

Machine Learning

Volker Roth

Department of Mathematics & Computer Science
University of Basel

Chapter 10: Linear latent variable models

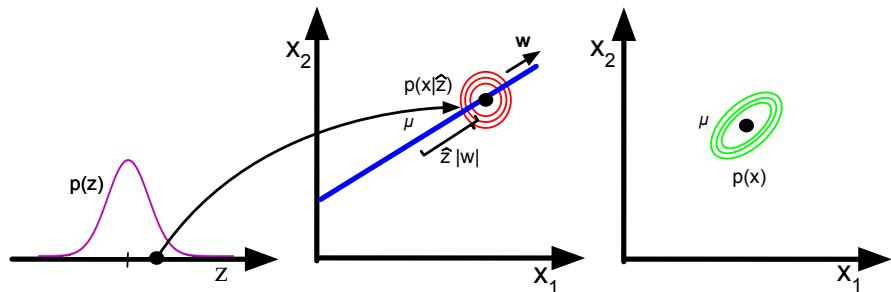
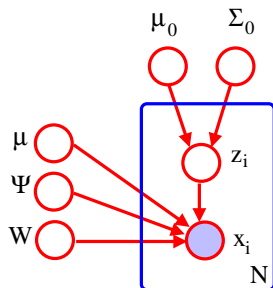


Figure 12.1 in K. Murphy: Machine Learning. MIT Press 2012.

Factor analysis

- One problem with mixture models: **only a single latent variable**. Each observation can only come from one of K prototypes.
- Alternative: $\mathbf{z}_i \in \mathbb{R}^k$. Gaussian prior:

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$



- For observations $\mathbf{x}_i \in \mathbb{R}^p$, we may use a **Gaussian likelihood**.
- As in linear regression, we assume the mean is a **linear** function:

$$p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(W\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}),$$

W : **factor loading matrix**, and $\boldsymbol{\Psi}$: **covariance matrix**.

- We take $\boldsymbol{\Psi}$ to be **diagonal**, since the whole point of the model is to “force” \mathbf{z}_i to **explain the correlation**.

Factor analysis: generative process

Generative process ($k = 1$, $p = 2$, diagonal Ψ):

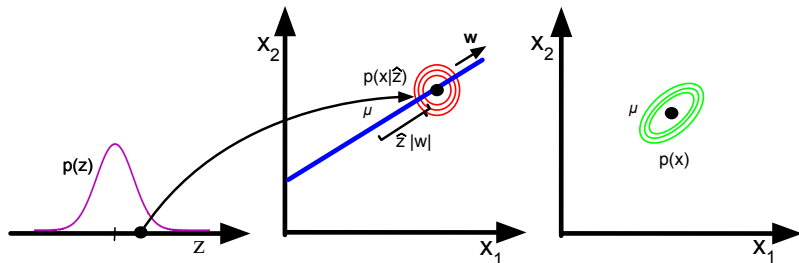


Figure 12.1 in K. Murphy: Machine Learning. MIT Press 2012.

We take an isotropic Gaussian “spray can” and slide it along the 1d line defined by $wz_i + \mu$. This induces a correlated Gaussian in 2d.

FA is a low rank parameterization of an MVN

- The induced marginal distribution $p(\mathbf{x}_i|\boldsymbol{\theta})$ is Gaussian:

$$\begin{aligned} p(\mathbf{x}_i|\boldsymbol{\theta}) &= \int \mathcal{N}(\mathbf{x}_i, W\mathbf{z}_i + \boldsymbol{\mu}, \Psi) \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \Sigma_0) d\mathbf{z}_i, \\ &= \mathcal{N}(\mathbf{x}_i|W\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \Psi + W\Sigma_0W^t) \end{aligned}$$

- We can set $\boldsymbol{\mu}_0 = \mathbf{0}$ without loss of generality (absorb $W\boldsymbol{\mu}_0$ into $\boldsymbol{\mu}$). Similarly, we can set $\Sigma_0 = I$ using $\tilde{W} = W\Sigma_0^{-1/2}$.

