



Language Technologies Institute



Multimodal Machine Learning

Lecture 4.2: Multivariate Statistics and Coordinated Representations Louis-Philippe Morency

* Original version co-developed with Tadas Baltrusaitis

Lecture Objectives

- Quick recap
- Multivariate statistical analysis
 - Basic concepts (multivariate, covariance,...)
- Canonical Correlation Analysis
- Deep Correlation Networks
 - Deep CCA, DCCA-AutoEncoder
- Multi-view clustering
 - Nonnegative Matrix Factorization
- Multi-view latent intact space
 - Autoencoder in Autoencoder networkds



Administrative Stuff



Language Technologies Institute

Carnegie Mellon University

Lecture Schedule

Classes	Tuesday Lectures	Thursday Lectures
Week 1 8/27 & 8/29	 Course introduction Research and technical challenges Course syllabus and requirements 	 Multimodal applications and datasets Research tasks and datasets Team projects
Week 2 9/3 & 9/5 ***	 Basic concepts: neural networks Language, visual and acoustic Loss functions and neural networks 	 Basic concepts: network optimization Gradients and backpropagation Practical deep model optimization
Week 3 9/10 & 9/12 * <i>Pre-proposal</i> * Week 4 9/17 & 9/19	 Convolutional neural networks Convolutional kernels and CNNs Residual networks Multimodal representation learning Multimodal auto-encoders Multimodal joint representations 	 Recurrent neural networks Gated networks and LSTM Backpropagation Through Time Coordinated representations Deep canonical correlation analysis Non-negative matrix factorization
Week 5 9/24 & 9/26	 Multimodal alignment Explicit - dynamic time warping Implicit - attention models 	 Structured representations Module networks Tree-based and stack models
Week 6 10/1 & 10/3	First project assignment - Presentations	First assignment due on Sunday 10/6



Lecture Schedule

Classes	Tuesday Lectures	Thursday Lectures
Week 7	Alignment and representation	Probabilistic graphical models
10/8 & 10/10	Multi-head attention	Dynamic Bayesian networks
	Multimodal transformers	Coupled and factor HMMs
Week 8	Reinforcement learning	Multimodal RL
10/15 & 10/17	 Markov decision process 	Deep Q learning
	 Q learning and policy gradients 	 Multimodal applications
Week 9	Discriminative graphical models	Generative Models and Translation
10/22 & 10/24	Boltzmann distribution and CRFs	Variational auto-encoder
	• Continuous and fully-connected CRFs	Generative adversarial approaches
Week 10	Multimodal fusion and co-learning	New directions in Multimodal ML
10/29 & 10/31	 Multi-kernel learning and fusion 	 Overview of recent approaches in
	 Multimodal transfer learning 	multimodal machine learning
Week 11 11/5 & 11/7	Mid-term project assignment - Presentation	Midterm due on 11/10.



Lecture Schedule

Classes	Tuesday Lectures	Thursday Lectures		
Week 12 11/12 & 11/14	 Multi-lingual representations Neural machine translation Guest lecture: Graham Neubig 	KnowledMult	ge representation imodal knowledge discovery	
Week 13 11/19 & 11/21	Thanksgiving week (no classes)	• Gues		
Week 14 11/26 & 11/28	 Language, vision and action Neural machine translation Guest lecture 	 Multimodal affective computing Emotion and sentiment analysis Guest lecture 		
Week 15 12/3 & 12/5	Final project assignment - Presentations		Final due on 12/8.	



Quick Recap



Language Technologies Institute



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

> Deep Multimodal Boltzmann machines





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

- Deep Multimodal Boltzmann machines
- Stacked Autoencoder





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

- Deep Multimodal Boltzmann machines
- Stacked Autoencoder
- Encoder-Decoder





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

- Deep Multimodal Boltzmann machines
- Stacked Autoencoder
- Encoder-Decoder
- Tensor Fusion representation



How Can We Learn Better Representations?

Coordinated Multimodal Representations

Coordinated multimodal embeddings

 Instead of projecting to a joint space enforce the similarity between unimodal embeddings







Coordinated Multimodal Representations

Learn (unsupervised) two or more coordinated representations from multiple modalities. A loss function is defined to bring closer these multiple representations.





