# Multimodal Machine Learning

## Lecture 4.2: Coordinated Representations

**Louis-Philippe Morency**

* Original version co-developed with Tadas Baltrusaitis

# Administrative Stuff

# Piazza Live Q&A – Reminder

Language Technologies Institute

Carnegie Mellon University

# Lecture Schedule

| Classes | Tuesday Lectures | Thursday Lectures |
|---|---|---|
| **Week 1** <br> 9/1 & 9/3 | **Course introduction** <br> • Research and technical challenges <br> • Course syllabus and requirements | **Multimodal applications and datasets** <br> • Research tasks and datasets <br> • Team projects |
| **Week 2** <br> 9/8 & 9/10 | **Basic concepts: neural networks** <br> • Language, visual and acoustic <br> • Loss functions and neural networks | **Basic concepts: network optimization** <br> • Gradients and backpropagation <br> • Practical deep model optimization |
| **Week 3** <br> 9/15 & 9/17 | **Visual unimodal representations** <br> • Convolutional kernels and CNNs <br> • Residual network and skip connection | **Language unimodal representations** <br> • Gated networks and LSTM <br> • Backpropagation Through Time |
| **Week 4** <br> 9/22 & 9/24 | **Multimodal representation learning** <br> • Multimodal auto-encoders <br> • Multimodal joint representations | **Coordinated representations** <br> • Deep canonical correlation analysis <br> • Non-negative matrix factorization |
| **Week 5** <br> 9/29 & 10/1 | **Multimodal alignment** <br> • Explicit - dynamic time warping <br> • Implicit - attention models | **Alignment and representation** <br> • Self-attention models <br> • Multimodal transformers |
| **Week 6** <br> 10/6 & 10/8 | ***First project assignment*** *(live working sessions instead of lectures)* | |

> **First project assignment**
> Presentations due Friday 10/9
> Reports due Sunday 10/11
> Peer feedback due Friday 10/16

Language Technologies Institute

Carnegie Mellon University

# Lecture Schedule

| Classes | Tuesday Lectures | Thursday Lectures |
|---|---|---|
| **Week 7**<br>10/13 & 10/15 | **Alignment and translation**<br>• Module networks<br>• Connectionist temporal classification | **Probabilistic graphical models**<br>• Dynamic Bayesian networks<br>• Coupled and factor HMMs |
| **Week 8**<br>10/20 & 10/22 | **Discriminative graphical models**<br>• Conditional random fields<br>• Continuous and fully-connected CRFs | **Neural Generative Models**<br>• Variational auto-encoder<br>• Generative adversarial networks |
| **Week 9**<br>10/27 & 10/29 | **Reinforcement learning**<br>• Markov decision process<br>• Q learning and policy gradients | **Multimodal RL**<br>• Deep Q learning<br>• Multimodal applications |
| **Week 10**<br>11/3 & 11/5 | **Fusion and co-learning**<br>• Multi-kernel learning and fusion<br>• Few shot learning and co-learning | **New research directions**<br>• Recent approaches in multimodal ML |
| **Week 11**<br>11/10 & 11/12 | ***Mid-term project assignment*** *(live working sessions instead of lectures)* | |

**Midterm project assignment**
Presentations due Friday 11/13
Reports due Sunday 11/15
Peer feedback due Friday 11/20

Language Technologies Institute

Carnegie Mellon University

# Lecture Schedule

| Classes | Tuesday Lectures | Thursday Lectures |
|---|---|---|
| **Week 12** <br> 11/17 & 11/19 | **Embodied Language Grounding** <br> • Connecting Language to Action <br> • Guest lecture: Yonatan Bisk | **Multi-lingual representations** <br> • Tentative topic <br> • Guest lecture: To be confirmed |
| **Week 13** <br> 11/24 & 11/26 | ***Thanksgiving week*** *(no lectures)* | |
| **Week 14** <br> 12/1 & 12/3 | **Learning to connect text and images** <br> • Discourse approaches, text & images <br> • Guest lecture: Malihe Alikhani | **Bias and fairness** <br> • Computational ethics <br> • Guest lecture: Yulia Tsvetkov |
| **Week 15** <br> 12/8 & 12/10 | ***Final project assignment*** *(live working sessions instead of lectures)* | |

> **Final project assignment**
> Presentations due Friday 12/11
> Reports due Sunday 12/13

Language Technologies Institute

Carnegie Mellon University

# GPU $50 Coupons - AWS

➡ First, create an account on AWS Educate portal:
https://aws.amazon.com/education/awseducate/

➡ Your account will need to be backed by your credit card

Be sure to setup billing alarms and monitor your spending!

➡ Refrain from including AWS credential in code/github

➡ To get your coupon, use your AndrewID and
the URL posted on Piazza

Language Technologies Institute

Carnegie Mellon University

# GPU $50 Coupons - GCP

➡️ Coupons can be redeemed at this address:
   https://console.cloud.google.com/education

Be sure to setup billing alarms and monitor your spending!

➡️ Refrain from including GCP credential in code/github

➡️ To get your coupon, use your AndrewID and
   the URL posted on Piazza

Language
Technologies
Institute

Carnegie
Mellon
University

# Multimodal Machine Learning

## Lecture 4.2: Coordinated Representations

**Louis-Philippe Morency**

*** Original version co-developed with Tadas Baltrusaitis**

# Lecture Objectives

- Quick recap
- Coordinated multimodal representations
- 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

