aligned}$$

# ML-PPCA

---

$$= \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - 2\bar{\mathbf{x}}_i^T \mathbf{W} \mathbf{y}_i + \mathbf{y}_i^T \mathbf{M} \mathbf{y}_i$$

$$= \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - 2(\mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)^T \mathbf{M} \mathbf{y}_i + \mathbf{y}_i^T \mathbf{M} \mathbf{y}_i + (\mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)^T \mathbf{M} \mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i - (\mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)^T \mathbf{M} \mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i$$

$$= \bar{\mathbf{x}}_i^T (\mathbf{I} - \mathbf{W} \mathbf{M}^{-1} \mathbf{W}^T) \bar{\mathbf{x}}_i + (\mathbf{y}_i - \mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)^T \mathbf{M} (\mathbf{y}_i - \mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)$$

# ML-PPCA

---

$$\begin{aligned} & \int_{\mathbf{y}_i} p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) d\mathbf{y}_i \propto \\ & \int_{\mathbf{y}_i} e^{-\frac{1}{\sigma^2} \bar{\mathbf{x}}_i^T (\mathbf{I} - \mathbf{W} \mathbf{M}^{-1} \mathbf{W}^T) \bar{\mathbf{x}}_i - \frac{1}{\sigma^2} (\mathbf{y}_i - \mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)^T \mathbf{M} (\mathbf{y}_i - \mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)} d\mathbf{y}_i \propto \\ & e^{-\frac{1}{\sigma^2} \bar{\mathbf{x}}_i^T (\mathbf{I} - \mathbf{W} \mathbf{M}^{-1} \mathbf{W}^T) \bar{\mathbf{x}}_i} \int_{\mathbf{y}_i} \underbrace{e^{-\frac{1}{\sigma^2} (\mathbf{y}_i - \mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)^T \mathbf{M} (\mathbf{y}_i - \mathbf{M}^{-1} \mathbf{W}^T \bar{\mathbf{x}}_i)}}_1 d\mathbf{y}_i \\ & = N(\mathbf{x}_i | \boldsymbol{\mu}, (\sigma^{-2} \mathbf{I} - \sigma^{-2} \mathbf{W} \mathbf{M}^{-1} \mathbf{W}^T)^{-1}) \end{aligned}$$

# ML-PPCA

---

Let's have a look at:

$$\sigma^{-2}\mathbf{I} - \sigma^{-2}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^T = \sigma^{-2}\mathbf{I} - \sigma^{-2}\mathbf{W}(\sigma^2\mathbf{I} + \mathbf{W}^T\mathbf{W})^{-1}\mathbf{W}^T$$

Woodbury formula:

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}$$

Can you prove that if:  $\mathbf{D} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$

then  $\mathbf{D}^{-1} = \sigma^{-2}\mathbf{I} - \sigma^{-2}\mathbf{W}(\sigma^2\mathbf{I} + \mathbf{W}^T\mathbf{W})^{-1}\mathbf{W}^T$



# ML-PPCA

---

$$p(\mathbf{x}_i | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \int_{\mathbf{y}_i} p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) d\mathbf{y}_i = N(\mathbf{x}_i | \boldsymbol{\mu}, \mathbf{D})$$
$$\mathbf{D} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$$

Using also the above we can easily show that the posterior :

$$p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \frac{p(\mathbf{y}_i)}{p(\mathbf{x}_i | \mathbf{W}, \boldsymbol{\mu}, \sigma^2)} p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) =$$
$$= N(\mathbf{y}_i | \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_i - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1})$$

# ML-PPCA

---

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta) = \prod_{i=1}^N N(\mathbf{x}_i | \boldsymbol{\mu}, \mathbf{D})$$

$$\theta = \operatorname{argmax}_{\theta} \ln p(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta)$$

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^N \ln N(\mathbf{x}_i | \boldsymbol{\mu}, \mathbf{D})$$

$$= -\frac{NF}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{D}| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{D}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$$

