



Language Technologies Institute



Multimodal Machine Learning

Lecture 4.1: Multimodal Representations

Louis-Philippe Morency

* Original version co-developed with Tadas Baltrusaitis

Objectives of today's class

- Unsupervised representation learning
 - Restricted Boltzmann Machines
 - Autoencoders
 - Deep Belief Nets, Stacked autoencoders
- Multi-modal representations
 - Coordinated vs. joint representations
 - Multimodal Deep Boltzmann Machines
 - Deep Multimodal autoencoders
 - Tensor Fusion representation
 - Low-rank fusion representations





Administrative Stuff



Language Technologies Institute

Upcoming Schedule

- First project assignment:
 - Proposal presentation (10/1 and 10/3)
 - First project report (Sunday 10/6)
- Midterm project assignment
 - Midterm presentations (11/5 and 11/7)
 - Midterm report (Sunday 11/10)
- Final project assignment
 - Final presentation (12/3 & 12/5)
 - Final report (Sunday 12/8)



Proposal Presentation

- 5 minutes (about 5-8 slides)
- All team members should be involved in the presentation
- Will receive feedback from instructors and other students
 - 1-2 minutes between presentations reserved for written feedback
- Main presentation points
 - Research problem and motivation
 - Prior work
 - New research ideas
- You need to submit a copy of your slides (PDF or PPT)
 - Deadline: Friday 10/4 (on Gradescope)



Project Proposal Report

Part 1 (updated version of your pre-proposal)

Introduction:

- Describe and motivate the research problem
- Define in generic terms the main computational challenges

Experimental Setup:

- 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

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

New Research Ideas

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



Project Proposal Report

Part 3

Language Modality Exploration:

- Explore neural language models on your dataset (e.g., using Keras)
- 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., VGG-Net) and process your test images.
- Visualize the visual representations (using t-sne visualization) with overlaid class labels with different colors.





Unsupervised representation learning



Language Technologies Institute

Unsupervised learning

- We have access to $X = \{x_1, x_2, ..., x_n\}$ and not $Y = \{y_1, y_2, ..., y_n\}$
- Why would we want to tackle such a task
- 1. Extracting interesting information from data
 - Clustering
 - Discovering interesting trends
 - Data compression
- 2. Learn better representations





Unsupervised representation learning

- Force our representations to better model input distribution
 - Not just extracting features for classification
 - Asking the model to be good at representing the data and not overfitting to a particular task
 - Potentially allowing for better generalizability
- Use for initialization of supervised task, especially when we have a lot of unlabeled data and much less labeled examples



Restricted Boltzmann Machines



Language Technologies Institute

Restricted Boltzmann Machine (RBM)

- Undirected Graphical Model
- A generative rather than discriminative model
- Connections from every hidden unit to every visible one
- No connections across units (hence Restricted), makes it easier to train and do inference on



[Smolensky, Information Processing in Dynamical Systems: Foundations of Harmony Theory, 1986]



Restricted Boltzmann Machine (RBM)

$$p(\mathbf{x}, \mathbf{h}; \theta) = \frac{\exp(-E(\mathbf{x}, \mathbf{h}; \theta))}{\sum_{\mathbf{x}'} \sum_{\mathbf{h}'} \exp(-E(\mathbf{x}', \mathbf{h}'; \theta))} \longleftarrow \operatorname{Partition}_{\text{function } \mathbf{Z}}$$

• Hidden and visible layers are binary (e.g. $x = \{0, ..., 1, 0, 1\}$)

• Model parameters
$$\theta = \{W, b, a\}$$

$$E = -xWh - bx - ah$$

$$E = -\sum_{i}\sum_{j}w_{i,j}x_{i}h_{j} - \sum_{i}b_{i}x_{i} - \sum_{j}a_{j}h_{j}$$

Interaction Bias terms
term Visible

$$x_{1}$$
 x_{2} x_{n} Visible
layer



Boltzmann Machine

$$p(\mathbf{x}, \mathbf{h}; \theta) = \frac{\exp(-E(\mathbf{x}, \mathbf{h}; \theta))}{\sum_{\mathbf{x}'} \sum_{\mathbf{h}'} \exp(-E(\mathbf{x}', \mathbf{h}'; \theta))}$$