# Quick Recap

# Multimodal Representation Learning

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

❑ Deep Multimodal Boltzmann machines

softmax

Text
$X$

Image
$Y$

# Multimodal Representation Learning

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

- ❑ Deep Multimodal Boltzmann machines
- ❑ Stacked Autoencoder

# Multimodal Representation Learning

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

- ❑ Deep Multimodal Boltzmann machines

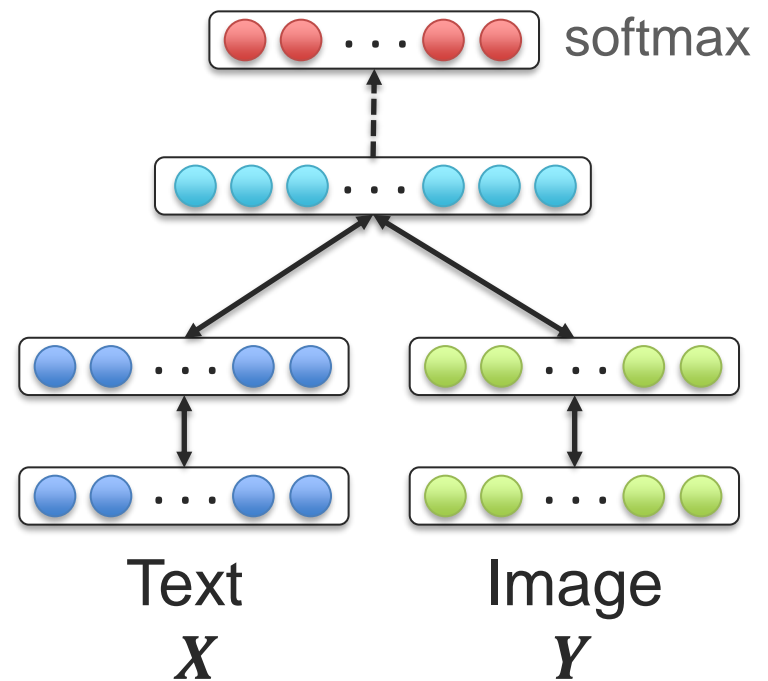- ❑ Stacked Autoencoder

- ❑ Encoder-Decoder

Text
*X*

Image
*Y*

# Multimodal Representation Learning

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

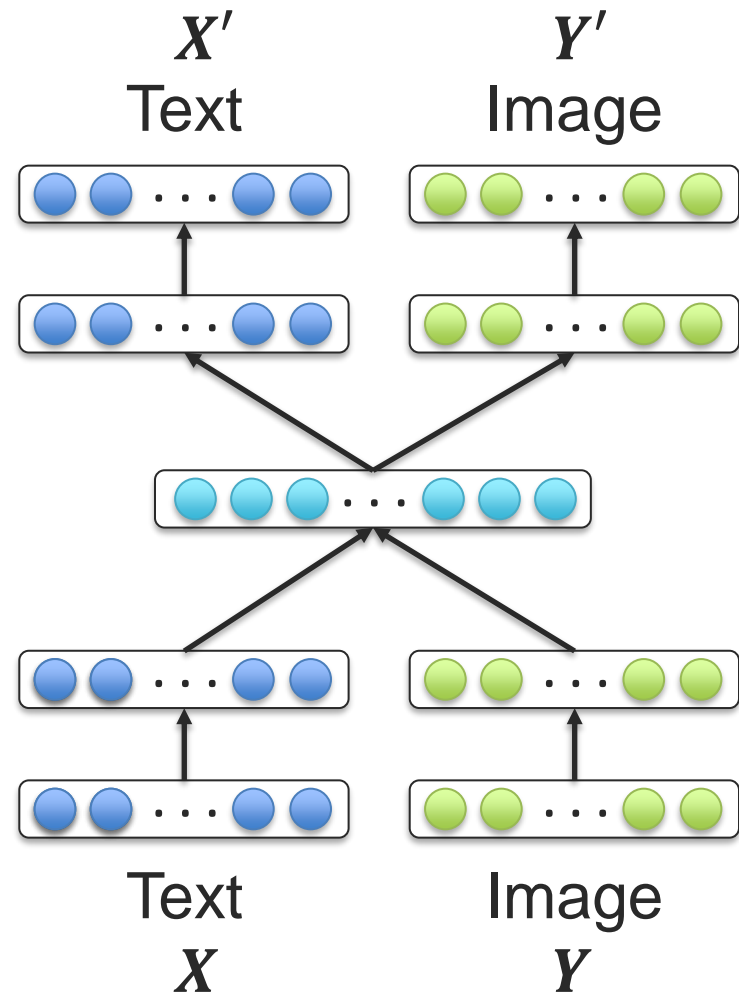- [ ] Deep Multimodal Boltzmann machines
- [ ] Stacked Autoencoder
- [ ] Encoder-Decoder
- [ ] Tensor Fusion representation

e.g. Sentiment

softmax

**Bimodal**

**Unimodal**

$h_m$

$h_x$  $h_y$

Text
$X$

Image
$Y$

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.



Similarity metric (e.g., cosine distance)

Text
$X$

Image
$Y$

# Coordinated Multimodal Embeddings



Distance(s,t)

Image features s

Text: *a parrot rides a tricycle*

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

Carnegie Mellon University

# Multimodal Vector Space Arithmetic

Nearest images



- blue + red =

- blue + yellow =

- yellow + red =

- white + red =

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

# Multimodal Vector Space Arithmetic



Nearest images

- day + night =

- flying + sailing =

- bowl + box =

- box + bowl =

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

Language Technologies Institute

Carnegie Mellon University

# 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]

# 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

$$\mu = E[X] = \sum_i x_i P(x_i)$$

➢ 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

➢ Variance is equal to the square of the standard deviation $\sigma$

**Covariance** measures how much two random variables change together

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

Language Technologies Institute

Carnegie Mellon University

# Definitions

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

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

Language Technologies Institute

Carnegie Mellon University

# Pearson Correlation Examples

Language Technologies Institute

Carnegie Mellon University

# Definitions

Multivariate (multidimensional) random variables

*(aka random vector)*

$$X = [X^1, X^2, X^3, \ldots, X^M]$$

$$Y = [Y^1, Y^2, Y^3, \ldots, Y^N]$$

**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

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

Language Technologies Institute

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

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

**Cross-covariance matrix** generalizes the notion of covariance

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

Language Technologies Institute

Carnegie Mellon University

# Definitions – Matrix Operations

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

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

Language Technologies Institute

Carnegie Mellon University

# 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

Language Technologies Institute

Carnegie Mellon University

# Eigenvalues and Eigenvectors

Eigenvalue decomposition

If $A$ is an $n \times n$ matrix, do there exist nonzero vectors $\mathbf{x}$

in $R^n$ 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

Geometric Interpretation

$A$: an $n \times n$ matrix

$\lambda$: a scalar (could be **zero**)

