





Advanced Multimodal Machine Learning

Lecture 7.1: Multivariate Statistics

Louis-Philippe Morency Tadas Baltrusaitis

Lecture Objectives

- Quick recap
- Multivariate statistical analysis
 - Basic concepts (multivariate, covariance,...)
 - Principal component analysis (+SVD)
- Canonical Correlation Analysis
- Deep Correlation Networks
 - Deep CCA, DCCA-AutoEncoder
 - (Deep) Correlational neural networks
- Matrix Factorization
 - Nonnegative Matrix Factorization

Administrative Stuff

Lecture Schedule

Classes	Lectures
Week 5	Multimodal representation learning
2/14 & 2/16	 Multimodal auto-encoders
	 Multimodal deep neural networks
Week 6 2/21 & 2/23	First project assignment - Presentations
Week 7	Multimodal component analysis
2/28 & 3/2	 Deep canonical correlation analysis
	 Non-negative matrix factorization
Week 8	Multimodal Optimization
3/7 & 3/9	 Optimization in deep neural networks
	 Variational approaches

Lecture Schedule

Classes	Lectures
Spring break 3/13 – 3/17	
Week 9	Multimodal alignment
3/21 & 3/23	 Attention models and multi-instance learning
	 Multimodal synchrony and prediction
Week 10	Markov Random Fields
3/28 & 3/30	 Boltzmann distribution and CRFs
	 Continuous and fully-connected CRFs
Week 11 4/4 & 4/6	Mid-term project assignment - Presentations

Lecture Schedule

Classes	Lectures
Week 12	Multimodal fusion
4/11 & 4/13	 Sample-based late fusion
	 Multi-kernel learning and fusion
Week 13	Application: Multilingual computational models
4/18 & 4/20	 Neural Machine Translation
	 Sub-word Models
Week 14	Application: Language and Vision
4/25 & 4/27	 Learned visual representations
	 Visual activity recognition
Week 15 5/2 & 5/4	Final project assignment - Presentations

Upcoming Schedule

- Pre-proposal (tomorrow Wednesday 2/5 at 9am)
- First project assignment:
 - Proposal presentation (2/21 and 2/23)
 - First project report (Sunday 3/5)
- Second project assignment
 - Midterm presentations (4/4 and 4/6)
 - Midterm report (Sunday 4/9)
- Final project assignment
 - Final presentation (5/2 & 5/4)
 - Final report (Sunday 5/7)

Project Proposal Report – Due on 3/5/17

- Part 1 (updated version of your pre-proposal)
 - Research problem:
 - Describe and motivate the research problem
 - Define in generic terms the main computational challenges
 - Dataset and Input Modalities:
 - Describe the dataset(s) you are planning to use for this project.
 - Describe the input modalities and annotations available in this dataset.

Project Proposal Report – Due on 3/5/17

Part 2

Related Work:

- Include 12-15 paper citations which give an overview of the prior work
- Present in more details the 3-4 research papers most related to your work

Research Challenges and Hypotheses:

- Describe your specific challenges and/or research hypotheses
- Highlight the novel aspect of your proposed research

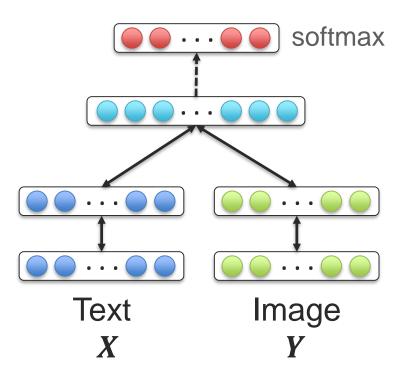
Project Proposal Report – Due on 3/5/17

- Part 3 (teams of 2 members can pick either one)
 - Language Modality Exploration:
 - Explore neural language models on your dataset (using Keras/Theano)
 - Train at least two different language models (e.g., using SimpleRNN, GRU or LSTM) on your dataset and compare their perplexity.
 - Include qualitative examples of successes and failure cases.
 - Visual Modality Exploration:
 - Explore pre-trained Convolutional Neural Networks (CNNs) on your dataset
 - Load a pre-existing CNN model trained for object recognition (e.g., AlexNet or VGG-Net) and process your test images.
 - Extract features at different network layers in the network and visualize them (using t-sne visualization) with overlaid class labels with different colors.

