Machine Learning

Volker Roth

Department of Mathematics & Computer Science University of Basel

Э

ヘロア 人間 アメヨア 人間 アー

Multiple Views: CCA

- Consider paired samples from different views.
- What is the dependency structure between the views ?
- Standard approach: global linear dependency detected by CCA.



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.

・ 同 ト ・ ヨ ト ・ ヨ ト

CCA vs PCA



E

・ロト ・四ト ・ヨト ・ヨト

CCA vs PCA



CCA: Setting

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

$$\boldsymbol{Z} = \begin{bmatrix} \boldsymbol{\Sigma}_{XX} \in \mathbb{R}^{p_X \times p_X} & | & \boldsymbol{\Sigma}_{XY} \in \mathbb{R}^{p_X \times p_y} \\ \boldsymbol{\Sigma}_{YX} \in \mathbb{R}^{p_y \times p_x} & | & \boldsymbol{\Sigma}_{YY} \in \mathbb{R}^{p_y \times p_y} \end{bmatrix}$$

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

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

CCA: Setting

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

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

Sample correlation:

$$\rho = \frac{\sum_{i} (x_{i} - \bar{x}) (y_{i} - \bar{y})^{t}}{\sqrt{\sum_{i} (x_{i} - \bar{x})^{2}} \sqrt{\sum_{i} (y_{i} - \bar{y})^{2}}} \stackrel{\text{centered observations}}{=} \frac{\boldsymbol{x}^{t} \boldsymbol{y}}{\sqrt{\boldsymbol{x}^{t} \boldsymbol{x}} \sqrt{\boldsymbol{y}^{t} \boldsymbol{y}}}.$$

- Want to find maximally correlated 1D-projections $x^t a$ and $y^t b$.
- Projected covariance: $cov(\mathbf{x}^t \mathbf{a}, \mathbf{y}^t \mathbf{b}) \stackrel{\text{zero means}}{=} \mathbf{a}^t \Sigma_{XY} \mathbf{b}$.

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

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

・ロト ・御 ト ・ ヨト ・ ヨト ・ ヨ

CCA: Setting

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

$$\begin{pmatrix} \boldsymbol{c}^{t} \underbrace{\boldsymbol{\Sigma}_{\boldsymbol{X}\boldsymbol{X}}^{-\frac{1}{2}} \boldsymbol{\Sigma}_{\boldsymbol{X}\boldsymbol{Y}} \boldsymbol{\Sigma}_{\boldsymbol{Y}\boldsymbol{Y}}^{-\frac{1}{2}}}_{\boldsymbol{H}} \end{pmatrix} \boldsymbol{d} \leq \begin{pmatrix} \boldsymbol{c}^{t} \underbrace{\boldsymbol{\Sigma}_{\boldsymbol{X}\boldsymbol{X}}^{-\frac{1}{2}} \boldsymbol{\Sigma}_{\boldsymbol{X}\boldsymbol{Y}} \boldsymbol{\Sigma}_{\boldsymbol{Y}\boldsymbol{Y}}^{-\frac{1}{2}} \boldsymbol{\Sigma}_{\boldsymbol{Y}\boldsymbol{X}} \boldsymbol{\Sigma}_{\boldsymbol{X}\boldsymbol{X}}^{-\frac{1}{2}}}_{\boldsymbol{G}:=\boldsymbol{H}\boldsymbol{H}^{t}} \boldsymbol{c} \end{pmatrix}^{\frac{1}{2}} (\boldsymbol{d}^{t}\boldsymbol{d})^{\frac{1}{2}}, \\ \rho \leq \frac{(\boldsymbol{c}^{t}\boldsymbol{G}\boldsymbol{c})^{\frac{1}{2}}}{(\boldsymbol{c}^{t}\boldsymbol{c})^{\frac{1}{2}}}, \\ \rho^{2} \leq \frac{\boldsymbol{c}^{t}\boldsymbol{G}\boldsymbol{c}}{\boldsymbol{c}^{t}\boldsymbol{c}}. \end{cases}$$

- Equality: vectors **d** and $\sum_{YY}^{-\frac{1}{2}} \sum_{YX} \sum_{XX}^{-\frac{1}{2}} c$ are collinear.
- Maximum: c is the eigenvector with the maximum eigenvalue of $G := \sum_{XX}^{-\frac{1}{2}} \sum_{XY} \sum_{YY}^{-1} \sum_{YX} \sum_{XX}^{-\frac{1}{2}}$. Subsequent pairs \rightsquigarrow using eigenvalues of decreasing magnitudes.