$\mathbf{x}$: a **nonzero** vector in $R^n$

# Singular Value Decomposition (SVD)

- SVD expresses any matrix $\mathbf{A}$ as

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

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

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

# Canonical Correlation Analysis

**Carnegie Mellon University**

# Multi-view Learning

$X$

$Y$

demographic properties

responses to survey

audio features at time $i$

video features at time $i$

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

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

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

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

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

Language Technologies Institute

Carnegie Mellon University

# Correlated Projection

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

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



Two views $\boldsymbol{X}, \boldsymbol{Y}$ where same instances have the same color

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

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

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

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

where
$$\boldsymbol{\Sigma}_{XY} = cov(\boldsymbol{X}, \boldsymbol{Y}) = \boldsymbol{X}\boldsymbol{Y}^T$$

if both $\boldsymbol{X}, \boldsymbol{Y}$ have 0 mean

$$\boldsymbol{\mu}_X = \boldsymbol{0} \quad \boldsymbol{\mu}_Y = \boldsymbol{0}$$

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

$$= \underset{\boldsymbol{u},\boldsymbol{v}}{\mathrm{argmax}}\ \frac{\boldsymbol{u}^T\boldsymbol{\Sigma}_{XY}\boldsymbol{v}}{\sqrt{\boldsymbol{u}^T\boldsymbol{\Sigma}_{XX}\boldsymbol{u}}\sqrt{\boldsymbol{v}^T\boldsymbol{\Sigma}_{YY}\boldsymbol{v}}}$$

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

We want to learn multiple projection pairs $\left( \boldsymbol{u}_{(i)} X, \boldsymbol{v}_{(i)} Y \right)$:

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

**②** We want these multiple projection pairs to be orthogonal ("canonical") to each other:

$$\boldsymbol{u}_{(i)}^T \boldsymbol{\Sigma}_{XY} \boldsymbol{v}_{(j)} = \boldsymbol{u}_{(j)}^T \boldsymbol{\Sigma}_{XY} \boldsymbol{v}_{(i)} = \boldsymbol{0} \qquad \text{for } i \neq j$$

$$|\boldsymbol{U} \boldsymbol{\Sigma}_{XY} \boldsymbol{V}| = tr(\boldsymbol{U} \boldsymbol{\Sigma}_{XY} \boldsymbol{V}) \quad \text{where } \boldsymbol{U} = [\boldsymbol{u}_{(1)}, \boldsymbol{u}_{(2)}, \ldots, \boldsymbol{u}_{(k)}]$$

$$\text{and } \boldsymbol{V} = [\boldsymbol{v}_{(1)}, \boldsymbol{v}_{(2)}, \ldots, \boldsymbol{v}_{(k)}]$$

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

$$(U^*, V^*) = \underset{U,V}{\mathrm{argmax}} \frac{tr(U^T \Sigma_{XY} V)}{\sqrt{U^T \Sigma_{XX} U} \sqrt{V^T \Sigma_{YY} V}}$$

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

subject to: $\quad U^T \Sigma_{XX} U = V^T \Sigma_{YY} V = I, \; u_{(j)}^T \Sigma_{XY} v_{(i)} = 0$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ for $i \neq j$

# Canonical Correlation Analysis

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

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

$$\Sigma = \begin{bmatrix} \boldsymbol{\Sigma_{XX}} & \boldsymbol{\Sigma_{YX}} \\ \hline \boldsymbol{\Sigma_{XY}} & \boldsymbol{\Sigma_{YY}} \end{bmatrix} \overset{\boldsymbol{U,V}}{\Longrightarrow} \begin{bmatrix} 1 & 0 & 0 & \lambda_1 & 0 & 0 \\ 0 & 1 & 0 & 0 & \lambda_2 & 0 \\ 0 & 0 & 1 & 0 & 0 & \lambda_3 \\ \hline \lambda_1 & 0 & 0 & 1 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 & 1 & 0 \\ 0 & 0 & \lambda_3 & 0 & 0 & 1 \end{bmatrix}$$

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

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

subject to: $\boxed{\boldsymbol{U^T\Sigma_{XX}U}} = \boxed{\boldsymbol{V^T\Sigma_{YY}V = I}}, \boldsymbol{u_{(j)}^T\Sigma_{XY}v_{(i)} = 0}$ for $i \neq j$

How to solve it?  ➤ Lagrange Multipliers!

Lagrange function

$$\boldsymbol{L} = tr(\boldsymbol{U^T\Sigma_{XY}V}) + \alpha(\boldsymbol{U^T\Sigma_{YY}U - I}) + \beta(\boldsymbol{V^T\Sigma_{YY}V - I})$$

➤ And then find stationary points of $L$:  $\dfrac{\partial L}{\partial \boldsymbol{U}} = 0$  $\dfrac{\partial L}{\partial \boldsymbol{V}} = 0$

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

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

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

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

subject to: $\boldsymbol{U^T \Sigma_{XX} U = V^T \Sigma_{YY} V = I}, \boldsymbol{u}_{(j)}^T \boldsymbol{\Sigma_{XY} v}_{(i)} = \boldsymbol{0}$ for $i \neq j$

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

$$(\boldsymbol{U^*, V^*}) = (\boldsymbol{\Sigma}_{XX}^{-1/2} \boldsymbol{U}_{SVD}, \boldsymbol{\Sigma}_{YY}^{-1/2} \boldsymbol{V}_{SVD})$$

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

Eigenvalues

Eigenvectors

Eigenvalue equations
$$\boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY} \boldsymbol{\Sigma}_{YY}^{-1} \boldsymbol{\Sigma}_{XY}^T \boldsymbol{U} = \lambda \boldsymbol{U}$$

$$\boldsymbol{\Sigma}_{YY}^{-1} \boldsymbol{\Sigma}_{XY}^T \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY} \boldsymbol{V} = \lambda \boldsymbol{V} \quad \text{where } \lambda = 4\alpha\beta$$