# ML-PPCA

---

$$L(\mathbf{W}, \sigma^2, \boldsymbol{\mu}) = -\frac{NF}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{D}| - N \text{tr}[\mathbf{D}^{-1} \mathbf{S}_t]$$

where  $\mathbf{S}_t = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$

$$\frac{dL}{d\boldsymbol{\mu}} = \mathbf{0} \Rightarrow \boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

# ML-PPCA

---

$$\frac{dL}{d\mathbf{W}} = N(\mathbf{D}^{-1}\mathbf{S}_t\mathbf{D}^{-1}\mathbf{W} - \mathbf{D}^{-1}\mathbf{W}) \Rightarrow \mathbf{S}_t\mathbf{D}^{-1}\mathbf{W} = \mathbf{W}$$

Three different solutions:

(1)  $\mathbf{W} = \mathbf{0}$  which is the minimum of the likelihood

$$(2) \mathbf{D} = \mathbf{S}_t \Rightarrow \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I} = \mathbf{S}_t$$

Assume the eigen-decomposition  $\mathbf{S}_t = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$

$\mathbf{U}$  square matrix of eigenvectors  $\mathbf{W} = \mathbf{U}(\mathbf{\Lambda} - \sigma^2\mathbf{I})^{1/2}\mathbf{R}$

$\mathbf{R}$  is a rotation matrix  $\mathbf{R}^T\mathbf{R} = \mathbf{I}$

# ML-PPCA

---

$$(3) \mathbf{D} \neq \mathbf{S}_t \text{ and } \mathbf{W} \neq 0 \quad d < q = \text{rank}(\mathbf{S}_t)$$

Assume the SVD of  $\mathbf{W} = \mathbf{U}\mathbf{L}\mathbf{V}^T$

$$\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_d] \text{ } F \times d \text{ matrix} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I} \quad \mathbf{V}^T \mathbf{V} = \mathbf{V}\mathbf{V}^T = \mathbf{I}$$

$$\mathbf{L} = \begin{bmatrix} l_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & l_d \end{bmatrix}$$

$$\mathbf{S}_t \mathbf{D}^{-1} \mathbf{U}\mathbf{L}\mathbf{V}^T = \mathbf{U}\mathbf{L}\mathbf{V}^T$$

Let's study  $\mathbf{D}^{-1} \mathbf{U}$

# ML-PPCA

---

$$\begin{aligned} D^{-1} &= (WW^T + \sigma^2 \mathbf{I})^{-1} \xrightarrow{W=ULV^T} \\ &= (\mathbf{U}L^2\mathbf{U}^T + \sigma^2\mathbf{I})^{-1} \end{aligned}$$

Assume a set of bases  $\mathbf{U}_{F-d}$  such that  $\mathbf{U}_{F-d}^T \mathbf{U} = 0$

and  $\mathbf{U}_{F-d}^T \mathbf{U}_{F-d} = \mathbf{I}$

$$\begin{aligned} &= \left( [\mathbf{U} \ \mathbf{U}_{F-d}] \begin{bmatrix} \mathbf{L}^2 & 0 \\ 0 & 0 \end{bmatrix} [\mathbf{U} \ \mathbf{U}_{F-d}]^T + [\mathbf{U} \ \mathbf{U}_{F-d}] \sigma^2 \mathbf{I} [\mathbf{U} \ \mathbf{U}_{F-d}]^T \right)^{-1} \\ &= [\mathbf{U} \ \mathbf{U}_{F-d}] \begin{bmatrix} \mathbf{L}^2 + \sigma^2 \mathbf{I} & 0 \\ 0 & \sigma^2 \mathbf{I} \end{bmatrix}^{-1} [\mathbf{U} \ \mathbf{U}_{F-d}]^T \\ &= [\mathbf{U} \ \mathbf{U}_{F-d}] \begin{bmatrix} (\mathbf{L}^2 + \sigma^2 \mathbf{I})^{-1} & 0 \\ 0 & \sigma^{-2} \mathbf{I} \end{bmatrix} [\mathbf{U} \ \mathbf{U}_{F-d}]^T \end{aligned}$$