Quick Recap

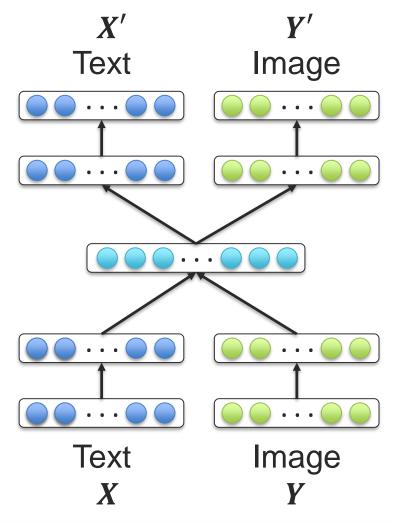
Learn (unsupervised) a joint representation between multiple modalities where similar unimodal concepts are closely projected.

Deep MultimodalBoltzmann machines



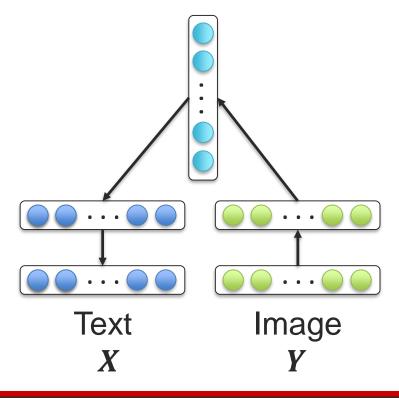
Learn (unsupervised) a joint representation between multiple modalities where similar unimodal concepts are closely projected.

- Deep MultimodalBoltzmann machines
- Stacked Autoencoder



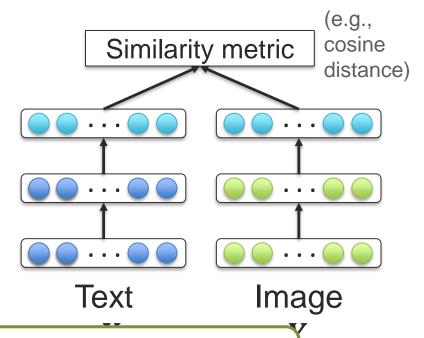
Learn (unsupervised) a joint representation between multiple modalities where similar unimodal concepts are closely projected.

- Deep MultimodalBoltzmann machines
- Stacked Autoencoder
- Encoder-Decoder



Learn (unsupervised) a joint representation between multiple modalities where similar unimodal concepts are closely projected.

- Deep MultimodalBoltzmann machines
- Stacked Autoencoder
- Encoder-Decoder
- "Minimum-distance"Multimodal Embedding



How Can We Learn Better Representations?



Multivariate Statistical Analysis

Multivariate Statistical Analysis

"Statistical approaches to understand the relationships in high dimensional data"

- Example of multivariate analysis approaches:
 - Multivariate analysis of variance (MANOVA)
 - Principal components analysis (PCA)
 - Factor analysis
 - Linear discriminant analysis (LDA)
 - Canonical correlation analysis (CCA)

Random Variables

Definition: A variable whose possible values are numerical outcomes of a random phenomenon.

- □ **Discrete** random variable is one which may take on only a countable number of distinct values such as 0,1,2,3,4,...
- ☐ Continuous random variable is one which takes an infinite number of possible values.

Examples of random variables:

- Someone's age
- Someone's height
- Someone's weight

Discrete or continuous?

Correlated?