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

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

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

1. Linear projections maximizing correlation

2. Orthogonal projections

3. Unit variance of the projection vectors



projection of Y

projection of X

$H_x$     $H_y$

$U$     $V$

Text     Image
$X$     $Y$

Language Technologies Institute

Carnegie Mellon University

# Exploring Deep Correlation Networks

# Deep Canonical Correlation Analysis

Same objective function as CCA:

$$\underset{V,U,W_x,W_y}{\text{argmax}} \; corr(H_x, H_y)$$

And need to compute gradients:

$$\frac{\partial corr(H_x, H_y)}{\partial U}$$

$$\frac{\partial corr(H_x, H_y)}{\partial V}$$

Andrew et al., ICML 2013

Language Technologies Institute

Carnegie Mellon University

# Deep Canonical Correlation Analysis

**Training procedure:**

1. Pre-train the models parameters using denoising autoencoders

$X'$
Text

$Y'$
Image

$H_x$

$H_y$

$U$

$V$

$W_x$

$W_y$

Text
$X$

Image
$Y$

Andrew et al., ICML 2013

Language Technologies Institute

Carnegie Mellon University
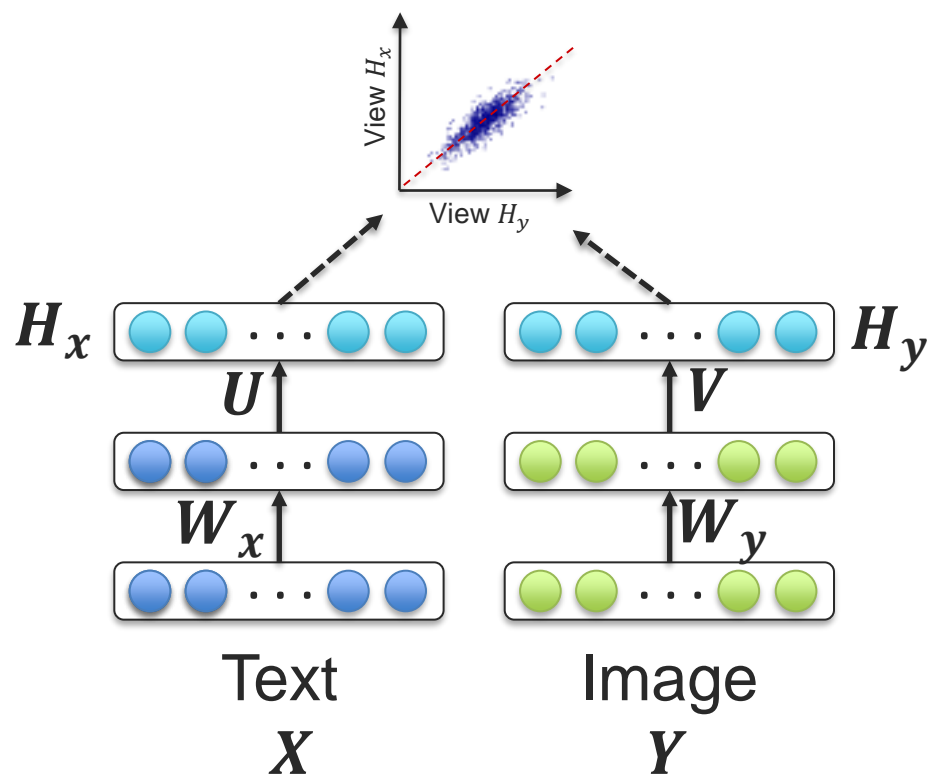
# Deep Canonical Correlation Analysis

**Training procedure:**

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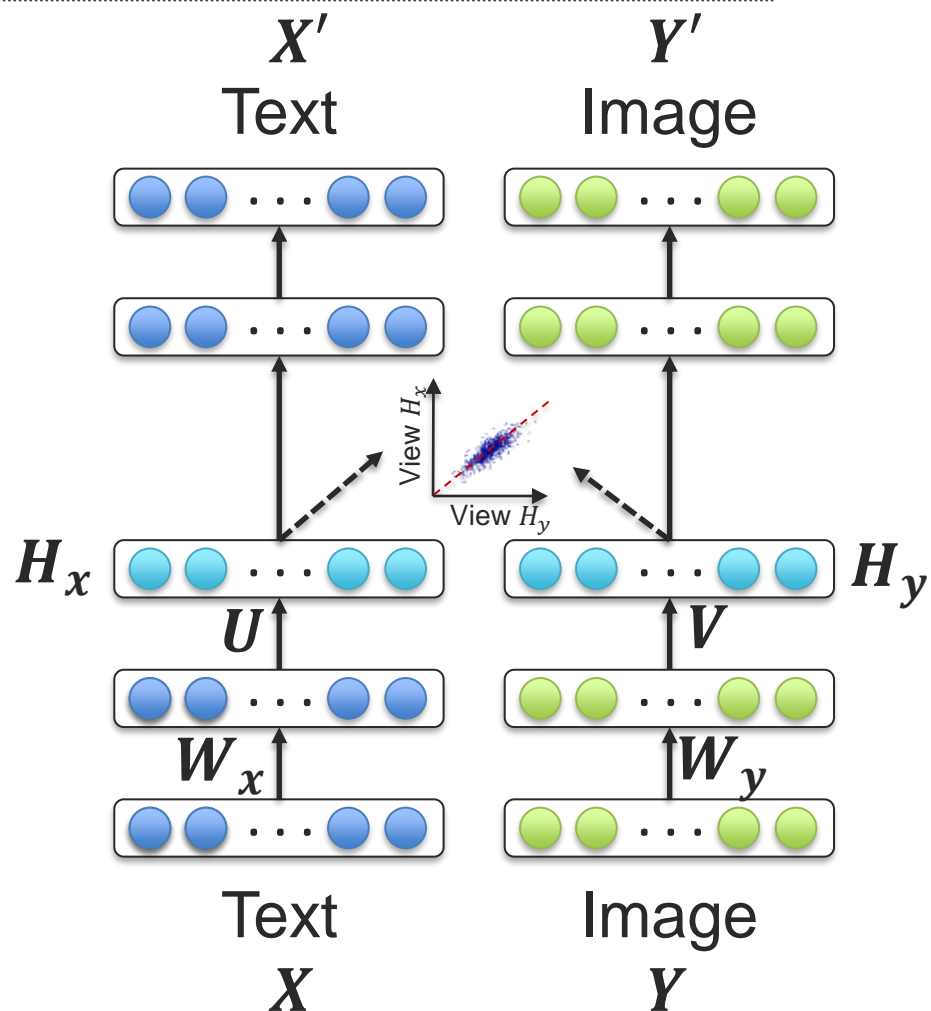2. Optimize the CCA objective functions using large mini-batches or full-batch (L-BFGS)

