





Lecture-07: Representation and Distributional Learning (Deep Learning & AI @LMU, Germany)

Speaker: Pankaj Gupta

PhD Student (Advisor: Prof. Hinrich Schütze) CIS, University of Munich (LMU) *Research Scientist* (NLP/Deep Learning), Machine Intelligence, Siemens AG | Nov 2018



Lecture Outline

- > Motivation: Distributed Feature Representations: What/Why?
- Strategies to obtain Feature Representation via Generative Models
 - Auto-encoders
 - Restricted Boltzmann Machines (RBMs) and Deep variants (DBM)
 - Neural Autoregressive Model: NADE, DocNADE, etc.
- > Language Modeling for Distributional Semantics:

(Tools: Word2vec, RNN-LM, RecvNN, textTOvec, etc.)

- > Metric Learning (e.g. Siamese Networks) for Textual Similarity
- > Compositionality: Recurrent vs Recursive (overview) (already covered in lecture 05)
- > Multi-task Learning (overview) (already covered in lecture 05)



What is Representation Learning ?

Why do we need Representation Learning ?

P How to form good Representations?

Intern © Siemens AG 2017

Seite 3 May 2017



Traditional Machine Learning

VISION



Intern © Siemens AG 2017

Seite 4 May 2017







Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Slide Credit: Marc'Aurelio Ranzato, Yann LeCun

Intern © Siemens AG 2017

Seite 5 May 2017



Exploiting compositionality gives an exponential gain in representation power-

(1) Distributed representations / embeddings

- \rightarrow feature learning
- (2) Deep architectures
 - → multiple levels of feature learning

Compositionality to describe the world around us efficiently



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Intern © Siemens AG 2017

Seite 6 May 2017







Intern © Siemens AG 2017

Seite 7 May 2017



Representation Learning (also known as feature learning)

Techniques to automatically discover the representations needed

- \rightarrow for feature detection that *explains the data*,
- \rightarrow classification from raw data

Feature

- \rightarrow measurable property or characteristic of a phenomenon being observed
- E.g.,
- shape, size, color, etc. of objects
- content, cuisine, color, etc. of food items
- syntactic and semantic relationships in words



Intern © Siemens AG 2017

Seite 8 May 2017



RAW DATA	Representation Learning:	Learning algorithm:	
e.g., images, vedios, text, etc.	Feature Disentangling	Performing tasks/actions	

Intern © Siemens AG 2017

Seite 9 May 2017







Collection of similar objects, based on shared [latent] features

Intern © Siemens AG 2017

Seite 10 May 2017



Word Analogy in distributional semantic vector space





Country-Capital

WWW.TENSORFLOW.ORG/IMAGES/LINEAR-RELATIONSHIPS.PNG

Intern © Siemens AG 2017

Seite 11 May 2017



Vector composition to determine semantic similarity in texts !!



HTTPS://WWW.SPRINGBOARD.COM/BLOG/INTRODUCTION-WORD-EMBEDDINGS/

Intern © Siemens AG 2017

Seite 12 May 2017



Vector composition to determine semantic similarity in phrases or sentences !!



HTTPS://WWW.SPRINGBOARD.COM/BLOG/INTRODUCTION-WORD-EMBEDDINGS/

Intern © Siemens AG 2017

Seite 13 May 2017



How about **CONTEXT** in feature representation learning?





How about **CONTEXT** in feature representation learning?



Intern © Siemens AG 2017

Seite 15 May 2017



How about **CONTEXT** in feature representation learning?

I am very <u>hungry</u>, I will <u>eat</u> HALWA

Intern © Siemens AG 2017

Seite 16 May 2017



How about **CONTEXT** in feature representation learning?

food

I am very <u>hungry</u>, I will <u>eat</u> HALWA

Intern © Siemens AG 2017

Seite 17 May 2017



How about **CONTEXT** in feature representation learning?

Is it a good idea to <u>eat HALWA after a meal</u>?

Intern © Siemens AG 2017

Seite 18 May 2017



How about **CONTEXT** in feature representation learning?

desert

Is it a good idea to <u>eat HALWA after a meal</u>?

Intern © Siemens AG 2017

Seite 19 May 2017



How about **CONTEXT** in feature representation learning?

desert

I put a cup of sugar too much in cooking HALWA?