• Collinearity:
$$\boldsymbol{d} \propto \Sigma_{YY}^{-\frac{1}{2}} \Sigma_{YX} \Sigma_{XX}^{-\frac{1}{2}} \boldsymbol{c}$$

• Transform back to original variables $\boldsymbol{a} = \sum_{XX}^{-\frac{1}{2}} \boldsymbol{c}, \ \boldsymbol{b} = \sum_{YY}^{-\frac{1}{2}} \boldsymbol{d}.$

Efficient computation of separate PCAs and joint CCA

- Separate SVD of X and Y: $X = U_x S_x V_x^t$, $Y = U_y S_y V_y^t$ \rightsquigarrow separate PCAs for both views.
- Consider sample covariance matrix

 $\Sigma_{XX} = 1/n \cdot X^t X = 1/n \cdot V_x S_x U_x^t U_x S_x V_x^t = 1/n \cdot V_x S_x^2 V_x^t.$

• Substitute in matrix G:

$$G = \Sigma_{XX}^{-1/2} \cdot \Sigma_{XY} \cdot \Sigma_{YY}^{-1} \cdot \Sigma_{YX} \cdot \Sigma_{XX}^{-1/2} = (V_x S_x^{-1} V_x^t) (V_x S_x U_x^t U_y S_y V_y^t) (V_y S_y^{-2} V_y^t) (V_y S_y U_y^t U_x S_x V_x^t) (V_x S_x^{-1} V_x^t) = V_x U_x^t U_y U_y^t U_x V_x^t.$$

イロト イポト イヨト イヨト 二日

Efficient computation of separate PCAs and joint CCA

• **c** is an eigenvector of **G**

 $\Rightarrow \tilde{\boldsymbol{c}} = V_x^t \boldsymbol{c}$ is an eigenvector of $U_x^t U_y U_y^t U_x$ with the same e.vals:

$$G\boldsymbol{c} = \lambda \boldsymbol{c}$$

$$V_{x} U_{x}^{t} U_{y} U_{y}^{t} U_{x} V_{x}^{t} \boldsymbol{c} = \lambda \boldsymbol{c},$$

$$U_{x}^{t} U_{y} U_{y}^{t} U_{x} V_{x}^{t} \boldsymbol{c} = \lambda V_{x}^{t} \boldsymbol{c},$$

$$U_{x}^{t} U_{y} U_{y}^{t} U_{x} \boldsymbol{\lambda}_{x}^{t} \boldsymbol{c} = \lambda \tilde{\boldsymbol{c}}.$$

• In matrix form (the columns of matrix C are the eigenvectors c_j)

 $GC = C\Lambda,$ $U_x^t U_y U_y^t U_x (V_x^t C) = (V_x^t C)\Lambda,$ $U_x^t U_y U_y^t U_x \tilde{C} = \tilde{C}\Lambda.$

• Columns of \tilde{C} are the columns of \tilde{U} in SVD of $U_x^t U_y = \tilde{U} \tilde{S} \tilde{V}^t$

・ロト ・四ト ・ヨト ・ ヨト

Efficient computation of separate PCAs and joint CCA

- Now recover C from \tilde{C} : $C = V_x \tilde{C} = V_x \tilde{U}$ $\Rightarrow A = \sum_{XX}^{-1/2} C = (\frac{1}{n} V_x S_x^{-1} V_x^t) (V_x \tilde{U}) = \frac{1}{n} V_x S_x^{-1} \tilde{U}.$
- $B = \frac{1}{n} V_y S_y^{-1} \tilde{V}$ analogous.
- Joint CCA: correlation is maximized between pairs of projections Xa_i and Yb_i, with a_i, b_i are the *i*-th columns of A, B. The corresponding correlations are ρ_i = S̃_{ii}.
- In summary:
 - Compute individual SVDs of X and $Y \rightsquigarrow V_x, S_X, V_y, S_y$
 - Compute "joint" SVD of $U_x^t U_y = \tilde{U}\tilde{S}\tilde{V}^t$.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

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

・ロト ・ 一 マ ・ コ ト ・ 日 ト

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

< ロ > < 同 > < 三 > < 三 >

Further connections

- If y is a discrete class label → CCA is (essentially) equivalent to Linear Discriminant Analysis (LDA), see (Hastie et al. 1994).
- Arbitrary y ~>> 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

 $l(\mathbf{x}; \mathbf{z}) - \beta l(\mathbf{z}; \mathbf{y})$

where $\beta \ge 0$ is some parameter controlling the trade-off between compression and predictive accuracy.

Arbitrary y, discrete shared latent z^s
 → dependency-seeking clustering (Klami and Kaski 2008): find clusters that "explain" the dependency between the two views.

イロト イポト イヨト イヨト 三日

The Information Bottleneck (Tishby et al., 1999)

FA is powerful, but still limited (Gaussian assumptions etc.). Alternatives?



Mutual Information

- A measure of **mutual dependence** between two random variables: reduction of uncertainty by knowing one variable.
- For continuous RVs:

$$\begin{split} l(\mathbf{x}; \mathbf{y}) &= \int \int p(\mathbf{x}, \mathbf{y}) \log \left(\frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x}) p(\mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \\ &= D_{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x}) p(\mathbf{y})) \\ &= \int p(\mathbf{x}) \int p(\mathbf{y} | \mathbf{x}) \log \left(\frac{p(\mathbf{y} | \mathbf{x})}{p(\mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \\ &= E_{p(\mathbf{x})} D_{KL}(p(\mathbf{y} | \mathbf{x}) \| p(\mathbf{y})). \end{split}$$

• **x** and **y** independent \rightsquigarrow knowing **x** does not give any information about $\mathbf{y} \rightsquigarrow l(\mathbf{x}; \mathbf{y}) = 0$.