# ML-PPCA

---

$$\begin{aligned} D^{-1}U &= [U \ U_{F-d}] \begin{bmatrix} (L^2 + \sigma^2 I)^{-1} & 0 \\ 0 & \sigma^{-2} I \end{bmatrix} [U \ U_{F-d}]^T U \\ &= [U \ U_{F-d}] \begin{bmatrix} (L^2 + \sigma^2 I)^{-1} & 0 \\ 0 & \sigma^{-2} I \end{bmatrix} [I \ 0]^T \\ &= [U \ U_{F-d}] \begin{bmatrix} (L^2 + \sigma^2 I)^{-1} \\ 0 \end{bmatrix} \\ &= U(L^2 + \sigma^2 I)^{-1} \end{aligned}$$

# ML-PPCA

---

$$\mathbf{S}_t \mathbf{D}^{-1} \mathbf{U} \mathbf{L} \mathbf{V}^T = \mathbf{U} \mathbf{L} \mathbf{V}^T$$

$$\mathbf{S}_t \mathbf{U} (\mathbf{L}^2 + \sigma^2 \mathbf{I})^{-1} = \mathbf{U}$$

$$\mathbf{S}_t \mathbf{U} = \mathbf{U} (\mathbf{L}^2 + \sigma^2 \mathbf{I})^{-1}$$

It means that

$$\mathbf{S}_t \mathbf{u}_i = (l_i^2 + \sigma^2) \mathbf{u}_i \quad \mathbf{S}_t = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$$

$\mathbf{u}_i$  are eigenvectors of  $\mathbf{S}_t$  and  $\lambda_i = l_i^2 + \sigma^2 \Rightarrow l_i = \sqrt{\lambda_i - \sigma^2}$

Unfortunately  $\mathbf{V}$  cannot be determined hence there is a rotation ambiguity.



# ML-PPCA

---

Hence the optimum is given by (keeping  $d$  eigenvectors)

$$\mathbf{W}_d = \mathbf{U}_d(\mathbf{\Lambda}_d - \sigma^2 \mathbf{I})\mathbf{V}^T$$

Having computed  $\mathbf{W}$  we need to compute the optimum  $\sigma^2$

$$\begin{aligned} L(\mathbf{W}, \sigma^2, \boldsymbol{\mu}) &= -\frac{NF}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{D}| - N\text{tr}[\mathbf{D}^{-1}\mathbf{S}_t] \\ &= -\frac{NF}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}| - N\text{tr}[(\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})^{-1}\mathbf{S}_t] \end{aligned}$$

# ML-PPCA

---

$$\begin{aligned}\mathbf{W}_d \mathbf{W}_d^T + \sigma^2 \mathbf{I} &= [\mathbf{U}_d \mathbf{U}_{F-d}] \begin{bmatrix} \boldsymbol{\Lambda}_d - \sigma^2 \mathbf{I} & 0 \\ 0 & 0 \end{bmatrix} [\mathbf{U}_d \mathbf{U}_{F-d}]^T \\ &\quad + [\mathbf{U}_d \mathbf{U}_{F-d}] \begin{bmatrix} \sigma^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma^2 \end{bmatrix} [\mathbf{U}_d \mathbf{U}_{F-d}]^T \\ &= [\mathbf{U}_d \mathbf{U}_{F-d}] \begin{bmatrix} \boldsymbol{\Lambda}_d & 0 \\ 0 & \sigma^2 \mathbf{I} \end{bmatrix} [\mathbf{U}_d \mathbf{U}_{F-d}]^T\end{aligned}$$

Hence  $|\mathbf{W}_d \mathbf{W}_d^T + \sigma^2 \mathbf{I}| = \prod_{i=1}^d \lambda_i \prod_{i=d+1}^F \sigma^2$

$$\ln |\mathbf{W}_d \mathbf{W}_d^T + \sigma^2 \mathbf{I}| = (F-d) \ln \sigma^2 + \sum_{i=1}^d \ln \lambda_i$$