Intern © Siemens AG 2017

Seite 20 May 2017



Unsupervised pre-training and Transfer learning



Loss function surfaces

Training deep networks can be challenging

- $\begin{array}{c} \rightarrow Non-convex \ loss \ function \\ \rightarrow \ in-appropriate \ initializers \end{array}$
- \rightarrow Initialisation is cricitial in Neural network training
- \rightarrow Using appropriate initializers for better convergence

Intern © Siemens AG 2017

Seite 21 May 2017



Unsupervised pre-training and Transfer learning

Transfer knowledge from previous learning:

- \rightarrow Representations or features
- \rightarrow Explanatory factors

↓ initialize hidden layers using UNSUPERVISED learning

Training deep networks can be challenging

 \checkmark

encode latent structure of input distribution in hidden layer

Previous learning from unlabeled data + labels for other tasks

 \rightarrow Prior. shared underlying explanatory factors, in particularly between P(x) and P(Y|x)

http://www.iro.umontreal.ca/~bengioy/ift6266/H16/representation-learning.pdf

Intern © Siemens AG 2017

Seite 22 May 2017



Unsupervised pre-training and Transfer learning

transfer knowledge from previous → Representations → Explanatory factors Example: Image recognition model Supervised Learning with unsupervised pre-training with unlabeled data to learn the representations of different levels of abstraction Transfer the knowledge Hu man

http://www.iro.umontreal.ca/~bengioy/ift6266/H16/representation-learning.pdf

Intern © Siemens AG 2017

Seite 23 May 2017



Motivation: Good Representation

- Good features for successful machine learning,
 - e.g., $man \leftrightarrow human$, $cat \leftrightarrow dog$, buy \leftrightarrow bought, buy $\leftarrow \Rightarrow$ acquired

buy $\leftarrow \rightarrow$ acquired Islam $\leftarrow \rightarrow$ Christianity, etc

Knowing features belief about objects in prior,

e.g., *features of car* → has_wheel, has_glasses, is_automobile, relatedto_manufacturer, etc.

Handcrafting features vs automatic feature learning

Representation learning: Estimate features / factors / causes that explains the data -> good representation i.e., good representation captures factors of variation that best explains the data

➤ Learning representations from representations → Representation learning e.g., autoencoders, RBMs, RSMs, NADE, DocNADE, iDocNADEe, generative RNNs, encoder-decoders, etc.

Intern © Siemens AG 2017

Seite 24 May 2017



Motivation: Distributed Representation Learning

Local Representations vs Distributed Representations?

Intern © Siemens AG 2017

Seite 25 May 2017



Consider a sequence **s** of words: "(a cat catches a mouse)"

A set of symbols is given by, **D** = {*mouse, cat, a, catches, (,)* }

Given a set of symbols **D**,

a local representation maps the i-th symbol in **D** to the i-th unit vector e_i of real values of n dimension, where n is the cardinality of **D**.

Hence, the i-th unit vector represents the i-th symbol.

```
\begin{aligned} \text{mouse} &\to \mathbf{e}_1 = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}^T \\ \text{cat} &\to \mathbf{e}_2 = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}^T \\ \text{a} &\to \mathbf{e}_3 = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 \end{pmatrix}^T \\ \text{catches} &\to \mathbf{e}_4 = \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}^T \\ \begin{pmatrix} &\to \mathbf{e}_5 = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix}^T \\ \end{pmatrix} &\to \mathbf{e}_6 = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}^T \end{aligned}
```

Seite 26 May 2017



Consider a sequence **s** of words: "(a cat catches a mouse)"

A set of symbols is given by, **D** = {*mouse, cat, a, catches, (,)* }

a sequence of vectors



 \rightarrow a sequence of vectors representing the symbols in the sequence

 \rightarrow used in recurrent neural networks

a bag-of-symbols

 \rightarrow a sequence is represented with one vector generally obtained with a weighted sum of vectors representing symbols, i.e., orderless

 \rightarrow SVM (used in Information Retrieval task)



Local Representations

 $mouse \rightarrow \mathbf{e}_1 = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \end{pmatrix}^T$