イロト イボト イヨト 一日

Information Bottleneck

- The IB principle: compress x into z, keep information about y.
- Assume y and z are conditionally independent given x an solve:

$$\min_{p(\boldsymbol{z}|\boldsymbol{x})} I(\boldsymbol{x}; \boldsymbol{z}) - \beta I(\boldsymbol{z}; \boldsymbol{y}).$$

The original IB formulation is not a generative model, nor do we actually condition: x, y are only used for estimating p(x, y).



IB as a latent variable model

Assume $\mathbf{z} = f(\mathbf{x}) + \boldsymbol{\xi}$ captures all relevant information about \mathbf{y} $\Rightarrow \mathbf{x} \perp \mathbf{y} | \mathbf{z}$

 \sim latent version IB ^(lat), basically an asymmetric CCA model.

CCA:
$$p(x|z)p(y|z)p(z)$$

 $x \perp y \mid z$

$$\begin{array}{c} \mathsf{IB}^{(\mathsf{lat})} \colon p(z|x)p(y|z)p(x) \\ x \perp y \mid z \end{array}$$





- 4 同 ト 4 ヨ ト

Gaussian IB (Chechnik et al. 2003)

• Assume x and y are jointly Gaussian-distributed.

$$(\mathbf{x}, \mathbf{y}) \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \Sigma_x & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_y \end{pmatrix}\right),$$

• The optimal z is a noisy projection of x:

 $\boldsymbol{z} = A\boldsymbol{x} + \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}) \Rightarrow \boldsymbol{z} | \boldsymbol{x} \sim \mathcal{N}(A\boldsymbol{x}, \boldsymbol{I}), \ \boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, A\Sigma_{x}A^{t} + \boldsymbol{I}).$

- Analytic form of mutual information: $l(\mathbf{x}; \mathbf{z}) = \frac{1}{2} \log |A\Sigma_x A^t + I|,$ $l(\mathbf{z}; \mathbf{y}) = l(\mathbf{x}; \mathbf{z}) - \frac{1}{2} \log |A\Sigma_{x|y} A^t + I|.$
- general y: rows of A are eigenvectors of $\Sigma_x^{-1}\Sigma_{x|y} \rightsquigarrow \mathbb{CCA}$
- one-dimensional y: ~> least squares regression
- y is noisy version of x: \rightsquigarrow PCA

Sparse Gaussian IB

- Recall: $\mathbf{z} = A\mathbf{x} + \boldsymbol{\xi}$, $\Rightarrow \mathbf{z} | \mathbf{x} \sim \mathcal{N}(A\mathbf{x}, I), \mathbf{z} \sim \mathcal{N}(\mathbf{0}, A\Sigma_{\mathbf{x}}A^t + I)$.
- Note that $|A\Sigma A^t + I| = |\Sigma A^t A + I| \rightsquigarrow$ only $A^t A$ identifiable.
- **Sparse version:** assume that $A^t A$ is diagonal.
- Intuition: for RVs x, x', any full-rank projection Ax' of x' would lead to the same mutual information since I(x, x') = I(x; Ax'), and a reduction can only be achieved by a rank-deficient matrix A.
- If we compress X, A cannot have full rank $\rightsquigarrow A^t A$ diagonal implies that A can be non-zero only in some subspace, $AP = \begin{pmatrix} A' & 0 \end{pmatrix}$.
- The mean $\mu = Ax$ is a sparse compression of x.
- Basically an information-theoretic lasso variant.

イロト イボト イヨト 一日

Sparse Gaussian IB: Experiments

Prognosis of cutaneous **malignant melanoma** (MM) using biomarkers (Fuchs & Buhmann, 2011; Rey & R. 2014).

- **Goal: identify biomarkers** relevant to the disease evolution.
- Crucial for cost-effective prognosis and therapy optimization

Data:

- Variable to compress, *x*: **immunohistochemical (IHC) expressions** of 70 candidate biomarkers measured for 364 patients.
- Relevance variable, **y**: 9 different clinical observations about the stage of the tumor & survival times.

Sparse Gaussian IB: Experiments



Multiple Views: Dependency Seeking Clustering

- Consider paired samples from different views.
- What is the dependency structure between the views ?
- Standard approach: global linear dependence detected by CCA.
- $\bullet~$ Generalization \rightarrow dependency-seeking clustering.
- The cluster structure captures dependencies.



Dependency Seeking Clustering



Volker Roth (University of Basel)

24 / 25

Э

< ∃ →

Dependency Seeking Clustering



Volker Roth (University of Basel)

25 / 25