Covariance structure:

$$\text{cov}[\mathbf{x}|\boldsymbol{\theta}] = (W\Sigma_0^{-1/2})\Sigma_0(W\Sigma_0^{-1/2})^t + \Psi = WW^t + \Psi.$$

- We thus see that FA approximates the covariance matrix of \mathbf{x} using a low-rank decomposition:

$$C := \text{cov}[\mathbf{x}] = WW^t + \Psi.$$

Only $O(p \cdot k)$ parameters \rightsquigarrow compromise between full covariance with $O(p^2)$ params, and a diagonal covariance $O(p)$ params.

Inference of the latent factors

- We hope that the latent factors z will reveal something interesting about the data \rightsquigarrow compute posterior over the latent variables:

$$\begin{aligned}p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{z}_i | \mathbf{m}_i, \Sigma) \\ \Sigma &= (I + W^t \Psi^{-1} W)^{-1} \\ \mathbf{m}_i &= \Sigma W^t \Psi^{-1} \mathbf{x}_i\end{aligned}$$

- The posterior means \mathbf{m}_i are called the **latent scores, or latent factors.**

Example

- Example from (Shalizi 2009). $p = 11$ variables and $n = 387$ cases describing aspects of cars: engine size, #(cylinders), miles per gallon (MPG), price, etc.
- Fit a $p = 2$ dim model. Plot m_i scores as points in \mathbb{R}^2 .
- To get a better understanding of the “meaning” of the latent factors, project unit vectors $e_1 = (1, 0, \dots, 0)$, $e_2 = (0, 1, 0, \dots, 0)$, etc. into the low dimensional space (blue lines)
- Horizontal axis represents price, corresponding to the features labeled “dealer” and “retail”, with expensive cars on the right. Vertical axis represents fuel efficiency (measured in terms of MPG) versus size: heavy vehicles are less efficient and are higher up, whereas light vehicles are more efficient and are lower down.
- Verify by finding the closest exemplars in the training set.

Example

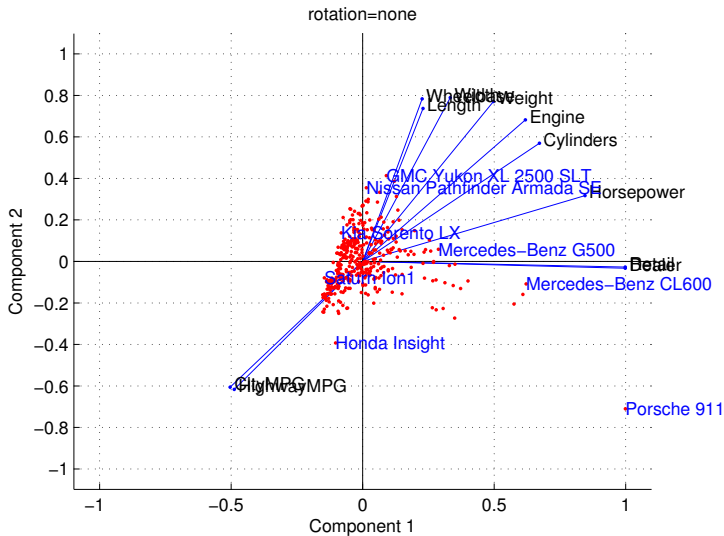
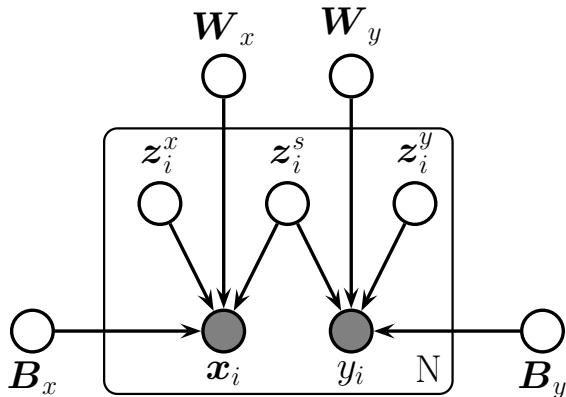


Figure 12.2 in K. Murphy: Machine Learning. MIT Press 2012.