Definitions

Given two random variables X and Y:

Expected value probability-weighted average of all possible values

$$\mu = E[X] = \sum_{i} x_i P(x_i)$$

 \triangleright If same probability for all observations x_i , then same as arithmetic mean

Variance measures the spread of the observations

$$\sigma^2 = Var(X) = E[(X - \mu)(X - \mu)] = E[\bar{X}\bar{X}]$$
 If data is centered

 \blacktriangleright Variance is equal to the square of the standard deviation σ

Covariance measures how much two random variables change together

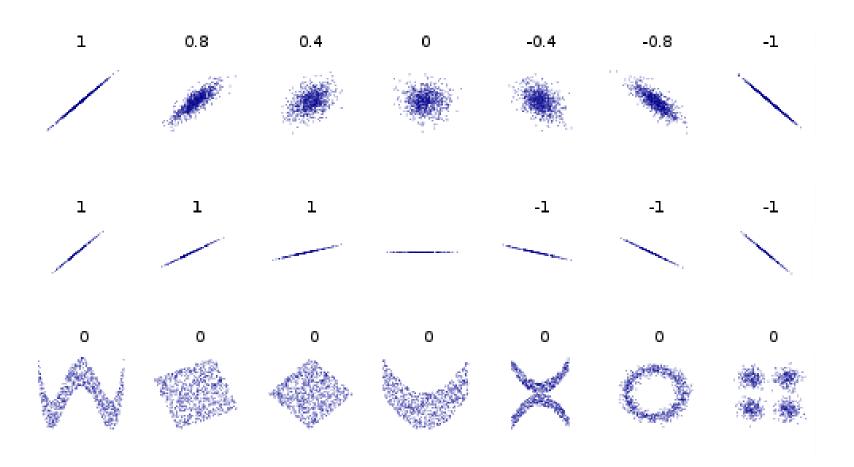
$$cov(X,Y) = E[(X - \mu_X)(Y - \mu_Y)] = E[\overline{X}\overline{Y}]$$

Definitions

Pearson Correlation measures the extent to which two variables have a linear relationship with each other

$$\rho_{X,Y} = corr(X,Y) = \frac{cov(X,Y)}{var(X)var(Y)}$$

Pearson Correlation Examples



Definitions

Multivariate (multidimensional) random variables

(aka random vector)

$$X = [X^1, X^2, X^3, ..., X^M]$$

 $Y = [Y^1, Y^2, Y^3, ..., Y^N]$

Covariance matrix generalizes the notion of variance

$$\Sigma_X = \Sigma_{X,X} = var(X) = E[(X - E[X])(X - E[X])^T] = E[\overline{X}\overline{X}^T]$$

Cross-covariance matrix generalizes the notion of covariance

$$\Sigma_{X,Y} = cov(X,Y) = E[(X - E[X])(Y - E[Y])^T] = E[\overline{X}\overline{Y}^T]$$

Definitions

Multivariate (multidimensional) random variables

(aka random vector)

$$X = [X^1, X^2, X^3, ..., X^M]$$

 $Y = [Y^1, Y^2, Y^3, ..., Y^N]$

Covariance matrix generalizes the notion of variance

$$\Sigma_X = \Sigma_{X,X} = var(X) = E[(X - E[X])(X - E[X])^T] = E[\overline{X}\overline{X}^T]$$

Cross-covariance matrix generalizes the notion of covariance

$$\Sigma_{\boldsymbol{X},\boldsymbol{Y}} = cov(\boldsymbol{X},\boldsymbol{Y}) = \begin{bmatrix} cov(X_1,Y_1) & cov(X_2,Y_1) & \cdots & cov(X_M,Y_1) \\ cov(X_1,Y_2) & cov(X_2,Y_2) & \cdots & cov(X_M,Y_2) \\ \vdots & & \vdots & \ddots & \vdots \\ cov(X_1,Y_N) & cov(X_2,Y_N) & \cdots & cov(X_M,Y_N) \end{bmatrix}$$

Definitions – Matrix Operations

Trace is defined as the sum of the elements on the main diagonal of any matrix *X*

$$tr(X) = \sum_{i=1}^{n} x_{ii}$$

Eigenvalues and Eigenvectors

Eigenvalue decomposition

If A is an $n \times n$ matrix, do there exist nonzero vectors \mathbf{x} in \mathbb{R}^n such that $A\mathbf{x}$ is a scalar multiple of \mathbf{x} ?