$H_x$ ⊙⊙ ⋯ ⊙⊙   ⊙⊙ ⋯ ⊙⊙ $H_y$

$U$   $V$

$W_x$   $W_y$

Text
$X$

Image
$Y$

View $H_x$

View $H_y$

Andrew et al., ICML 2013

Language Technologies Institute

Carnegie Mellon University

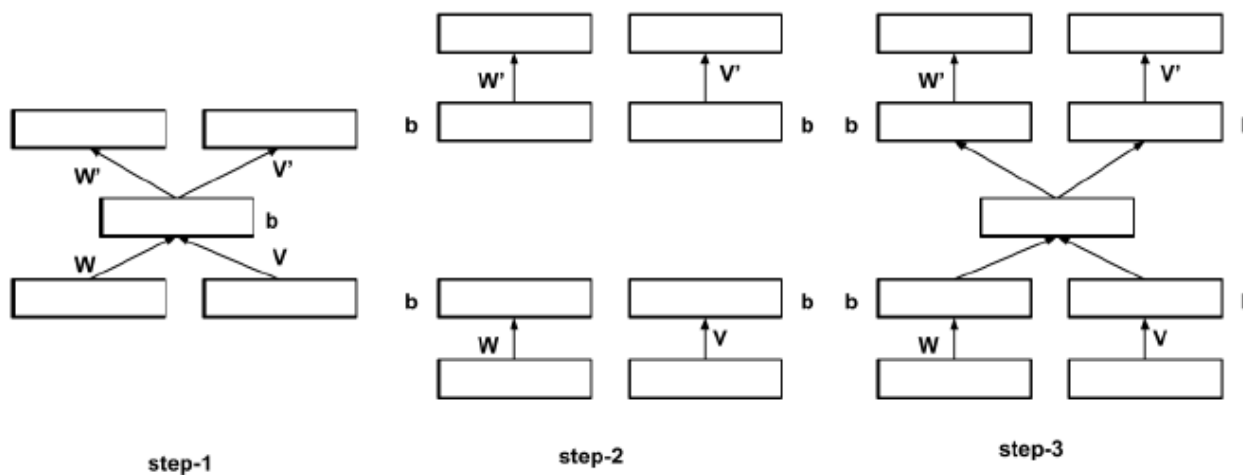# Deep Canonically Correlated Autoencoders (DCCAE)

Jointly optimize for DCCA and autoencoders loss functions

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

$X'$ Text       $Y'$ Image

$H_x$        $H_y$

$U$       $V$

$W_x$       $W_y$

Text $X$       Image $Y$

View $H_x$ / View $H_y$

Wang et al., ICML 2015

Language Technologies Institute

Carnegie Mellon University

# 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

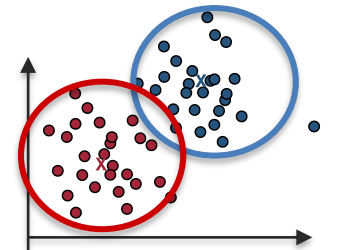# Multi-View Clustering

# Data Clustering

**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"

Image

# "Soft" Clustering: Nonnegative Matrix Factorization

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

$$\left( \quad X \quad \right) = \left( \; F \; \right) \left( \quad G \quad \right)$$

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

Language Technologies Institute

Carnegie Mellon University

# Semi-NMF and Other Extensions

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

$$\text{NMF:} \qquad X_+ \approx F_+ G_+^T$$

$$\boxed{\text{Semi-NMF:} \qquad X_\pm \approx F_\pm G_+^T}$$

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

Image

Ding et al., TPAMI2015

# Deep Semi-NMF Model



Trigerous et al., TPAMI 2015
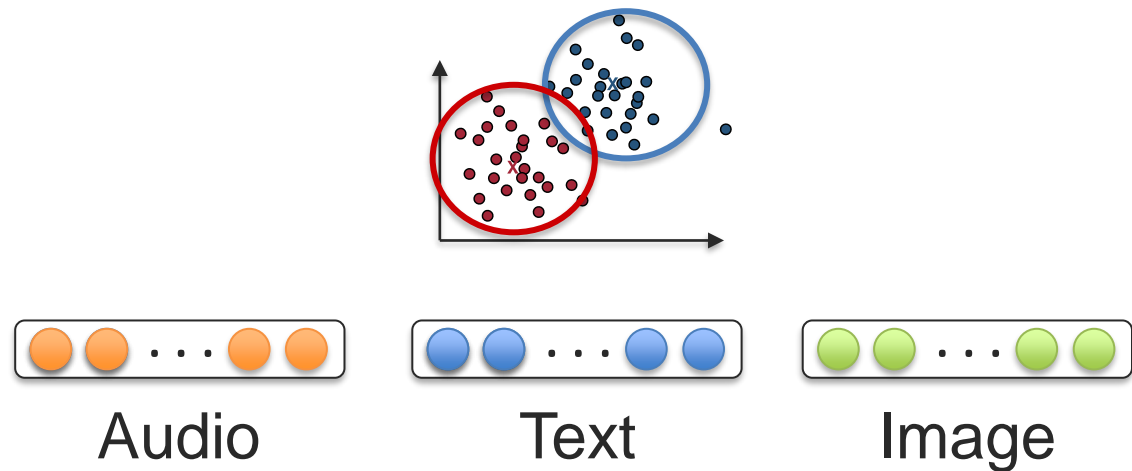
Language Technologies Institute

Carnegie Mellon University

# Multi-View Clustering

Learn data partitioning from multiple views (modalities)

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

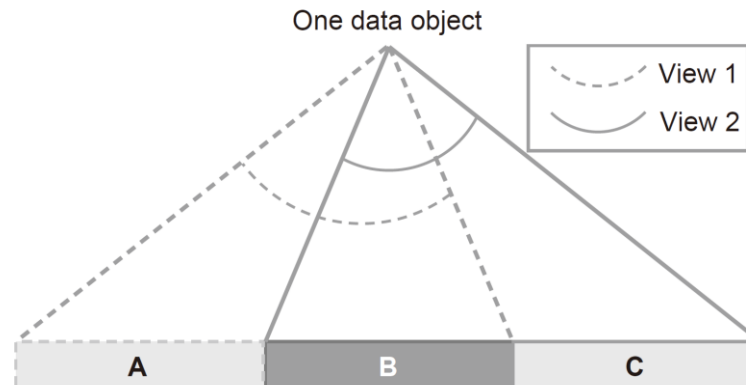Example: Multiple views in computer vision - LBP, SIFT, HOG