# ML-PPCA

---

$$\mathbf{D}^{-1}\mathbf{S}_t = [\mathbf{U}_d \mathbf{U}_{F-d}] \begin{bmatrix} \boldsymbol{\Lambda}_d & 0 \\ 0 & \sigma^2 \mathbf{I} \end{bmatrix}^{-1} [\mathbf{U}_d \mathbf{U}_{F-d}]^T [\mathbf{U}_d \mathbf{U}_{F-d}] \begin{bmatrix} \boldsymbol{\Lambda}_d & 0 & 0 \\ 0 & \boldsymbol{\Lambda}_{q-d} & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$= [\mathbf{U}_d \mathbf{U}_{F-d}] \begin{bmatrix} \mathbf{I} & 0 & 0 \\ 0 & \frac{1}{\sigma^2} \boldsymbol{\Lambda}_{q-d} & 0 \\ 0 & 0 & 0 \end{bmatrix} [\mathbf{U}_d \mathbf{U}_{F-d}]^T$$

$$\Rightarrow \text{tr}(\mathbf{D}^{-1}\mathbf{S}_t) = \frac{1}{\sigma^2} \sum_{i=d+1}^q \lambda_i + d$$

# ML-PPCA

---

$$L(\sigma^2) = -\frac{N}{2} \left\{ F \ln 2\pi + \sum_{j=1}^d \ln \lambda_j + \frac{1}{\sigma^2} \sum_{j=d+1}^q \lambda_j + (F - d) \ln \sigma^2 + d \right\}$$

$$\frac{\partial L}{\partial \sigma} = 0 \Rightarrow -2\sigma^{-3} \sum_{j=d+1}^q \lambda_j + \frac{F - d}{\sigma} = 0 \Rightarrow \sigma^2 = \frac{1}{F - d} \sum_{j=d+1}^q \lambda_j$$

If I put the solution back

$$L(\sigma^2) = -\frac{N}{2} \left\{ \sum_{j=1}^d \ln \lambda_j + (F - q) \ln \frac{1}{F - d} \sum_{j=d+1}^q \lambda_j + F \ln 2\pi + F \right\}$$

# ML-PPCA

---

$$L(\sigma^2) = -\frac{N}{2} \left\{ \underbrace{\sum_{j=1}^d \ln \lambda_j + \sum_{j=d}^q \ln \lambda_j}_{\ln |\mathbf{S}_t|} - \sum_{j=d}^q \ln \lambda_j + (F-d) \ln \frac{1}{F-d} \sum_{j=d+1}^F \lambda_j + F \ln 2\pi + F \right\}$$

$$\max \frac{N}{2} \left\{ \frac{1}{F-d} \ln |\mathbf{S}_t| - \frac{1}{F-d} \sum_{j=d}^q \ln \lambda_j + \ln \left( \frac{1}{F-d} \sum_{j=d+1}^F \lambda_j \right) + \text{cost} \right\}$$

$$\Rightarrow \min \ln \left( \frac{1}{F-d} \sum_{j=d}^q \ln \lambda_j \right) - \frac{1}{F-d} \sum_{j=d}^q \ln \lambda_j$$

# ML-PPCA

---

Jensen inequality  $\ln \left( \frac{\sum_{i=1}^n r_i}{n} \right) \geq \frac{1}{n} \sum_{i=1}^n \ln r_i$

$$\ln \left( \frac{1}{F-d} \sum_{j=d+1}^q \lambda_j \right) \geq \frac{1}{F-d} \sum_{j=d+1}^q \ln \lambda_j \text{ hence}$$

$$\Rightarrow \ln \left( \frac{1}{F-d} \sum_{j=d}^q \ln \lambda_j \right) - \frac{1}{F-d} \sum_{j=d}^q \ln \lambda_j \geq 0$$

Hence, the function is minimized when the discarded eigenvectors are the ones that correspond to the  $q - d$  eigenvalues

# ML-PPCA

---

Summarize:

$$\sigma^2 = \frac{1}{F - d} \sum_{j=d+1}^q \lambda_j$$

$$\mathbf{W}_d = \mathbf{U}_d(\mathbf{\Lambda}_d - \sigma^2 \mathbf{I})\mathbf{V}^T$$