Special Cases: PCA and CCA

- Covariance matrix $\Psi = \sigma^2 I \rightsquigarrow$ (probabilistic) **PCA**.
- Two-view version involving \mathbf{x} and $\mathbf{y} \rightsquigarrow$ **CCA**.



From figure 12.19 in K. Murphy: Machine Learning. MIT Press 2012.

PCA and dimensionality reduction

Given n data points in p dimensions:

$$X = \begin{bmatrix} - & \mathbf{x}_1 & - \\ - & \mathbf{x}_2 & - \\ - & \vdots & - \\ - & \mathbf{x}_n & - \end{bmatrix} \in \mathbb{R}^{n \times p}$$

Want to reduce dimensionality from p to k . Choose k directions $\mathbf{w}_1, \dots, \mathbf{w}_k$, arrange them as columns in matrix W :

$$W = [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \dots \quad \mathbf{w}_k] \in \mathbb{R}^{p \times k}$$

For each \mathbf{w}_j , compute **similarity** $z_j = \mathbf{w}_j^t \mathbf{x}$, $j = 1 \dots k$.

Project \mathbf{x} down to $\mathbf{z} = (z_1, \dots, z_k)^t = W^t \mathbf{x}$. How to choose W ?

Encoding–decoding model

The projection matrix W serves two functions:

- **Encode:** $\mathbf{z} = W^t \mathbf{x}$, $\mathbf{z} \in \mathbb{R}^k$, $z_j = \mathbf{w}_j^t \mathbf{x}$.
 - ▶ The vectors \mathbf{w}_j form a basis of the projected space.
 - ▶ We will require that this basis is orthonormal, i.e. $W^t W = I$.
- **Decode:** $\tilde{\mathbf{x}} = W \mathbf{z} = \sum_{j=1}^k z_j \mathbf{w}_j$, $\tilde{\mathbf{x}} \in \mathbb{R}^p$.
 - ▶ If $k = p$, the above orthonormality condition implies $W^t = W^{-1}$, and encoding can be undone without loss of information.
 - ▶ If $k < p$, we use the pseudo-inverse
 - ↪ the reconstruction error will be nonzero.
- Above we assumed that the origin of the coordinate system is in the sample mean, i.e. $\sum_i \mathbf{x}_i = 0$.

Principal Component Analysis (PCA)

In the general case, we want the reconstruction error $\|\mathbf{x} - \tilde{\mathbf{x}}\|$ to be small.

Objective: minimize $\min_{W \in \mathbb{R}^{p \times k}: W^t W = I} \sum_{i=1}^n \|\mathbf{x}_i - W W^t \mathbf{x}_i\|^2$

Finding the principal components

Projection vectors are orthogonal \rightsquigarrow can treat them separately:

$$\begin{aligned} \min_{\mathbf{w}: \|\mathbf{w}\|=1} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{w}\mathbf{w}^t \mathbf{x}_i\|^2 \\ \sum_i \|\mathbf{x}_i - \mathbf{w}\mathbf{w}^t \mathbf{x}_i\|^2 &= \sum_{i=1}^n [\mathbf{x}_i^t \mathbf{x}_i - 2\mathbf{x}_i^t \mathbf{w}\mathbf{w}^t \mathbf{x}_i + \underbrace{\mathbf{x}_i^t \mathbf{w}\mathbf{w}^t \mathbf{w}\mathbf{w}^t \mathbf{x}_i}_{=1}] \\ &= \sum_i [\mathbf{x}_i^t \mathbf{x}_i - \mathbf{x}_i^t \mathbf{w}\mathbf{w}^t \mathbf{x}_i] \\ &= \sum_i \mathbf{x}_i^t \mathbf{x}_i - \mathbf{w}^t \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^t \mathbf{w} \\ &= \underbrace{\sum_i \mathbf{x}_i^t \mathbf{x}_i}_{\text{const.}} - \mathbf{w}^t \mathbf{X}^t \mathbf{X} \mathbf{w}. \end{aligned}$$