Audio            Text            Image

Language Technologies Institute

Carnegie Mellon University

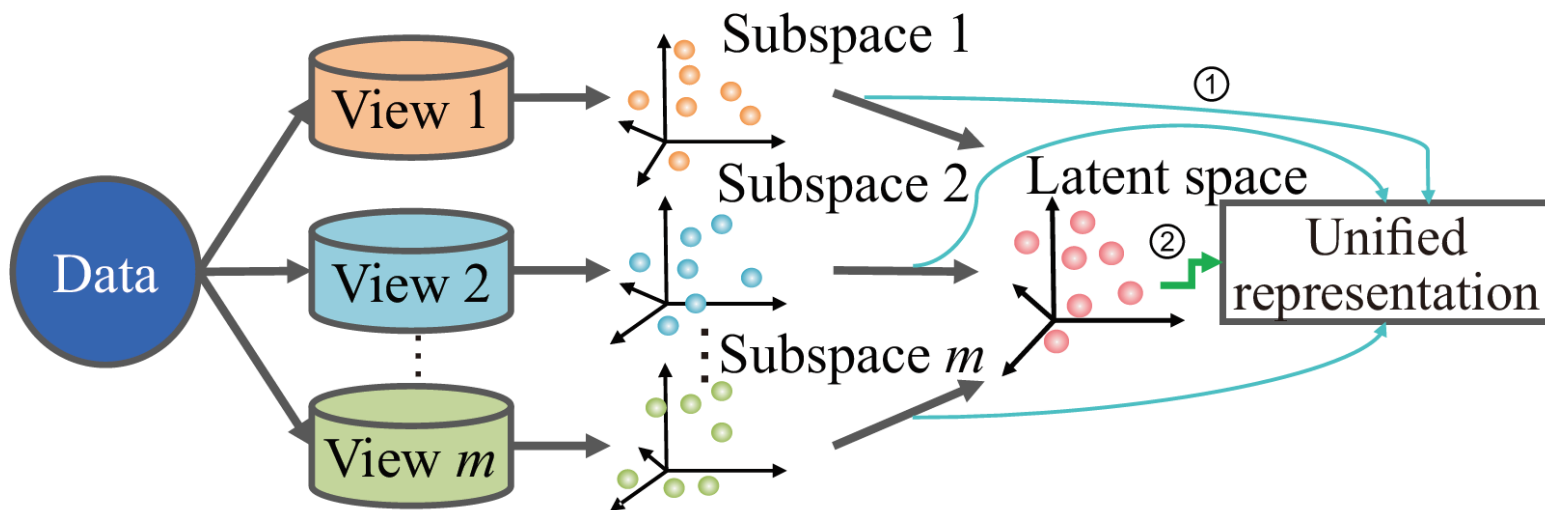# Principles of Multi-View Clustering

Two important principles:

(1) **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
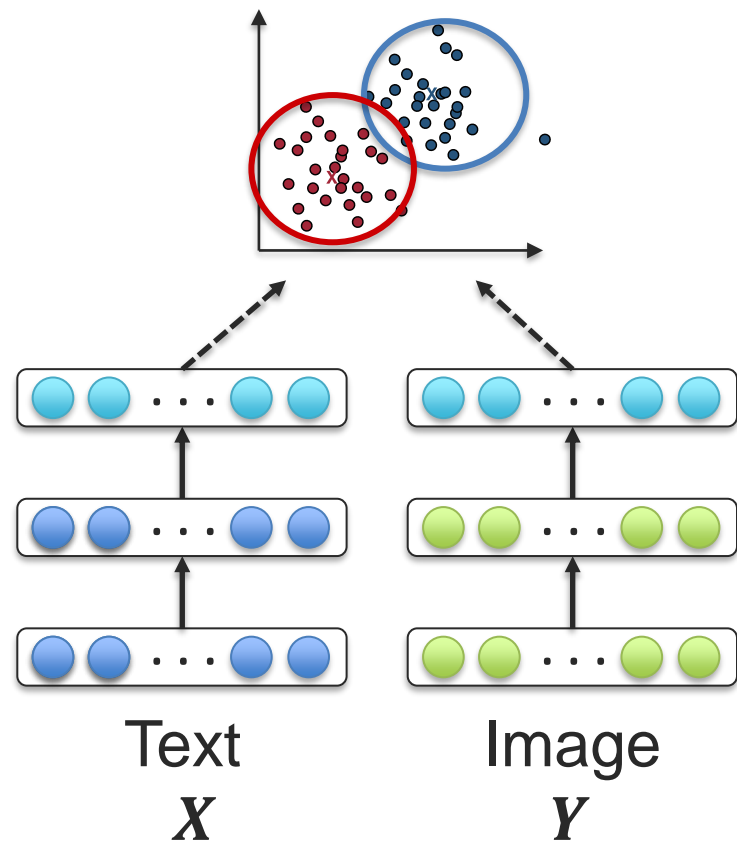
# Multi-view subspace clustering