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

We no longer have a projection but:

$$\mathbf{E}_{p(\mathbf{y}_i|\mathbf{x}_i)}\{\mathbf{y}_i\} = \mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x}_i - \boldsymbol{\mu})$$

and a reconstruction

$$\hat{\mathbf{x}}_i = \mathbf{W}\mathbf{E}_{p(\mathbf{y}_i|\mathbf{x}_i)}\{\mathbf{y}_i\} + \boldsymbol{\mu}$$

# ML-PPCA

---

$$\lim_{\sigma^2 \rightarrow 0} \mathbf{W}_d = \mathbf{U}_d \mathbf{\Lambda}_d$$

$$\lim_{\sigma^2 \rightarrow 0} \mathbf{M} = \mathbf{W}_d^T \mathbf{W}_d$$

Hence

$$\begin{aligned} \lim_{\sigma^2 \rightarrow 0} \mathbf{E}_{p(\mathbf{y}_i | \mathbf{x}_i)} \{\mathbf{y}_i\} &= \mathbf{M}^{-1} \mathbf{W}_d^T (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= \mathbf{\Lambda}_d^{-1} \mathbf{U}_d (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned}$$

which gives the whitened PCA



# EM-PPCA

---

First step formulate the joint likelihood

$$p(\mathbf{X}, \mathbf{Y} | \theta) = p(\mathbf{X}, \mathbf{Y} | \theta) p(\mathbf{Y}) = \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{y}_i, \theta) p(\mathbf{y}_i)$$

$$\ln p(\mathbf{X}, \mathbf{Z} | \theta) = \sum_{i=1}^N \{ \ln p(\mathbf{x}_i | \mathbf{y}_i, \theta) + \ln p(\mathbf{y}_i) \}$$

$$\begin{aligned} \ln p(\mathbf{x}_i | \mathbf{y}_i, \theta) &= \ln \frac{1}{\sqrt{(2\pi)^F (\sigma^2)^F}} e^{-\frac{1}{2\sigma^2} (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu})^T (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu})} \\ &= -\frac{F}{2} \ln(\sqrt{2\pi}\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu})^T (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu}) \end{aligned}$$

$$\ln p(\mathbf{y}_i) = -\frac{1}{2} \mathbf{y}_i^T \mathbf{y}_i - \frac{D}{2} \ln 2\pi$$

# EM-PPCA

---

$$\begin{aligned} \ln p(\mathbf{X}, \mathbf{Y} | \theta) &= \sum_{i=1}^N \left\{ -\frac{F}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu})^T (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i \right. \\ &\quad \left. - \boldsymbol{\mu}) - \frac{1}{2} \mathbf{y}_i^T \mathbf{y}_i - \frac{D}{2} \ln 2\pi \right\} \end{aligned}$$

I need now to optimize it with regards to  $\mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma^2$

We can't hence we need to take the expectations

# EM-PPCA

---

But first we need to expand a bit:

$$\begin{aligned}\ln p(\mathbf{X}, \mathbf{Y}|\theta) &= -\frac{NF}{2} \ln 2\pi\sigma^2 - \frac{ND}{2} \ln 2\pi \\ &\quad - \sum_{i=1}^N \left\{ \frac{1}{2\sigma^2} [(\mathbf{x}_i - \boldsymbol{\mu})^T (\mathbf{x}_i - \boldsymbol{\mu}) - 2(\mathbf{W}^T \mathbf{x}_i)^T \mathbf{y}_i \right. \\ &\quad \left. + \mathbf{y}_i^T \mathbf{W}^T \mathbf{W} \mathbf{y}_i] \right\} + \frac{1}{2} \mathbf{y}_i^T \mathbf{y}_i\end{aligned}$$