• Hidden and visible layers are binary (e.g. $x = \{0, ..., 1, 0, 1\}$)





Statistical Mechanics: Boltzmann Distribution

[also called Gibbs measure]

$$p(\boldsymbol{h};\theta) = \frac{\exp(-E(\boldsymbol{h};\theta)/kT)}{\sum_{\boldsymbol{h}'} \exp(-E(\boldsymbol{h}';\theta)/kT)}$$

probability distribution that gives the probability that a system will be in a certain state h

 $E(h; \theta)$: Energy of state h

- k: Boltzmann constant
- *T*: Thermodynamic temperature





RBM inference (have a trained θ)

- For inference
 - $p(h_j = 1 | \mathbf{x}; \theta) = \sigma(\sum_i x_i \mathbf{w}_{ij} + \mathbf{a}_j),$
 - $p(x_i = 1 | \mathbf{h}; \theta) = \sigma(\sum_j h_j w_{ij} + \mathbf{b}_i)$
 - derived from the joint probability definition
- Conditional inference is easy and of sigmoidal form
 - Given a trained model θ and an observed value x can easily infer h
 - Given a trained model θ and an hidden layer value h can easily infer x
- Need to sample as we get probabilities rather than values





C<mark>arnegie Mellon University</mark>

RBM training (learning the θ)

- Want to have a model that leads to good likelihood of training data
- First express the data likelihood (through marginal probability):

•
$$p(\mathbf{x};\theta) = \frac{\sum_{h} \exp(-E(\mathbf{x},h;\theta))}{Z}$$
 $Z = \sum_{x} \sum_{h} \exp(-E(\mathbf{x},h;\theta))$

- Want to optimize:
 - $\operatorname{argmin}_{\theta} \left[\sum_{t} \log \left(p(x^{(t)}; \theta) \right) \right]$, where *t* is a data sample
 - sum across all samples
 - minimizing negative log likelihood instead of maximizing the likelihood
- To Approximate computation of model term using Contrastive Divergence
 - Based on Markov Chain Monte Carlo (Gibbs) sampling

[G. Hinton, Training Products of Experts by Minimizing Contrastive Divergence, 2002]

See <u>http://www.iro.umontreal.ca/~lisa/twiki/bin/view.cgi/Public/DBNEquations</u> for more details



RBM extensions

- So far have only modeled binary input and hidden states
- Gaussian-Bernoulli RBM allows for real value modeling
 - Changes the inference and training only very slightly
 - Visible units are modeled as real values (under a Gaussian distribution), but hidden units are still binary
 - [Hinton and Salakhutdinov, Reducing the Dimensionality of Data with Neural Networks, 2006]
- Only requires a small change in some of the equations
- Can also introduce sparsity in hidden layers (sometimes helps)
 - [Lee et al., Sparse deep belief net model for visual area V2, 2007]



Examples of what the model learns



Learned W terms for each hidden unit





Deep Restricted Boltzmann Machines (DBMs)

- Can stack RBMs together to lead do deep versions of them
- The visible layer can be binary, Gaussian or Bernoulli
- Training fully end to end is very difficult
- Greedy layer-wise training
- Combine the RBMs layer by layer







Deep Belief Networks (DBN)

- To make it easier used Deep Belief Networks
 - Actually came before Deep RBMs
- Simplifies model training
- Turn the undirected model to directed one, making the interaction simpler



For more details see [Salakhutdinov and Hinton, Deep Boltzmann Machines, 2009]



Autoencoders



Language Technologies Institute



Autoencoders – an alternative to RBM

- What does auto mean?
 - Greek for self self encoding
- Feed forward network intended to reproduce the input
- Two parts encoder/decoder
 - x' = f(g(x)) -score function
 - g encoder
 - f decoder





Autoencoders

- Mostly follows Neural Network structure
 - Typically a matrix multiplication followed by a nonlinearity (e.g sigmoid)
- Activation will depend on type of x
 - Sigmoid for binary
 - Linear for real valued
- Often we use *tied weights* to force the sharing of weights in encoder/decoder