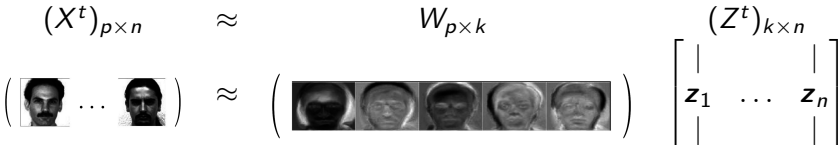
Finding the principal components

- Want to maximize $\mathbf{w}^t X^t X \mathbf{w}$ under the constraint $\|\mathbf{w}\| = 1$
- Can also maximize the ratio $J(\mathbf{w}) = \frac{\mathbf{w}^t X^t X \mathbf{w}}{\mathbf{w}^t \mathbf{w}}$ and rescale $\hat{\mathbf{w}}$.
- Optimal projection u is the eigenvector of $X^t X$ with largest eigenvalue (J is a Rayleigh quotient).
- Note that we assumed that $\sum_i \mathbf{x}_i = 0$. Thus, the columns of X are assumed to sum to zero.
 - ↪ compute SVD of “centered” matrix X
 - ↪ column vectors in W are eigenvectors of $X^t X$
 - ↪ they are the principal components.

Eigen-faces [Turk and Pentland, 1991]

- p = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^p$ is a face image
- x_{ji} = intensity of the j -th pixel in image i

$$(\mathbf{X}^t)_{p \times n} \approx \mathbf{W}_{p \times k} (\mathbf{Z}^t)_{k \times n}$$

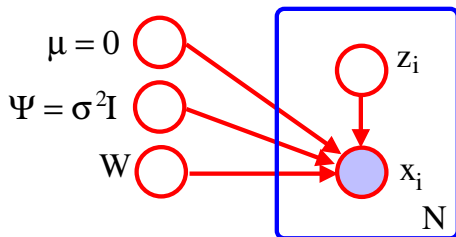


Idea: \mathbf{z}_i more 'meaningful' representation of i -th face than \mathbf{x}_i

Can use \mathbf{z}_i for nearest-neighbor classification

Much faster when $p \gg k$.

Probabilistic PCA



- Assuming $\Psi = \sigma^2 I$ and centered data in the FA model
 \rightsquigarrow likelihood and marginal likelihood

$$p(\mathbf{x}_i | z_i, \theta) = \mathcal{N}(W \mathbf{z}_i, \sigma^2 I),$$

$$\begin{aligned} p(\mathbf{x}_i | \theta) &= \int \mathcal{N}(\mathbf{x}_i, W \mathbf{z}_i, \sigma^2 I) \mathcal{N}(\mathbf{z}_i | \mathbf{0}, I) d\mathbf{z}_i, \\ &= \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \sigma^2 I + W W^t) \end{aligned}$$

Probabilistic PCA

- (Tipping & Bishop 1999): Maxima of the marg. likelihood given by

$$\hat{W} = V(\Lambda - \sigma^2 I)^{\frac{1}{2}} R,$$

where R is an arbitrary orthogonal matrix,
columns of V : first k eigenvectors of $S = \frac{1}{n} X^t X$,
 Λ : diagonal matrix of eigenvalues.

- As $\sigma^2 \rightarrow 0$, we have $\hat{W} \rightarrow V$, **as in classical PCA** (for $R = \Lambda^{-\frac{1}{2}}$).
- Projections \mathbf{z}_i : Posterior over the latent factors:

$$p(\mathbf{z}_i | \mathbf{x}_i, \hat{\boldsymbol{\theta}}) = \mathcal{N}(\mathbf{z}_i | \hat{\mathbf{m}}_i, \sigma^2 \hat{F}^{-1})$$

$$\hat{F} = \sigma^2 I + \hat{W}^t \hat{W}$$

$$\mathbf{m}_i = \hat{F}^{-1} \hat{W}^t \mathbf{x}_i$$

For $\sigma^2 \rightarrow 0$, $\mathbf{z}_i \rightarrow \mathbf{m}_i$ and $\mathbf{m}_i \rightarrow V^t \mathbf{x}_i \rightsquigarrow$ orthogonal projection of the data onto the column space of V , **as in classical PCA**.