# EM-PPCA

---

$$\mathbf{E}_{P(\mathbf{Y}|\mathbf{X})}\{\log P(\mathbf{X}, \mathbf{Y}|\theta)\} = -\frac{NF}{2} \ln 2\pi\sigma^2 - \frac{ND}{2} \ln 2\pi$$
$$- \sum_{i=1}^N \left\{ \frac{1}{2\sigma^2} [(\mathbf{x}_i - \boldsymbol{\mu})^T (\mathbf{x}_i - \boldsymbol{\mu}) - 2(\mathbf{W}^T \mathbf{x}_i)^T \mathbf{E}\{\mathbf{y}_i\} + \text{tr}[\mathbf{E}\{\mathbf{y}_i \mathbf{y}_i^T\} \mathbf{W}^T \mathbf{W}]] \right\}$$

# EM-PPCA

---

$$\begin{aligned} E\{\mathbf{y}_i\} &= \int \mathbf{y}_i p(\mathbf{y}_i | \mathbf{x}_i, \theta) d\mathbf{y}_n \\ &= E\{N(\mathbf{y}_i | \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_i - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1})\} \\ &= \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned}$$

$$\begin{aligned} E\{(\mathbf{y}_i - E\{\mathbf{y}_i\})(\mathbf{y}_i - E\{\mathbf{y}_i\})^T\} &= \\ &= E\{(\mathbf{y}_i \mathbf{y}_i^T)\} - E\{E\{\mathbf{y}_i\} \mathbf{y}_i^T\} - E\{\mathbf{y}_i E\{\mathbf{y}_i\}^T\} + E\{\mathbf{y}_i\} E\{\mathbf{y}_i\}^T \\ &= E\{\mathbf{y}_i \mathbf{y}_i^T\} - E\{\mathbf{y}_i\} E\{\mathbf{y}_i\}^T - E\{\mathbf{y}_i\} E\{\mathbf{y}_i\}^T + E\{\mathbf{y}_i\} E\{\mathbf{y}_i\}^T \\ &= E\{\mathbf{y}_i \mathbf{y}_i^T\} - E\{\mathbf{y}_i\} E\{\mathbf{y}_i\}^T = \sigma^2 \mathbf{M}^{-1} \\ \text{hence } E\{\mathbf{y}_i \mathbf{y}_i^T\} &= \sigma^2 \mathbf{M}^{-1} + E\{\mathbf{y}_i\} E\{\mathbf{y}_i\}^T \end{aligned}$$

# EM-PPCA

---

Expectation Step:

Given  $\mathbf{W}$ ,  $\boldsymbol{\mu}$ ,  $\sigma$ :

$$E\{\mathbf{y}_i\} = \mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x}_i - \boldsymbol{\mu})$$

$$E\{\mathbf{y}_i\mathbf{y}_i^T\} = \sigma^2\mathbf{M}^{-1} + E\{\mathbf{y}_i\}E\{\mathbf{y}_i\}^T$$

# EM-PPCA

---

Maximization step

$$\frac{\partial E\{\boldsymbol{\mu}\}}{\partial \boldsymbol{\mu}} = \mathbf{0} \Rightarrow \boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

$$\begin{aligned} \frac{\partial E\{\mathbf{W}\}}{\partial \mathbf{W}} = \mathbf{0} &\Rightarrow \sum_{i=1}^N \left( -\frac{1}{\sigma^2} (\mathbf{x}_i - \boldsymbol{\mu}) E\{\mathbf{y}_i\}^T + \frac{2}{2\sigma^2} \mathbf{W} E\{\mathbf{y}_i \mathbf{y}_i^T\} \right) = \mathbf{0} \\ &\Rightarrow \mathbf{W} = \left[ \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) E\{\mathbf{y}_i\}^T \right] \left[ \sum_{i=1}^N E\{\mathbf{y}_i \mathbf{y}_i^T\} \right]^{-1} \end{aligned}$$

# EM-PPCA

---

$$\frac{\partial \mathbf{E}\{\sigma^2\}}{\partial \sigma} = \mathbf{0} \Rightarrow$$