(The term eigenvalue is from the German word Eigenwert, meaning "proper value")

Eigenvalue equation:

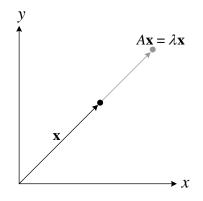
 $A\mathbf{x} = \lambda \mathbf{x}$ Eigenvector Eigenvalue

A: an $n \times n$ matrix

 λ : a scalar (could be **zero**)

x: a **nonzero** vector in R^n

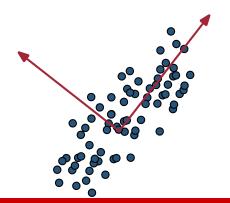
Geometric Interpretation

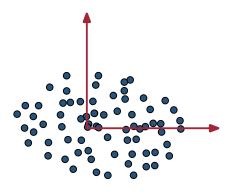


Principal component analysis

PCA converts a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called *principal components*

- Eigenvectors are orthogonal towards each other and have length one
- The first couple of eigenvectors explain the most of the variance observed in the data
- Low eigenvalues indicate little loss of information if omitted





Singular Value Decomposition (SVD)

SVD expresses any matrix A as

$$A = USV^T$$

• The columns of \mathbf{U} are eigenvectors of $\mathbf{A}\mathbf{A}^T$, and the columns of \mathbf{V} are eigenvectors of $\mathbf{A}^T\mathbf{A}$.

$$\mathbf{A}\mathbf{A}^T\mathbf{u}_i = s_i^2\mathbf{u}_i$$
$$\mathbf{A}^T\mathbf{A}\mathbf{v}_i = s_i^2\mathbf{v}_i$$

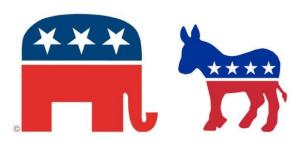
Multi-view Learning

X

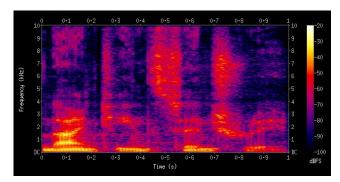


demographic properties

Y



responses to survey



audio features at time i



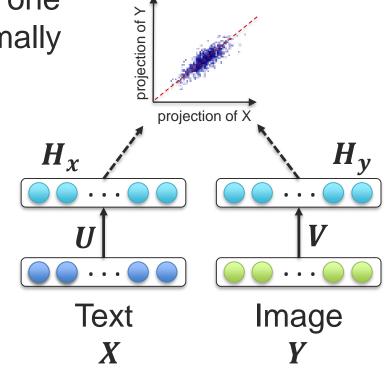
video features at time i

"canonical": reduced to the simplest or clearest schema possible

1 Learn two linear projections, one for each view, that are maximally correlated:

$$(u^*, v^*) = \underset{u,v}{\operatorname{argmax}} corr(H_x, H_y)$$

= $\underset{u,v}{\operatorname{argmax}} corr(u^T X, v^T Y)$



Correlated Projection

1 Learn two linear projections, one for each view, that are maximally correlated:

$$(\boldsymbol{u}^*, \boldsymbol{v}^*) = \underset{\boldsymbol{u}, \boldsymbol{v}}{\operatorname{argmax}} corr(\boldsymbol{u}^T \boldsymbol{X}, \boldsymbol{v}^T \boldsymbol{Y})$$



Two views *X*, *Y* where same instances have the same color

1 Learn two linear projections, one for each view, that are maximally correlated:

$$(u^*, v^*) = \underset{u,v}{\operatorname{argmax}} \operatorname{corr}(u^T X, v^T Y)$$

$$= \underset{u,v}{\operatorname{argmax}} \frac{\operatorname{cov}(u^T X, v^T Y)}{\operatorname{var}(u^T X) \operatorname{var}(v^T Y)}$$

$$= \underset{u,v}{\operatorname{argmax}} \frac{u^T X Y^T v}{\sqrt{u^T X X^T u} \sqrt{v^T Y Y^T v}}$$

$$= \underset{u,v}{\operatorname{argmax}} \frac{u^T \Sigma_{XY} v}{\sqrt{u^T \Sigma_{XY} u} \sqrt{v^T \Sigma_{YY} v}}$$

where

$$\Sigma_{XY} = cov(X, Y) = XY^T$$

if both X, Y have 0 mean

$$\mu_X = 0$$
 $\mu_Y = 0$

We want to learn multiple projection pairs $(u_{(i)}X, v_{(i)}Y)$:

$$(\boldsymbol{u}_{(i)}^*, \boldsymbol{v}_{(i)}^*) = \underset{\boldsymbol{u}_{(i)}, \boldsymbol{v}_{(i)}}{\operatorname{argmax}} \frac{\boldsymbol{u}_{(i)}^T \boldsymbol{\Sigma}_{XY} \boldsymbol{v}_{(i)}}{\sqrt{\boldsymbol{u}_{(i)}^T \boldsymbol{\Sigma}_{XX} \boldsymbol{u}_{(i)}} \sqrt{\boldsymbol{v}_{(i)}^T \boldsymbol{\Sigma}_{YY} \boldsymbol{v}_{(i)}} }$$

We want these multiple projection pairs to be orthogonal ("canonical") to each other:

$$u_{(i)}^T \Sigma_{XY} v_{(j)} = u_{(j)}^T \Sigma_{XY} v_{(i)} = 0$$
 for $i \neq j$

$$m{U}m{\Sigma}_{m{X}m{Y}}m{V}=tr(m{U}m{\Sigma}_{m{X}m{Y}}m{V}) \qquad ext{where } m{U}=[m{u}_{(1)},m{u}_{(2)},...,m{u}_{(k)}] \ ext{and } m{V}=[m{v}_{(1)},m{v}_{(2)},...,m{v}_{(k)}]$$

$$(\boldsymbol{U}^*, \boldsymbol{V}^*) = \underset{\boldsymbol{U}, \boldsymbol{V}}{\operatorname{argmax}} \frac{tr(\boldsymbol{U}^T \boldsymbol{\Sigma}_{\boldsymbol{X} \boldsymbol{Y}} \boldsymbol{V})}{\sqrt{\boldsymbol{U}^T \boldsymbol{\Sigma}_{\boldsymbol{X} \boldsymbol{X}} \boldsymbol{U}} \sqrt{\boldsymbol{V}^T \boldsymbol{\Sigma}_{\boldsymbol{Y} \boldsymbol{Y}} \boldsymbol{V}}}$$

3 Since this objective function is invariant to scaling, we can constraint the projections to have unit variance:

$$U^T \Sigma_{XX} U = I \qquad V^T \Sigma_{YY} V = I$$

Canonical Correlation Analysis:

maximize: $tr(U^T\Sigma_{XY}V)$

subject to: $U^T \Sigma_{YY} U = V^T \Sigma_{YY} V = I$

maximize: $tr(U^T \Sigma_{XY} V)$

subject to: $U^T \Sigma_{YY} U = V^T \Sigma_{YY} V = I$

$$\Sigma = \begin{bmatrix} \Sigma_{XX} & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_{YY} \end{bmatrix} \stackrel{u,v}{\Longrightarrow} \begin{bmatrix} 1 & 0 & 0 & \lambda_1 & 0 & 0 \\ 0 & 1 & 0 & 0 & \lambda_2 & 0 \\ 0 & 0 & 1 & 0 & 0 & \lambda_3 \\ \lambda_1 & 0 & 0 & 1 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 & 1 & 0 \\ 0 & 0 & \lambda_3 & 0 & 0 & 1 \end{bmatrix}$$

maximize:
$$tr(U^T\Sigma_{XY}V)$$
 subject to: $U^T\Sigma_{YY}U = V^T\Sigma_{YY}V = I$ How to solve it? Lagrange Multipliers! Lagrange function $L = tr(U^T\Sigma_{XY}V) + \alpha(U^T\Sigma_{YY}U - I) + \beta(V^T\Sigma_{YY}V - I)$ \Rightarrow And then find stationary points of L : $\frac{\partial L}{\partial U} = 0$ $\frac{\partial L}{\partial V} = 0$ $\Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{XY}^{T}U = \lambda U$

 $\Sigma_{YY}^{-1}\Sigma_{XY}^T\Sigma_{XX}^{-1}\Sigma_{XY}V = \lambda V$ where $\lambda = 4\alpha\beta$