Canonical Correlation Analysis [Hotelling, 1936]

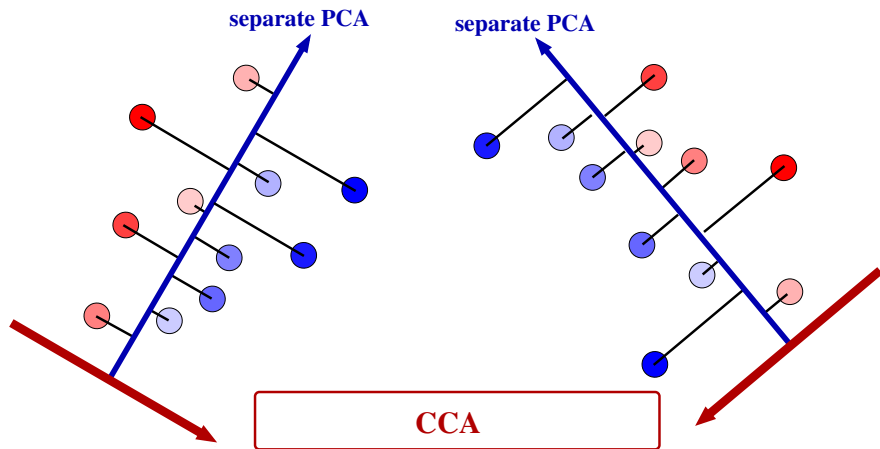
Often, each data point consists of **two views**:

- Image retrieval: for each image, have the following:
 - ▶ X : Pixels (or other visual features) Y : Text around the image
- Time series:
 - ▶ X : Signal at time t
 - ▶ Y : Signal at time $t + 1$
- Two-view learning: divide features into two sets
 - ▶ X : Features of a word/object, etc.
 - ▶ Y : Features of the context in which it appears

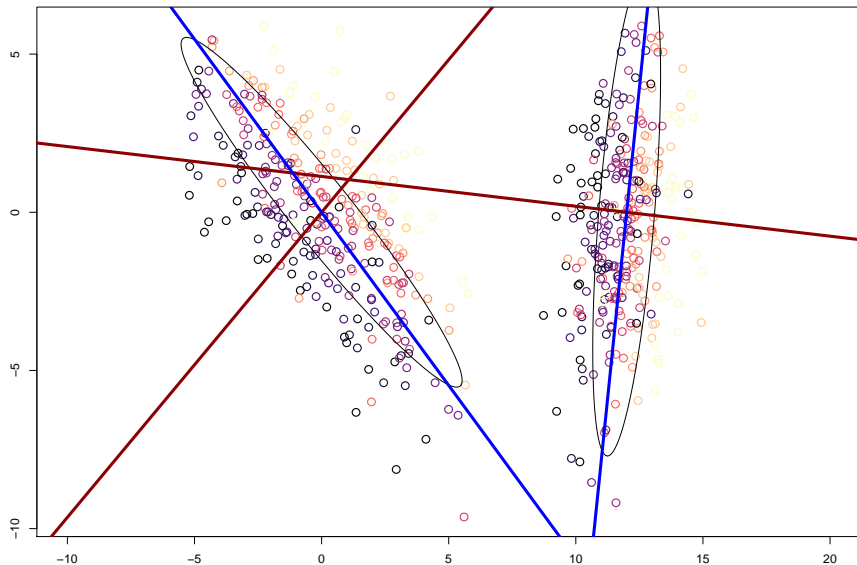
Goal: reduce the dimensionality of the two views **jointly**.

Find projections such that projected views are **maximally correlated**.