$\sigma^2$

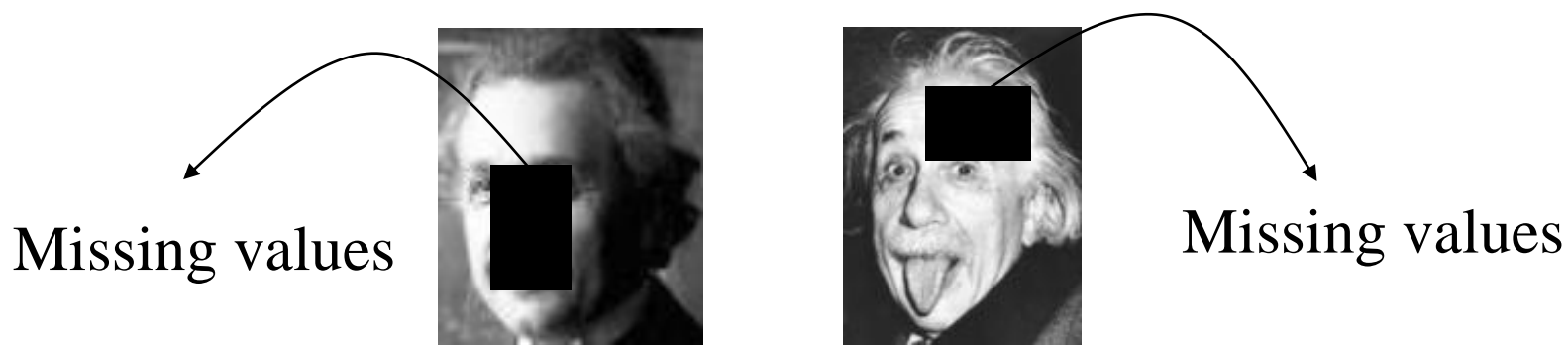
$$= \frac{1}{NF} \sum_{i=1}^N \{ \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 - 2E\{\mathbf{y}_i\}^T \mathbf{W}^T (\mathbf{x}_i - \boldsymbol{\mu}) + \text{tr}[E\{\mathbf{y}_i \mathbf{y}_i^T\} \mathbf{W}^T \mathbf{W}] \}$$



# Why EM-PPCA?

---

- A complexity of  $O(NFd)$  which can be significant smaller than  $O(NF^2)$  (computation of the covariance)
- A combination of a probabilistic model and EM allows us to deal with missing values in the dataset.



# Why EM-PPCA?

---

A combination of a probabilistic model and EM allows us to deal with missing values in the dataset.

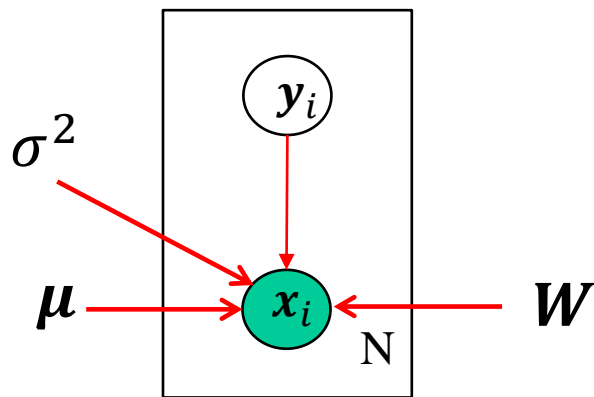


$$\begin{aligned} p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}_1, \boldsymbol{\mu}_1, \sigma) \\ \vdots \\ p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}_9, \boldsymbol{\mu}_9, \sigma) \end{aligned}$$

# Summary

---

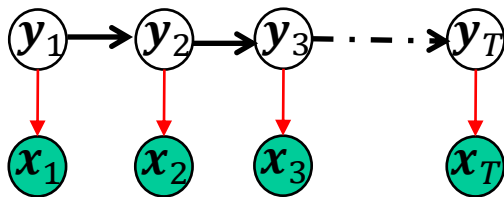
- We saw how to perform parameter estimation and inference using ML and EM in the following graphical model



# What will we see next?

---

Dynamic data



EM in the following graphical models

Spatial data