•
$$W^* = W^T$$

 word2vec is actually a bit similar to an autoencoder (except for the auto part)





Loss function

- Any differentiable similarity function
- Cross-entropy for binary x

•
$$L = -\sum_{k} (x_k \log(x'_k) + (1 - x_k) \log(1 - x'_k))$$

- Euclidean for real valued x
 - $L = \frac{1}{2} \sum_{k} (x_k x'_k)^2$
- Cosine similarity etc.
- Depends on the data being modeled





Learning

- To learn the model parameters (W*, W), we use back-propagation
- In case of Euclidean (with linear act) and Cross-entropy (with sigmoid act), we just have (x' - x) error to propagate
- If we're using *tied* weights, gradients need to be summed (like back propagation through time in RNN)
- Can use batch/stochastic gradient descent as before







Denoising autoencoder

- Simple idea
 - Add noise to input *x* but learn to reconstruct original
- Leads to a more robust representation and prevents copying
- Learns what the relationship is to represent a certain *x*
- Different noise added during each epoch





Autoencoder vs denoising autoencoder

MNIST data (as before)



Qualitatively denoising autoencoder leads to more meaningful features



- Can stack autoencoders as well
- Each encoding unit has a corresponding decoder
- As before, inference is feedforward, but now with more hidden layers







- Greedy layer-wise training
- Start with training first layer
 - Learn to encode x to h₁ and to decode x from h₁
 - Use backpropagation





- Greedy layer-wise training
- Start with training first layer
 - Learn to encode x to h₁ and to decode x from h₁
 - Use backpropagation
- Map from all x's to h₁'s
 - Discard decoder for now
- Train the second layer
 - Learn to encode *h*₁ to *h*₂ and to decode *h*₂ from *h*₁
 - Repeat for as many layers





- Greedy layer-wise training
- Start with training first layer
 - Learn to encode x to h₁ and to decode x from h₁
 - Use backpropagation
- Map from all x's to h₁'s
 - Discard decoder for now
- Train the second layer
 - Learn to encode *h*₁ to *h*₂ and to decode *h*₂ from *h*₁
 - Repeat for as many layers
- Reconstruct using previously learned decoders mappings
- Fine-tune the full network end-to-end





Stacked denoising autoencoders

- Can extend this to a denoising model
- Add noise when training each of the layers
 - Often with increasing amount of noise per layer
 - 0.1 for first, 0.2 for second,
 0.3 for third





Deep representations

- What can we do with them?
- Compression
 - Can work better than PCA
 - [Hinton and Salatkhudinov, Reducing the dimensionality of data with neural networks, 2006]





Deep representations

- What can we do with them?
- Compression
 - Can work better than PCA
 - [Hinton and Salatkhudinov, Reducing the dimensionality of data with neural networks, 2006]
- Discarding the decoder and using the middle layer as a representation
- Finetuning the autoencoder for a task





Multimodal representations



Language Technologies Institute

Multimodal representations

- What do we want from multi-modal representation
 - Similarity in that space implies similarity in corresponding *concepts*
 - Useful for various discriminative tasks – retrieval, mapping, fusion etc.
 - Possible to obtain in absence of one or more modalities
 - Fill in missing modalities given others (map between modalities)





Core Challenge: Multimodal Representation

Definition: Learning how to represent and summarize multimodal data in away that exploits the complementarity and redundancy.







Joint Multimodal Representation





Definition: Learning how to represent and summarize multimodal data in away that exploits the complementarity and redundancy.