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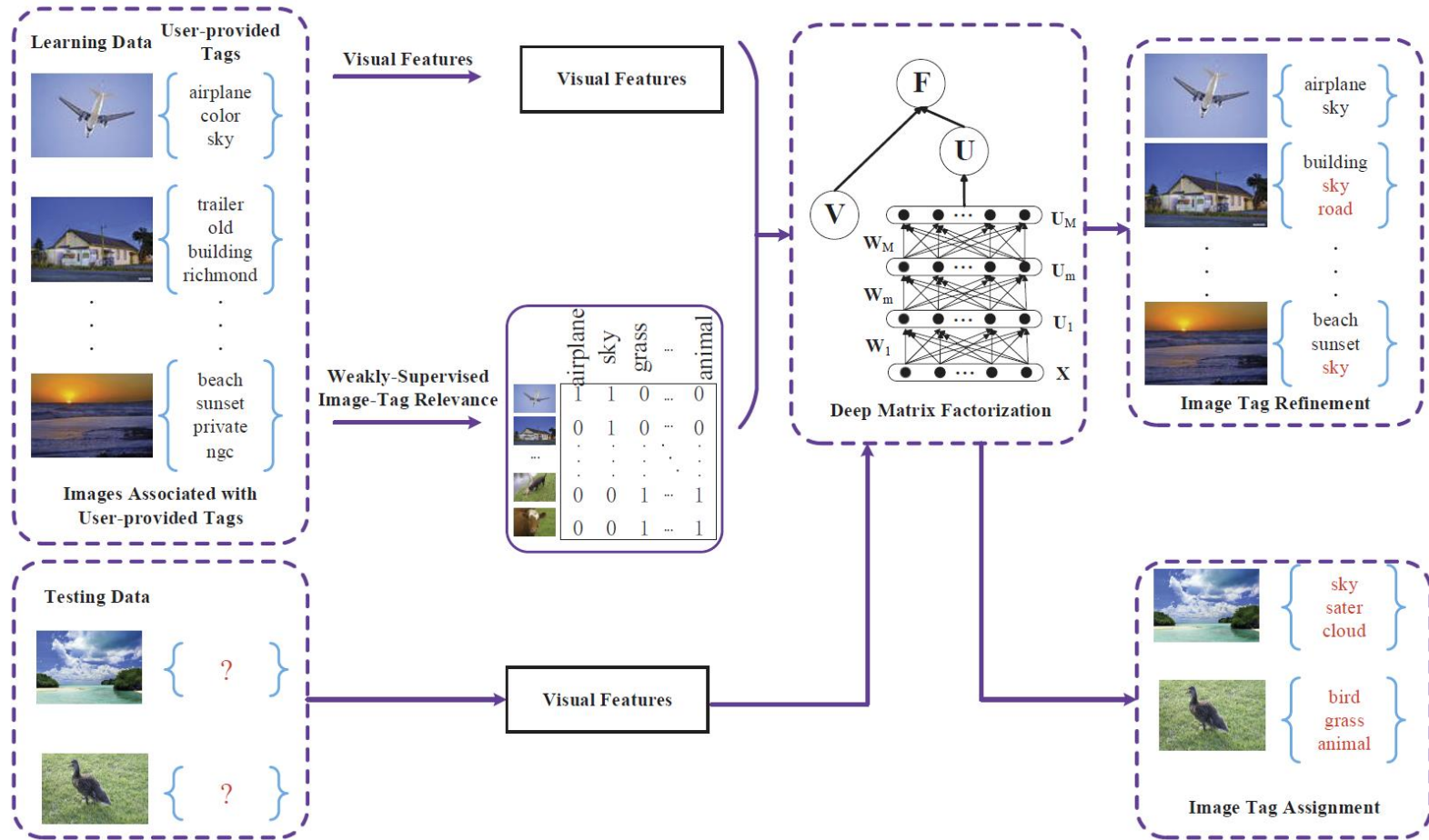
Language Technologies Institute

Carnegie Mellon University

# Enforcing Data Clustering in Deep Networks

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

Language Technologies Institute

Carnegie Mellon University

# Deep Matrix Factorization



Li and Tang, MMML 2015

Language Technologies Institute

Carnegie Mellon University
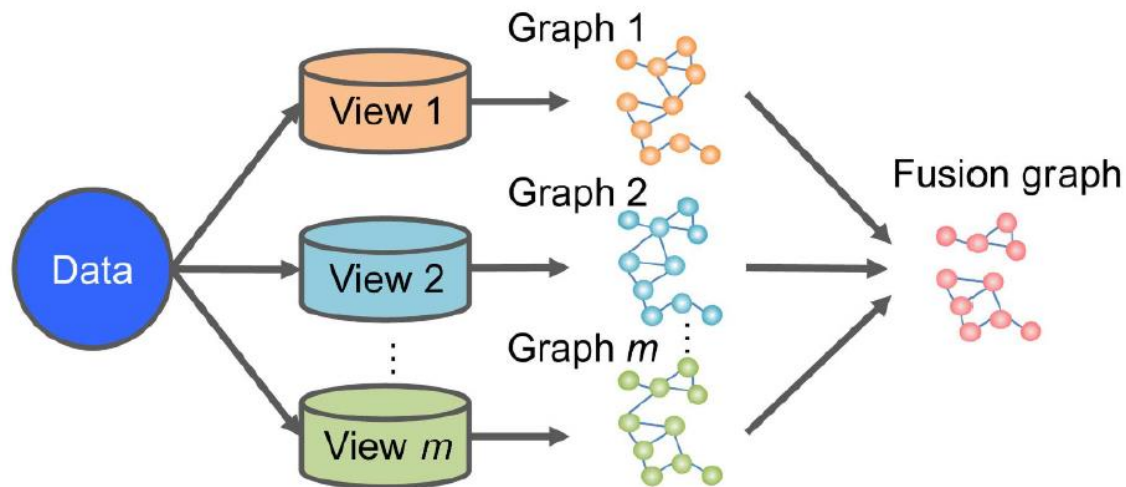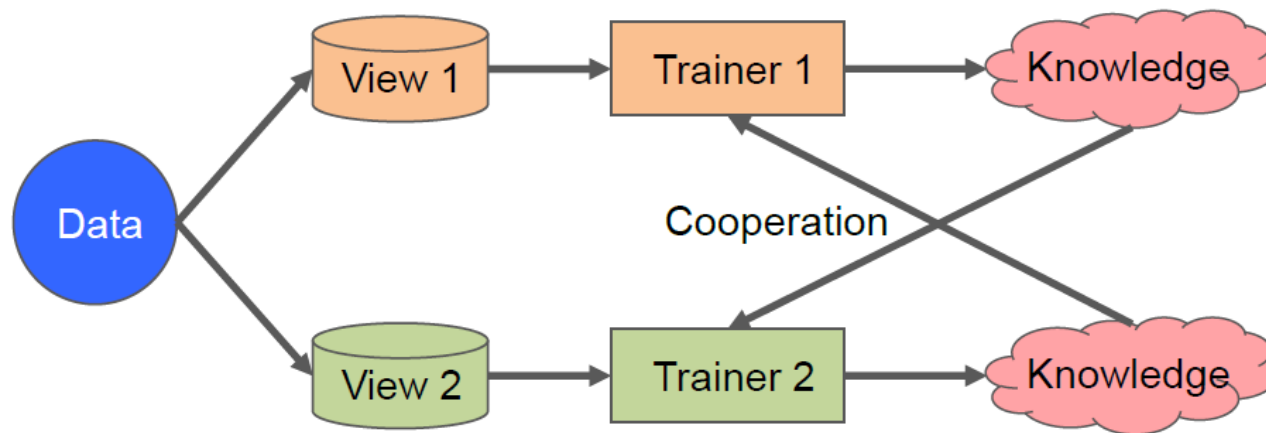
# Other Multi-View Clustering Approaches