Canonical Correlation Analysis

maximize: $tr(U^T \Sigma_{XY} V)$

subject to: $U^T \Sigma_{VV} U = V^T \Sigma_{VV} V = I$

$$T \triangleq \mathbf{\Sigma}_{XX}^{-1/2} \mathbf{\Sigma}_{XY} \mathbf{\Sigma}_{YY}^{-1/2}$$

$$(U^*, V^*) = (\Sigma_{XX}^{-1/2} U_{SVD}, \Sigma_{YY}^{-1/2} V_{SVD})$$

Can solve these eigenvalue equations with Singular Value Decomposition (SVD)

Eigenvalues

Eigenvectors

$$\Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{XY}^TU = \lambda U$$

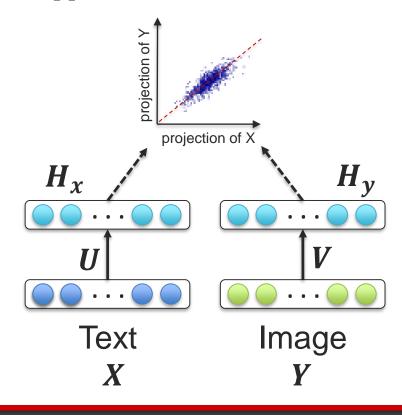
Eigenvalue
$$\begin{cases} \Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{XY}^TU = \lambda U \\ \text{equations} \end{cases} \quad \Sigma_{YY}^{-1}\Sigma_{XY}^T\Sigma_{XX}^{-1}\Sigma_{XY}V = \lambda V \quad \text{where } \lambda = 4\alpha\beta \end{cases}$$

Canonical Correlation Analysis

maximize: $tr(U^T \Sigma_{XY} V)$

subject to: $U^T \Sigma_{YY} U = V^T \Sigma_{YY} V = I$

- Linear projections maximizing correlation
- Orthogonal projections
- Unit variance of the projection vectors



Exploring Deep Correlation Networks

Deep Canonical Correlation Analysis

Same objective function as CCA:

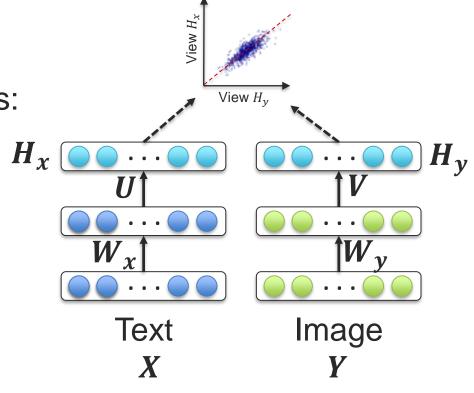
$$\underset{v,u,w_x,w_y}{\operatorname{argmax}} \ \operatorname{corr} \big(\boldsymbol{H}_x, \boldsymbol{H}_y \big)$$

But need to compute gradients:

$$\frac{\partial corr(\boldsymbol{H}_x, \boldsymbol{H}_y)}{\partial U}$$

$$\frac{\partial corr(\boldsymbol{H}_{x}, \boldsymbol{H}_{y})}{\partial V}$$

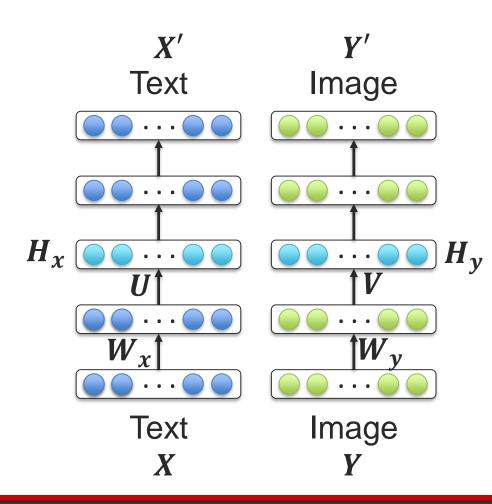
Andrew et al., ICML 2013



Deep Canonical Correlation Analysis

Training procedure:

 Pre-train the models parameters using denoising autoencoders

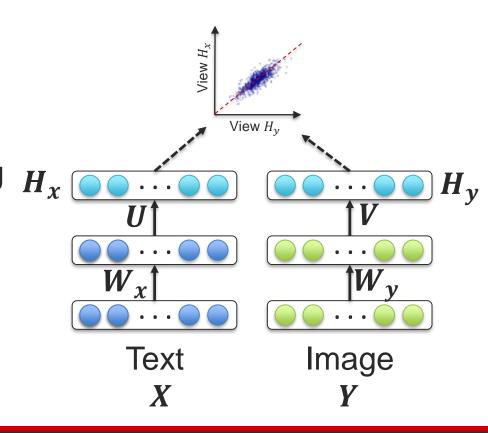


Andrew et al., ICML 2013

Deep Canonical Correlation Analysis

Training procedure:

- Pre-train the models parameters using denoising autoencoders
- Optimize the CCA
 objective functions using
 large mini-batches or
 full-batch (L-BFGS)

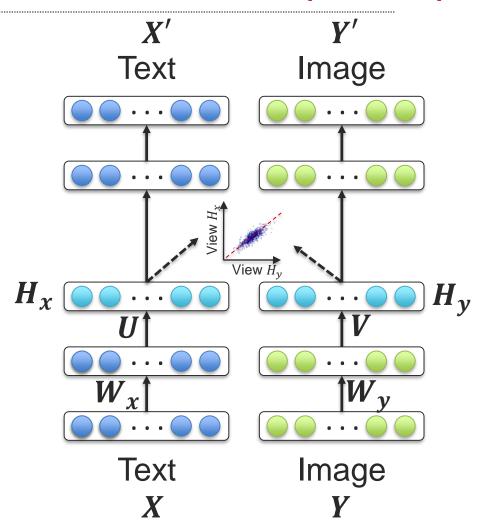


Andrew et al., ICML 2013

Deep Canonically Correlated Autoencoders (DCCAE)

Jointly optimize for DCCA and autoencoders loss functions

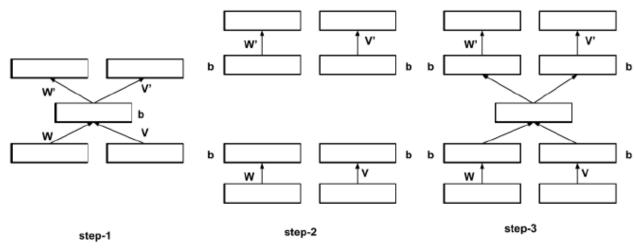
A trade-off between multi-view correlation and reconstruction error from individual views



Wang et al., ICML 2015

Deep Correlational Neural Network

- Learn a shallow CCA autoencoder (similar to 1 layer DCCAE model)
- 2. Use the learned weights for initializing the autoencoder layer
- 3. Repeat procedure



Chandar et al., Neural Computation, 2015

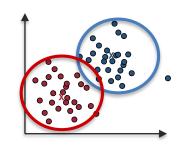
Matrix Factorization

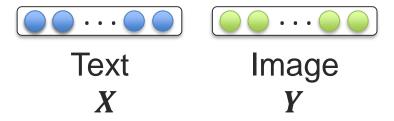
Data Clustering

How to discover groups in your data?

K-mean is a simple clustering algorithm based on competitive learning

- Iterative approach
 - Assign each data point to one cluster (based on distance metric)
 - Update cluster centers
 - Until convergence
- "Winner takes all"





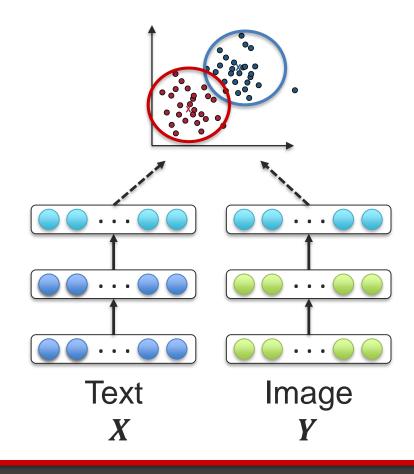


Enforcing Data Clustering in Deep Networks

How to enforce data clustering in our

(multimodal) deep learning

algorithms?



Nonnegative Matrix Factorization (NMF)

Given: Nonnegative n x m matrix M (all entries ≥ 0)

$$\begin{pmatrix} & & \\ & & \\ & & \\ & & \end{pmatrix} = \begin{pmatrix} & & \\ & & \\ & & \\ & & \end{pmatrix}$$

Want: Nonnegative matrices F (n x r) and G (r x m), s.t. X = FG.