Coordinated Multimodal Embeddings



[Huang et al., Learning Deep Structured Semantic Models for Web Search using Clickthrough Data, 2013]





Multimodal Vector Space Arithmetic

Nearest images



[Kiros et al., Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, 2014]





Multimodal Vector Space Arithmetic

Nearest images



[Kiros et al., Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, 2014]





Structured coordinated embeddings

Instead of or in addition to similarity add alternative structure





[Vendrov et al., Order-Embeddings of Images and Language, 2016]

[Jiang and Li, Deep Cross-Modal Hashing]



Carnegie Mellon University

Multivariate Statistical Analysis



Language Technologies Institute

Carnegie Mellon University

"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)



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

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

> If same probability for all observations x_i , then same as arithmetic mean **Variance** measures the spread of the observations

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

 \succ 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





Carnegie Mellon University

Definitions

Multivariate (multidimensional) random variables

(aka random vector)

 $\boldsymbol{X} = [X^{1}, X^{2}, X^{3}, \dots, X^{M}]$ $\boldsymbol{Y} = [Y^{1}, Y^{2}, Y^{3}, \dots, Y^{N}]$

Covariance matrix generalizes the notion of variance

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

Cross-covariance matrix generalizes the notion of covariance

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



Definitions

Multivariate (multidimensional) random variables

(aka random vector)

$$\boldsymbol{X} = [X^{1}, X^{2}, X^{3}, \dots, X^{M}]$$
$$\boldsymbol{Y} = [Y^{1}, Y^{2}, Y^{3}, \dots, Y^{N}]$$

Covariance matrix generalizes the notion of variance

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

Cross-covariance matrix generalizes the notion of covariance

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

26



Definitions – Matrix Operations

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

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





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







Eigenvalues and Eigenvectors

Eigenvalue decomposition

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

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

Eigenvalue equation:



- λ : a scalar (could be **zero**)
- **x**: a **nonzero** vector in R^n







Carnegie Mellon University

Singular Value Decomposition (SVD)

SVD expresses any matrix A as

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

 The columns of U are eigenvectors of AA^T, and the columns of V are eigenvectors of A^TA.

$$\mathbf{A}\mathbf{A}^{T}\mathbf{u}_{i} = s_{i}^{2}\mathbf{u}_{i}$$
$$\mathbf{A}^{T}\mathbf{A}\mathbf{v}_{i} = s_{i}^{2}\mathbf{v}_{i}$$







Language Technologies Institute

Carnegie Mellon University

Multi-view Learning





demographic properties



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:

$$(\boldsymbol{u}^*, \boldsymbol{v}^*) = \operatorname*{argmax}_{\boldsymbol{u}, \boldsymbol{v}} corr(\boldsymbol{H}_{\boldsymbol{x}}, \boldsymbol{H}_{\boldsymbol{y}})$$

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





Correlated Projection

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

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



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



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

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

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

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

$$= \operatorname*{argmax}_{\boldsymbol{u}, \boldsymbol{v}} \frac{\boldsymbol{u}^{T}\boldsymbol{\Sigma}_{\boldsymbol{X}\boldsymbol{Y}}\boldsymbol{v}}{\sqrt{\boldsymbol{u}^{T}\boldsymbol{\Sigma}_{\boldsymbol{X}\boldsymbol{X}}\boldsymbol{u}}\sqrt{\boldsymbol{v}^{T}\boldsymbol{\Sigma}_{\boldsymbol{Y}\boldsymbol{Y}}\boldsymbol{v}}}$$
where

$$\Sigma_{XY} = cov(X,Y) = XY^{T}$$
if both X, Y have 0 mean

$$\mu_{X} = \mathbf{0} \quad \mu_{Y} = \mathbf{0}$$



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

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

2

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

$$\begin{aligned} u_{(i)}^{T} \Sigma_{XY} v_{(j)} &= u_{(j)}^{T} \Sigma_{XY} v_{(i)} = \mathbf{0} & \text{for } i \neq j \\ |U \Sigma_{XY} V| &= tr(U \Sigma_{XY} V) & \text{where } U = [u_{(1)}, u_{(2)}, \dots, u_{(k)}] \\ & \text{and } V = [v_{(1)}, v_{(2)}, \dots, v_{(k)}] \end{aligned}$$



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

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: $U^T \Sigma_{XX} U = V^T \Sigma_{YY} V = I, u^T_{(j)} \Sigma_{XY} v_{(i)} = 0$