CCA vs PCA



CCA vs PCA



CCA: Setting

- Let \mathbf{X} be a random vector $\in \mathbb{R}^{p_x}$ and \mathbf{Y} be a random vector $\in \mathbb{R}^{p_y}$. Consider the combined ($p := p_x + p_y$)-dimensional random vector $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})^t$. Let its ($p \times p$) covariance matrix be partitioned into blocks according to:

$$\mathbf{Z} = \begin{bmatrix} \Sigma_{XX} \in \mathbb{R}^{p_x \times p_x} & | & \Sigma_{XY} \in \mathbb{R}^{p_x \times p_y} \\ \Sigma_{YX} \in \mathbb{R}^{p_y \times p_x} & | & \Sigma_{YY} \in \mathbb{R}^{p_y \times p_y} \end{bmatrix}$$

- Assuming centered data, the blocks in the covariance matrix can be estimated from observed data sets $X \in \mathbb{R}^{n \times p_x}$, $Y \in \mathbb{R}^{n \times p_y}$:

$$\mathbf{Z} \approx \frac{1}{n} \begin{bmatrix} X^t X & | & X^t Y \\ Y^t X & | & Y^t Y \end{bmatrix}$$

CCA: Setting

- $\text{Correlation}(x, y) = \frac{\text{covariance}(x, y)}{\text{standard deviation}(x) \cdot \text{standard deviation}(y)}$

$$\rho = \text{cor}(x, y) = \frac{\text{cov}(x, y)}{\sigma(x)\sigma(y)}.$$

- Sample correlation:

$$\rho = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})^t}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}} \quad \text{centered observations} \quad \frac{\mathbf{x}^t \mathbf{y}}{\sqrt{\mathbf{x}^t \mathbf{x}} \sqrt{\mathbf{y}^t \mathbf{y}}}.$$

- Want to find maximally correlated 1D-projections $\mathbf{x}^t \mathbf{a}$ and $\mathbf{y}^t \mathbf{b}$.

- Projected covariance: $\text{cov}(\mathbf{x}^t \mathbf{a}, \mathbf{y}^t \mathbf{b}) \stackrel{\text{zero means}}{=} \mathbf{a}^t \Sigma_{XY} \mathbf{b}$.

- Define $\mathbf{c} = \Sigma_{XX}^{-\frac{1}{2}} \mathbf{a}$, $\mathbf{d} = \Sigma_{YY}^{-\frac{1}{2}} \mathbf{b}$.

- Thus, the projected correlation coefficient is: $\rho = \frac{\mathbf{c}^t \Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}} \mathbf{d}}{\sqrt{\mathbf{c}^t \mathbf{c}} \sqrt{\mathbf{d}^t \mathbf{d}}}$.

CCA: Setting

- By the Cauchy-Schwarz inequality ($\mathbf{x}^t \mathbf{y} \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|$), we have

$$\left(\mathbf{c}^t \underbrace{\Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}}}_H \mathbf{d} \right) \leq \left(\mathbf{c}^t \underbrace{\Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma_{YX} \Sigma_{XX}^{-\frac{1}{2}} \mathbf{c}}_{G := HH^t} \right)^{\frac{1}{2}} (\mathbf{d}^t \mathbf{d})^{\frac{1}{2}},$$

$$\rho \leq \frac{(\mathbf{c}^t G \mathbf{c})^{\frac{1}{2}}}{(\mathbf{c}^t \mathbf{c})^{\frac{1}{2}}},$$

$$\rho^2 \leq \frac{\mathbf{c}^t G \mathbf{c}}{\mathbf{c}^t \mathbf{c}}.$$

- Equality: vectors \mathbf{d} and $\Sigma_{YY}^{-\frac{1}{2}} \Sigma_{YX} \Sigma_{XX}^{-\frac{1}{2}} \mathbf{c}$ are collinear.
- Maximum: \mathbf{c} is the eigenvector with the maximum eigenvalue of $G := \Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma_{YX} \Sigma_{XX}^{-\frac{1}{2}}$.
Subsequent pairs \rightsquigarrow using eigenvalues of decreasing magnitudes.
- Collinearity: $\mathbf{d} \propto \Sigma_{YY}^{-\frac{1}{2}} \Sigma_{YX} \Sigma_{XX}^{-\frac{1}{2}} \mathbf{c}$
- Transform back to original variables $\mathbf{a} = \Sigma_{XX}^{-\frac{1}{2}} \mathbf{c}$, $\mathbf{b} = \Sigma_{YY}^{-\frac{1}{2}} \mathbf{d}$.