Unsupervised Joint representations



Language Technologies Institute

Shallow multimodal representations

- Want deep multimodal representations
 - Shallow representations do not capture complex relationships
 - Often shared layer only maps to the shared section directly



Deep Multimodal autoencoders

- A deep representation learning approach
- A bimodal auto-encoder
 - Used for Audio-visual speech recognition

[Ngiam et al., Multimodal Deep Learning, 2011]

Deep Multimodal autoencoders - training

- Individual modalities can be pre-trained
 - Denoising Autoencoders
- To train the model to reconstruct the other modality
 - Use both
 - Remove audio

Deep Multimodal autoencoders - training

- Individual modalities can be pretrained
 - RBMs
 - Denoising Autoencoders
- To train the model to reconstruct the other modality
 - Use both
 - Remove audio
 - Remove video

Deep Multimodal autoencoders

- Can now discard the decoder and use it for the AVSR task
- Interesting experiment
 - "Hearing to see"

Deep Multimodal Boltzmann machines

- Generative model
- Individual modalities trained like a DBN
- Multimodal representation trained using Variational approaches
- Used for image tagging and crossmedia retrieval
- Reconstruction of one modality from another is a bit more "natural" than in autoencoder representation
- Can actually sample text and images

 [Srivastava and Salakhutdinov, Multimodal Learning with Deep Boltzmann Machines, 2012, 2014]

Deep Multimodal Boltzmann machines

- Pre-training on unlabeled data helps
- Can use generative models

Model	MAP	Prec@50
Random	0.124	0.124
SVM (Huiskes et al., 2010)	0.475	0.758
LDA (Huiskes et al., 2010)	0.492	0.754
DBM	0.526 ± 0.007	0.791 ± 0.008
DBM (using unlabelled data)	0.585 ± 0.004	$\textbf{0.836} \pm 0.004$

Image

kangarooisland, southaustralia. sa, australia, australiansealion, sand, ocean, 300mm

Given Tags

pentax, k10d,

<no text>

unseulpixel naturey crap fall, autumn, trees, leaves, foliage, forest, woods. branches, path

Generated Tags

beach, sea,

surf. strand.

shore, wave,

seascape,

waves night, lights, christmas,

nightshot,

woman,

people, faces,

girl, blackwhite, person, man

nacht. nuit.notte.

longexposure, noche, nocturna portrait, bw, blackandwhite,

flower, nature. areen, flowers, petal, petals, bud

nature, hill

clouds

scenery, green

blue, red, art, artwork, painted, paint, artistic surreal, gallery bleu

noiretblanc.

biancoenero

blancovnegro

bw, blackandwhite,

Code is available

http://www.cs.toronto.edu/~nitish/multimodal/

aheram, 0505 sarahc, moo

Deep Multimodal Boltzmann Machines

- Text information can help visual predictions!
 - Image retrieval task on MIR Flickr dataset

Model	MAP	Prec@50
Image LDA (Huiskes et al., 2010)	0.315	-
Image SVM (Huiskes et al., 2010)	0.375	-
Image DBN	0.463 ± 0.004	0.801 ± 0.005
Image DBM	0.469 ± 0.005	0.803 ± 0.005
Multimodal DBM (generated text)	$\textbf{0.531}\pm\textbf{0.005}$	$\textbf{0.832}\pm\textbf{0.004}$

Analyzing Intermediate Representations

Comparing deep multimodal representations

- Difference between them and the RBMs and the autoencoders
- Overall very similar behavior

Model	DBN	DAE	DBM
Logistic regression on joint layer features	0.599 ± 0.004	0.600 ± 0.004	0.609 ± 0.004
Sparsity + Logistic regression on joint layer features	0.626 ± 0.003	0.628 ± 0.004	0.631 ± 0.004
Sparsity + discriminative fine-tuning	0.630 ± 0.004	0.630 ± 0.003	0.634 ± 0.004
Sparsity + discriminative fine-tuning + dropout	0.638 ± 0.004	0.638 ± 0.004	$\textbf{0.641} \pm \textbf{0.004}$

Supervised Joint representations

Language Technologies Institute

Multimodal Joint Representation

- For supervised learning tasks
- Joining the unimodal representations:
 - Simple concatenation
 - Element-wise multiplication or summation
 - Multilayer perceptron
- How to explicitly model both unimodal and bimodal interactions?

e.g. Sentiment

Multimodal Sentiment Analysis

MOSI dataset (Zadeh et al, 2016)

- 2199 subjective video segments
- Sentiment intensity annotations
- 3 modalities: text, video, audio