3

$$\Sigma = \begin{bmatrix} \mathbf{\Sigma}_{XX} & \mathbf{\Sigma}_{YX} \\ \mathbf{\Sigma}_{XY} & \mathbf{\Sigma}_{YY} \end{bmatrix} \stackrel{U,V}{\Rightarrow} \begin{bmatrix} 1 & 0 & 0 & \lambda_1 & 0 & 0 \\ 0 & 1 & 0 & 0 & \lambda_2 & 0 \\ 0 & 0 & 1 & 0 & 0 & \lambda_3 \\ 0 & 0 & 1 & 0 & 0 & \lambda_3 \\ \lambda_1 & 0 & 0 & 1 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 & 1 & 0 \\ 0 & \lambda_2 & 0 & 0 & 1 & 0 \\ 0 & \lambda_2 & 0 & 0 & 1 & 0 \\ 0 & \lambda_2 & 0 & 0 & 1 & 0 \\ 0 & 0 & \lambda_3 & 0 & 0 & 1 \end{bmatrix}$$



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

subject to: $U^T \Sigma_{XX} U = V^T \Sigma_{YY} V = I$, $u_{(j)}^T \Sigma_{XY} v_{(i)} = 0$
for $i \neq j$
How to solve it? > Lagrange Multipliers!
agrange function
 $L = tr(U^T \Sigma_{XY} V) + \alpha (U^T \Sigma_{YY} U - I) + \beta (V^T \Sigma_{YY} V - I)$
> 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 $\lambda = 4\alpha\beta$





• Language Technologies Institute

Text

X

Image

Y



Exploring Deep Correlation Networks



Language Technologies Institute

Carnegie Mellon University

Same objective function as CCA:



43

Andrew et al., ICML 2013



44

Training procedure:

 Pre-train the models parameters using denoising autoencoders



Andrew et al., ICML 2013



Training procedure:

- Pre-train the models parameters using denoising autoencoders
- Optimize the CCA objective functions using *I* 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

- 1. 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



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
- Multidimensional scaling
- Multivariate regression
- Discriminant analysis



Multi-View Clustering



Language Technologies Institute





Clustering definition: partition a set of data samples such that similar samples are grouped, and dissimilar samples are divided

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"







"Soft" Clustering: Nonnegative Matrix Factorization

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



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 Semi-NMF Model



Trigerous et al., TPAMI 2015





Learn data partitioning from multiple views (modalities)

Views: different sources in diverse domains or obtained from various feature collectors or modalities

Example: Multiple views in computer vision - LBP, SIFT, HOG



Yan Yang and Hao Wang, Multi-view Clustering: A Survey, Big data mining and analytics, Volume 1, Number 2, June 2018



Principles of Multi-View Clustering

Two important principles:



Consensus principle: maximize consistency across multiple distinct views

2 **Complementarity principle:** multiple views needed to get more comprehensive and accurate descriptions



Yan Yang and Hao Wang, Multi-view Clustering: A Survey, Big data mining and analytics, Volume 1, Number 2, June 2018



Multi-view subspace clustering

Definition: learns a unified feature representation from all the view subspaces by assuming that all views share this representation







Enforcing Data Clustering in Deep Networks

How to enforce data clustering in our (multimodal) deep learning algorithms?





Carnegie Mellon University

Deep Matrix Factorization



Li and Tang, MMML 2015





Other Multi-View Clustering Approaches

Graph-based clustering: search for a fusion graph (or network) across all views and then perform clustering



Yan Yang and Hao Wang, Multi-view Clustering: A Survey, Big data mining and analytics, Volume 1, Number 2, June 2018



Other Multi-View Clustering Approaches

Co-training: bootstraps the clustering of different views by using the learning knowledge from other views



Yan Yang and Hao Wang, Multi-view Clustering: A Survey, Big data mining and analytics, Volume 1, Number 2, June 2018





Auto-Encoder in Auto-Encoder Network



Language Technologies Institute

Carnegie Mellon University





Carnegie Mellon University

Xu et al., TPAMI 2015

Given multiple views z_i from the same "object":



1) There is an "intact" representation which is *complete* and *not damaged*

2) The views z_i are partial (and possibly degenerated) representations of the intact representation



Auto-Encoder in Auto-Encoder Network

Zhang et al., CVPR 2019



Carnegie Mellon University