 $\operatorname{catches} \to \mathbf{e}_4 = \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}^T$

 $\operatorname{cat} \to \mathbf{e}_2 = (0 \ 1 \ 0 \ 0 \ 0 \ 0)^T$

 $a \to e_3 = (0 \ 0 \ 1 \ 0 \ 0 \ 0)^T$

 $(\rightarrow \mathbf{e}_5 = (0 \ 0 \ 0 \ 0 \ 1 \ 0)^T$

 $) \rightarrow \mathbf{e}_6 = (0 \ 0 \ 0 \ 0 \ 0 \ 1)^T$









https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/slides/L1_intro.pdf Intern © Siemens AG 2017









Seite 31 May 2017



What is Distributed Representation Learning?

Distributed: "information is encoded by a multiplicity of features / factors / causes"

Distributed representations:

- \rightarrow vectors or tensors in metric spaces
- → transformations of the data that *compactly* capture many different factors of variations
- \rightarrow underlying learning models are neural networks
- E.g. Distributed word representations:

Each word is represented as a dense and real-valued vector in low dimensional space,

and each latent feature encodes syntactic and semantic relatedness information

 \rightarrow addresses the Curse of Dimensionality



Need for Distributed Representation Learning

The need for distributed representations

- Factor models, PCA, RBMs, Neural Nets, Sparse Coding, Deep Learning, etc.
- Each parameter influences many regions, not just local neighbors
- # of distinguishable regions grows almost exponentially with # of parameters
- GENERALIZE NON-LOCALLY
 TO NEVER-SEEN REGIONS







(b)

no pattern

Motivation: Distributed Representation Learning

 \rightarrow dense vectors in distributed representations

→ one concept represented by the dense vector
→ one dimension per property

 \rightarrow enable to share similarity between more than two concepts





Motivation: Distributed Representation Learning

Power of Distributed Representations



Distributed Representation by Hinton (1984)

Intern © Siemens AG 2017

Seite 35 May 2017



Representations Learning



→ Can we interpret each dimension/property?
 → Lack of interpretability in dense representations
 learned by Deep Learning
 Leard to track down what's failing

 \rightarrow Hierarchical composition

 \rightarrow Deep Learning is very good !!!

 \rightarrow Hard to track down what's failing

Intern © Siemens AG 2017

Seite 36 May 2017


Distributed Representation Learning in Deep Networks

Distributed Representations in Deep Leaning !!! Yes it works, but how?

Representation learning: Attempts to *automatically learn good features or representations* Deep Learning: Attempts to learn *multiple levels of representations* of increasing abstraction

Intern © Siemens AG 2017

Seite 37 May 2017



Distributed Representation Learning in Neural Networks

- → Key concept in neural networks: *Distributed Representation Learning*
- \rightarrow Key questions:
- How can a neural network be so effective representing objects when it has only a few hidden units (i.e. much fewer units than possible objects)?
- What is each hidden unit actually representing?
- How can a neural network generalize to objects that is has never seen before?



Distributed Representations in Deep Learning

Deep Learning = Hierarchical Compositionality

- Cascade of non-linear transformations
- Multiple layers of representations



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



Distributed Representations in Deep Learning

Deep Learning = Hierarchical Compositionality

- Cascade of non-linear transformations
- Multiple layers of representations

 \rightarrow No single neuron "encodes" everything

 \rightarrow Groups of neurons work together





Supervised Learning

Data: (x, y) x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



Classification

https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/

Intern © Siemens AG 2017

Seite 41 May 2017



Unsupervised Learning

Data: x Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



K-means clustering

https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/

Intern © Siemens AG 2017

Seite 42 May 2017



Unsupervised Learning

Data: x Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc. original data space component space PCA PC 1 PC 2 N Ч С PC 1 3-d 2-d

> Principal Component Analysis (Dimensionality reduction)

https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/

Intern © Siemens AG 2017

Seite 43 May 2017



Unsupervised Learning

Data: x Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



1-d density estimation





https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/

Intern © Siemens AG 2017

Seite 44 May 2017



Supervised Learning

Data: (x, y) x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.