Pixels That Sound [Kidron, Schechner, Elad, 2005]

“People and animals fuse auditory and visual information to obtain robust perception. A particular benefit of such cross-modal analysis is the ability to localize visual events associated with sound sources. We aim to achieve this using computer-vision aided by a single microphone”.



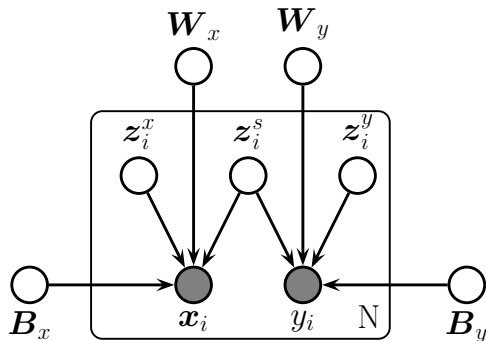
<https://webee.technion.ac.il/~yoav/research/pixels-that-sound.html>

Probabilistic CCA

(Bach and Jordan 2005): With Gaussian priors

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}^s | \mathbf{0}, I) \mathcal{N}(\mathbf{z}^x | \mathbf{0}, I) \mathcal{N}(\mathbf{z}^y | \mathbf{0}, I),$$

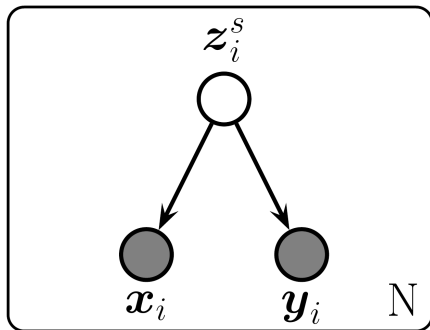
the MLE in the two-view FA model is equivalent to classical CCA (up to rotation and scaling).



From figure 12.19 in K. Murphy: Machine Learning. MIT Press 2012.

Probabilistic CCA: Essential Model Structure

Shared \mathbf{z}^S decorrelates the two views (\mathbf{x}, \mathbf{y}) :



$$\begin{aligned}\mathbf{x} &\in \mathbb{R}^p, \mathbf{y} \in \mathbb{R}^q \\ \mathbf{z}^S &\sim \mathcal{N}(\mathbf{z}^S | \mathbf{0}, I) \\ (\mathbf{x}, \mathbf{y}) | \mathbf{z} &\sim \mathcal{N}_{p+q}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}) \\ \boldsymbol{\Sigma} &= \begin{pmatrix} \boldsymbol{\Sigma}_x & 0 \\ 0 & \boldsymbol{\Sigma}_y \end{pmatrix}\end{aligned}$$

Further connections

- If y is a **discrete class label** \rightsquigarrow CCA is (essentially) equivalent to Linear Discriminant Analysis (LDA), see (Hastie et al. 1994).
- Arbitrary y \rightsquigarrow CCA is (essentially) equivalent to the **Gaussian Information Bottleneck** (Chechik et al. 2005)
 - ▶ Basic idea: **compress x** into compact latent representation while **preserving information about y** .
 - ▶ Information theoretic motivation:
Find encoding distribution $p(z|x)$ by minimizing
$$I(x; z) - \beta I(z; y)$$
where $\beta \geq 0$ is some parameter controlling the trade-off between compression and predictive accuracy.
- Arbitrary y , **discrete shared latent z^s**
 \rightsquigarrow **dependency-seeking clustering** (Klami and Kaski 2008): find clusters that “explain” the dependency between the two views.