- easier to interpret
- provide better results in information retrieval, clustering

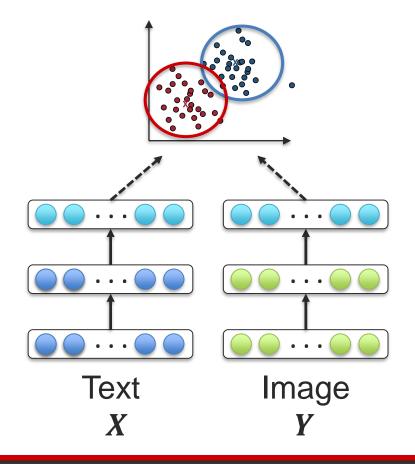
Semi-NMF and Other Extensions

SVD: $X_{\pm} \approx F_{\pm} G_{\pm}^T$

NMF: $X_+ \approx F_+ G_+^T$

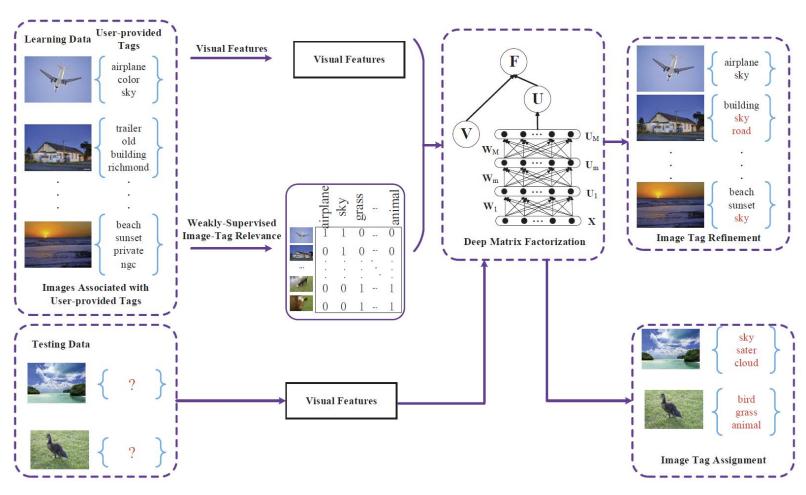
Semi-NMF: $X_{\pm} \approx F_{\pm} G_{\pm}^T$

Convex-NMF: $X_{\pm} \approx X_{\pm} W_{+} G_{+}^{T}$



Ding et al., TPAMI2015

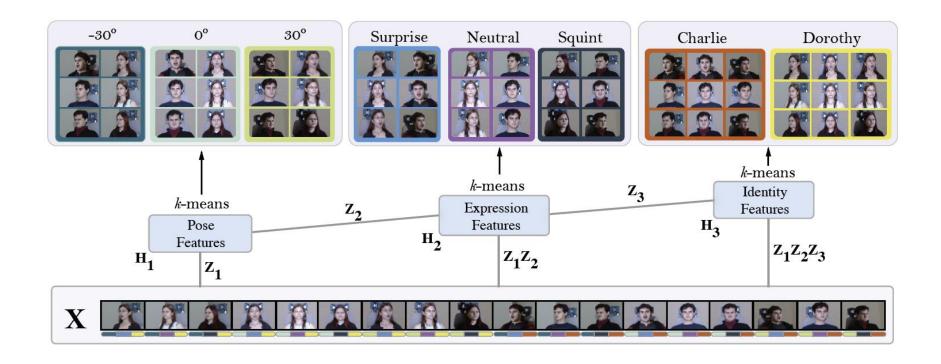
Deep Matrix Factorization



Li and Tang, MMML 2015



Deep Semi-NMF Model



Trigerous et al., TPAMI 2015

Multivariate Statistics

- Multivariate analysis of variance (MANOVA)
- Principal components analysis (PCA)
- Factor analysis
- Linear discriminant analysis (LDA)
- Canonical correlation analysis (CCA)
- Correspondence analysis
- Canonical correspondence analysis
- Multidimension scaling
- Multivariate regression
- Discriminant analysis