Unsupervised Learning

Data: x Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.

https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/

Intern © Siemens AG 2017

Seite 45 May 2017



Generative Classification

- \rightarrow Model **p(x, y)**; estimate p(x|y) and p(y)
- \rightarrow Use Bayes Rule to predict y, e.g. *Naïve Bayes*

Discriminative Classification

→ Estimate **p(y|x)** directly, e.g. *Logistic Regression, CNN, RNN, etc.*



→ Model p(x), e.g. RBMs, VAEs, NADE, RSM, DocNADE, etc.
This lecture



Given training data, generate new samples from same distribution



Several flavors:

→ Explicit density estimation: explicitly define and solve for $p_{model}(x)$ →Implicit density estimation: learn model that can sample from $p_{model}(x)$ w/o explicitly defining it

Training generative models can enable inference of latent representations, used as general features

Intern © Siemens AG 2017

Seite 47 May 2017



Given training data, generate new samples from same distribution

Unsupervised Representation Learning to exploits tons to unlabeled data

Training data ~ $p_{data}(x)$

Generated samples ~ $p_{model}(x)$

Want to learn $p_{model}(x)$ similar to $p_{data}(x)$

Density estimation: a core problem in unsupervised learning

Several flavors:

 \rightarrow Explicit density estimation: explicitly define and solve for $p_{model}(x)$

 \rightarrow Implicit density estimation: learn model that can sample from $p_{model}(x)$ w/o explicitly defining it

Training generative models can enable inference of latent representations, used as general features

Intern © Siemens AG 2017

Seite 48 May 2017



Why Latent Factors & Unsupervised Representation Learning? Because of *Causality*.

On causal and anticausal learning, (Janzing et al ICML 2012)

- If Ys of interest are among the causal factors of X, then $P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$

is tied to P(X) and P(X|Y), and P(X) is defined in terms of P(X|Y), i.e.

- The best possible model of X (unsupervised learning) MUST involve Y as a latent factor, implicitly or explicitly.
- Representation learning SEEKS the latent variables H that explain the variations of X, making it likely to also uncover Y.

http://www.iro.umontreal.ca/~bengioy/ift6266/H16/representation-learning.pdf Intern © Siemens AG 2017

Intern © Siemens AG 2017



Manifolds

TO DO: Motivation: image manifold

 \rightarrow probability mass concentrates near regions that have a much smaller dimensionality than the original space where the data lives



http://www.iro.umontreal.ca/~bengioy/ift6266/H16/representation-learning.pdf

Intern © Siemens AG 2017

Seite 50 May 2017



Reconstruction Objective

→ train a network to learn weights such as the network can *reconstruct* the input, given the output (e.g., hidden vector encoding the input)!!!





Intern © Siemens AG 2017

Seite 52 May 2017

Corporate Technology



Learning Distributed Word Representations (i.e., word embeddings)

Intern © Siemens AG 2017

Seite 53 May 2017

Corporate Technology

Distributional Language Semantics: Tools for Language Modeling

Distributed word representation

- → represent (embed) words in a continuous vector space where semantically similar words are mapped to nearby points
- \rightarrow describe meaning of words and sentences with vectorial representations

Idea1: "you shall judge a word by the company it keeps"

Idea2: Distributional hypothesis: "words have similar meaning if used in similar contexts "



- \rightarrow words sharing similar attributes are similar
- E.g., *dog* is more similar to *cat* than to *car* as *dog* and *cat* share more attributes than *dog* and *car*

 s_2

 s_3

1 0 /

Local representation: BoW

 s_1

 \rightarrow word-to-word matrices obtained by observing n-word windows of target words

eats

s_1	a cat catches a mouse	
s_2	$a \ dog \ eats \ a \ mouse$	X =
s_3	$a \ dog \ catches \ a \ cat$	

Small corpus of 3 sentence/documents

a	$\binom{2}{2}$	2	2	a	0	1
cat	1	0	1	cat	2	0
dog	0	1	1	dog	2	0
mouse	1	1	0	$=$ $_{mouse}$	2	0
catches	1	0	1	catches	2	1
	4				•	

co-occurrence matrix

2

0

0

0

a cat dog mouse catches eats

2

0

0

0

0

0

2

0

eats $\setminus 1 \ 0$



WWW.TENSORFLOW.ORG/IMAGES/LINEAR-RELATIONSHIPS.PNG

Intern © Siemens AG 2017

Seite 55 May 2017 $\mathbf{2}$

0

1

0

0

0



co-occurrence matrix

a cat dog mouse catches eats

Why: Distributed Word Representations

 \rightarrow words sharing similar attributes are similar.