Multimodal joint representation:

$$\boldsymbol{h}_{m} = \boldsymbol{f} \big(\boldsymbol{W} \cdot \big[\boldsymbol{h}_{x}, \boldsymbol{h}_{y}, \boldsymbol{h}_{z} \big] \big)$$

Unimodal, Bimodal and Trimodal Interactions

Bilinear Pooling

Models bimodal interactions:

 $h_m = h_x \otimes h_y = h_x \otimes h_y$

[Tenenbaum and Freeman, 2000]

Multimodal Tensor Fusion Network (TFN)

Models both unimodal and bimodal interactions:

$$h_{m} = \begin{bmatrix} h_{x} \\ 1 \end{bmatrix} \otimes \begin{bmatrix} h_{y} \\ 1 \end{bmatrix} = \begin{bmatrix} h_{x} \\ 1 \end{bmatrix} \begin{bmatrix} h_{x} \otimes h_{y} \\ h_{y} \end{bmatrix}$$
Important !

[Zadeh, Jones and Morency, EMNLP 2017]

Multimodal Tensor Fusion Network (TFN)

Can be extended to three modalities:

 $\boldsymbol{h}_{m} = \begin{bmatrix} \boldsymbol{h}_{x} \\ 1 \end{bmatrix} \otimes \begin{bmatrix} \boldsymbol{h}_{y} \\ 1 \end{bmatrix} \otimes \begin{bmatrix} \boldsymbol{h}_{z} \\ 1 \end{bmatrix}$

Explicitly models unimodal, bimodal and trimodal interactions !

[Zadeh, Jones and Morency, EMNLP 2017]

Language Technologies Institute

Experimental Results – MOSI Dataset

Multimodal Baseline	Bin	Binary		Regression	
	Acc(%)	F1	Acc(%)	MAE	r
Random	50.2	48.7	23.9	1.88	-
C-MKL	73.1	75.2	35.3	-	-
SAL-CNN	73.0	-	-	-	-
SVM-MD	71.6	72.3	32.0	1.10	0.53
RF	714	72.1	31.9	1 1 1	0.51
TFN	77.1	77.9	42.0	0.87	0.70
Human	85.7	87.5	53.9	0.71	0.82
Δ^{SOTA}	† 4.0	↑ 2.7	† 6.7	↓ 0.23	↑ 0.17

Improvement over State-Of-The-Art

Baseline	Binary		5-class	Regression	
	Acc(%)	F1	$\overline{\operatorname{Acc}(\%)}$	MAE	r
TFN _{language}	74.8	75.6	38.5	0.99	0.61
TFN _{visual}	66.8	70.4	30.4	1.13	0.48
$\mathrm{TFN}_{a coustic}$	65.1	67.3	27.5	1.23	0.36
$\mathrm{TFN}_{bimodal}$	75.2	76.0	39.6	0.92	0.65
$\mathrm{TFN}_{trimodal}$	74.5	75.0	38.9	0.93	0.65
$\mathrm{TFN}_{notrimodal}$	75.3	76.2	39.7	0.919	0.66
TFN	77.1	77.9	42.0	0.87	0.70
TFN_{early}	75.2	76.2	39.0	0.96	0.63

Language Technologies Institute

From Tensor Representation to Low-rank Fusion

1 Decomposition of weight tensor W

3 Rearranging computation

64

Multimodal Encoder-Decoder

- Visual modality often encoded using CNN
- Language modality will be decoded using LSTM
 - A simple multilayer perceptron will be used to translate from visual (CNN) to language (LSTM)

65

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

What should be the loss function?

[Frome et al., DeViSE: A Deep Visual-Semantic Embedding Model, NIPS 2013]

Max-Margin Loss – Multimodal Embeddings

Max-margin:

What should be the loss function?

[Frome et al., DeViSE: A Deep Visual-Semantic Embedding Model, NIPS 2013]

Structure-preserving Loss – Multimodal Embeddings

Symmetric max-margin:

[Wang et al., Learning Deep Structure-Preserving Image-Text Embeddings, CVPR 2016]