**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

Language Technologies Institute

Carnegie Mellon University

# Other Multi-View Clustering Approaches

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



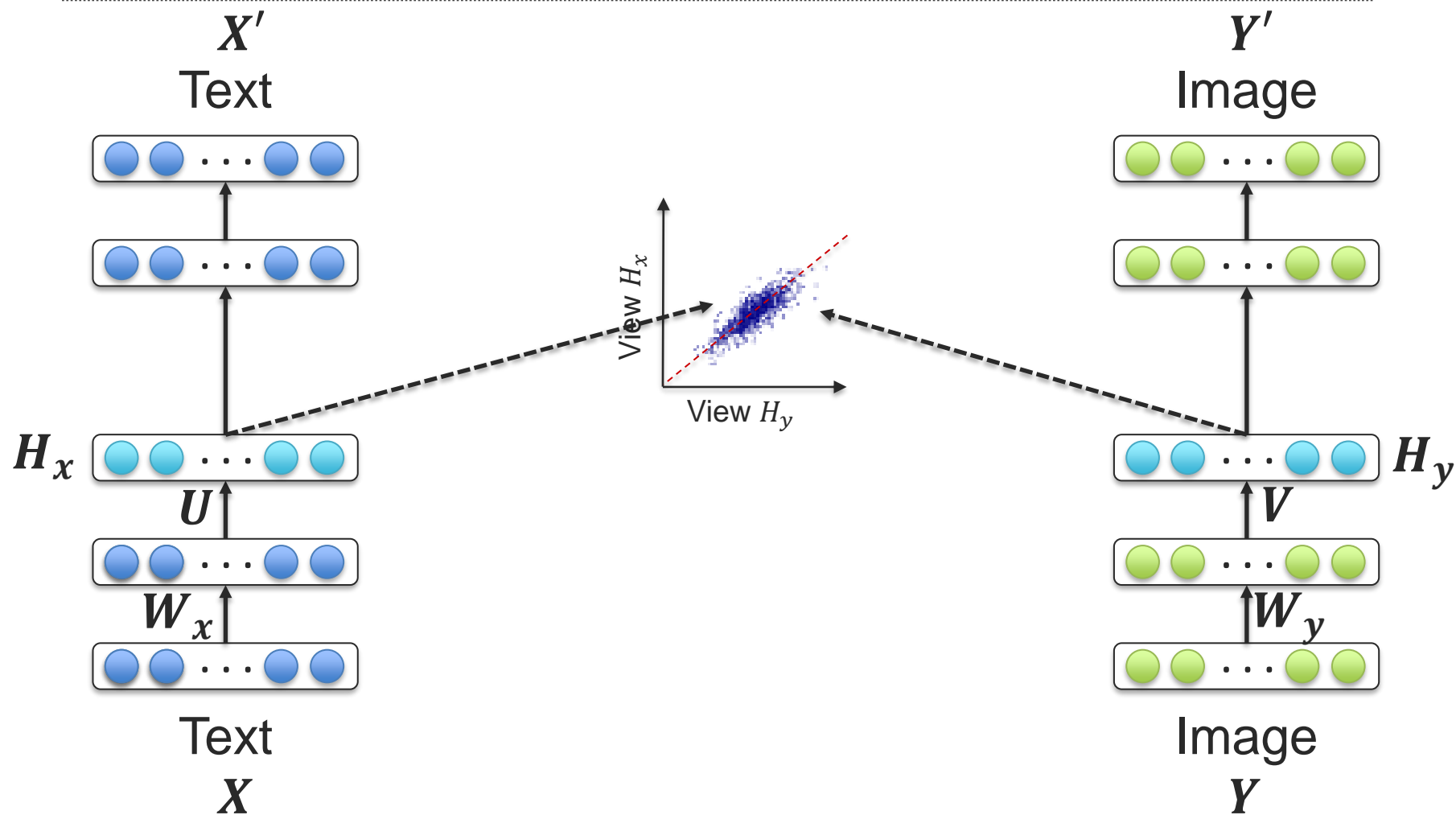Yan Yang and Hao Wang, Multi-view Clustering: A Survey, Big data mining and analytics, Volume 1, Number 2, June 2018

Language Technologies Institute

Carnegie Mellon University

# Auto-Encoder in Auto-Encoder Network

# Deep Canonically Correlated Autoencoders (DCCAE)

# Multi-view Latent "Intact" Space

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

Reconstructed Text $Z^{(M,1)}$

Latent *Intact* Representation $H$

Reconstructed Image $Z^{(M,2)}$

$G^{(L,1)}$

$G^{(L,2)}$

Loss

Loss

Loss

Loss

$Z^{(\frac{M}{2},1)}$

$Z^{(\frac{M}{2},2)}$

Degradation network

Degradation network

Input Text $X^{(1)}$

Input Image $X^{(2)}$

Total Loss:

$$\min_{\{\Theta_{ae}^{(v)}, \Theta_{dg}^{(v)}\}_{v=1}^{V}, \mathbf{H}} \frac{1}{2} \sum_{v=1}^{V} \left( \left\| \mathbf{X}^{(v)} - \mathbf{Z}^{(M,v)} \right\|_{F}^{2} + \lambda \left\| \mathbf{Z}^{(\frac{M}{2},v)} - \mathbf{G}^{(L,v)} \right\|_{F}^{2} \right)$$

Latent variable