E.g., dog is more similar to cat than to car as dog and cat share more attributes than dog and car.

 \rightarrow word-to-word matrices obtained by observing n-word windows of target words

 $\begin{array}{cccc} s_1 & a & cat \ catches \ a \ mouse \\ s_2 & a \ dog \ eats \ a \ mouse \\ s_3 & a \ dog \ catches \ a \ cat \\ \end{array} X =$

Small corpus of 3 sentence/documents

a	(2	2	2	\	a	/ 0	1	2°	2	2	2
cat		1	0	1		cat	2	0	$\overline{0}$	$\overline{0}$	1	0
dog		0	1	1		dog	2	0	0	0	1	1
mouse		1	1	0	=	mouse	2	0	0	0	0	0
catches		1	0	1		catches	2	1	1	0	0	0
eats		0	1	0 /	/	eats	$\setminus 1$	0	1	0	0	0

Local representation: BoW

 $s_1 \ s_2 \ s_3$

Distributed representation: a word-to-word matrix considering a 1-word window of context

WWW.TENSORFLOW.ORG/IMAGES/LINEAR-RELATIONSHIPS.PNG

Intern © Siemens AG 2017

Seite 56 May 2017



 \rightarrow words sharing similar attributes are similar.

E.g., dog is more similar to cat than to car as dog and cat share more attributes than dog and car.

 \rightarrow word-to-word matrices obtained by observing n-word windows of target words

 $\begin{array}{cccc} s_1 & a & cat \ catches \ a \ mouse \\ s_2 & a & dog \ eats \ a \ mouse \\ s_3 & a & dog \ catches \ a \ cat \end{array} \quad X =$

Small corpus of 3 sentence/documents

a	$\binom{2}{2}$	2	2		a
cat	1	0	1		cat
dog	0	1	1		dog
mouse	1	1	0	=	mouse
catches	1	0	1		catches
eats	0	1	0	/	eats

 s_3

 s_2

 s_1

Local representation: BoW

co-occurrence matrix

a cat dog mouse catches eats

0

 $\mathbf{2}$

 $\mathbf{2}$

2

 $\mathbf{2}$

0

Distributed representation: a word-to-word matrix considering a 1-word window of context

WWW.TENSORFLOW.ORG/IMAGES/LINEAR-RELATIONSHIPS.PNG

Intern © Siemens AG 2017

Seite 57 May 2017

0

 $\mathbf{2}$

0

1

0

0

0



 \rightarrow words sharing similar attributes are similar.

E.g., dog is more similar to cat than to car as dog and cat share more attributes than dog and car.

 \rightarrow word-to-word matrices obtained by observing n-word windows of target words

$$\begin{array}{cccc} s_1 & a \ cat \ catches \ a \ mouse \\ s_2 & a \ dog \ eats \ a \ mouse \\ s_3 & a \ dog \ catches \ a \ cat \end{array} X$$

Small corpus of 3 sentence/documents

$$\begin{array}{c} a \\ cat \\ dog \\ mouse \\ catches \\ eats \end{array} \begin{pmatrix} 0 & 1 & 2 & 2 & 2 & 2 \\ 2 & 0 & 0 & 0 & 1 & 0 \\ 2 & 0 & 0 & 0 & 1 & 1 \\ 2 & 0 & 0 & 0 & 0 & 1 \\ 2 & 0 & 0 & 0 & 0 & 0 \\ 2 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ \end{array}$$

a cat dog mouse catches eats

co-occurrence matrix

similar in distributed representation space

Distributed representation: a word-to-word matrix considering a 1-word window of context

WWW.TENSORFLOW.ORG/IMAGES/LINEAR-RELATIONSHIPS.PNG

Intern © Siemens AG 2017

Seite 58 May 2017



However, original co-occurrence matrix is very costly to obtain and store

Need compact distributed vectors

 $\begin{array}{cccc} s_1 & a \ cat \ catches \ a \ mouse \\ s_2 & a \ dog \ eats \ a \ mouse \\ s_3 & a \ dog \ catches \ a \ cat \end{array} X =$

Small corpus of 3 sentence/documents

$$\begin{array}{c} a \\ cat \\ dog \\ mouse \\ catches \\ eats \end{array} \begin{pmatrix} 0 & 1 & 2 & 2 & 2 & 2 \\ 2 & 0 & 0 & 0 & 1 & 0 \\ 2 & 0 & 0 & 0 & 1 & 1 \\ 2 & 0 & 0 & 0 & 0 & 1 \\ 2 & 0 & 0 & 0 & 0 & 0 \\ 2 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ \end{array} \right)$$

a cat dog mouse catches cate

Distributed representation: a word-to-word matrix considering a 1-word window of context

WWW.TENSORFLOW.ORG/IMAGES/LINEAR-RELATIONSHIPS.PNG

Intern © Siemens AG 2017

 \rightarrow

E.g

 \rightarrow

Seite 59 May 2017

SIEMENS

Learning Compact Distributed Word Representations

Word2vec: Tool to generate Word embeddings (i.e., distributed word representation) using large corpus

 \rightarrow represent (embed) words in a

continuous vector space

where semantically similar words are mapped to nearby points

→ uses contextual information to learn word vectors

→ neural network predicts a target word from the words surrounding it (context)

 \rightarrow no explicitly co-occurrence matrix

- \rightarrow no explicit association between word pairs
- \rightarrow distribution of the words learned implicitly
- \rightarrow compact distributed representation

 \rightarrow dimensions of vectors not interpretable

Intern © Siemens AG 2017





Why: Power of Distributed Word Representation



Further reading: doc2vec, LDA, LSA, etc.

Intern © Siemens AG 2017

Seite 61 May 2017



Family of Generative Models

Strategies to obtaining Distributed Representations

Approximate Density

e.g.,
→ Variational Autoencoders,
→ RBMs, DBN, DBM,
→ RSM

 \rightarrow RNN-RSM

Tractable Density

↓ e.g., → NADE, MADE → DocNADE → iDocNADEe

→ PixelRNN/CNN





Intern © Siemens AG 2017

Seite 63 May 2017



→ a feed-forward neural network trained to reproduce its input at the output layer



Intern © Siemens AG 2017





Intern © Siemens AG 2017

Seite 66 May 2017





\rightarrow Difficult to interpret

Intern © Siemens AG 2017

Seite 67 May 2017

Benefits of Auto-encoders (AEs)

 \rightarrow a feed-forward neural network trained to **reproduce its input** at the output layer **Key Facts about AEs**:

- \rightarrow unsupervised ML algorithm, similar to PCA
- \rightarrow neural network's target output is its input
- → learn latent features/encoding of input (no manual feature engineering)
- → represent both linear and non-linear transformation in encoding
- \rightarrow layered to form deep learning network, i.e., distributed representations
- \rightarrow tractable / easier optimization
- → applications in denoising and dimensionality reduction (dense representation)
- \rightarrow powerful non-linear (i.e., non-linear encoding and decoding) generalization of PCA







Limitations of (regular) Autoencoders

 \rightarrow Need to Conceptualize, not only compress.....

Different variants to rescue.....!!!

Intern © Siemens AG 2017

Seite 69 May 2017

Intern © Siemens AG 2017

Auto-encoders (AEs) variants

Autoencoder variants:

- 1. Denoising Autoencoders
- 2. Sparse Autoencoders
- 3. Convolutional Autoencoders
- 4. Variational Autoencoders (VAE)
- 5. Contractive Autoencoders (CAE)
- 6. Stacked Autoencoders, etc...

.....details later in the lecture series (Lecture on "Generative Models")



Idea: Constraint the reconstruction of an autoencoder



Family of Neural Generative Models

Intern © Siemens AG 2017

Seite 71 May 2017

Corporate Technology



Family of Generative Models

Directed graphical models

- define prior over top-most latent representation
- define conditionals from top latent representation to observation

 $p(\mathbf{x}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = p(\mathbf{x} | \mathbf{h}^{(1)}) p(\mathbf{h}^{(1)} | \mathbf{h}^{(2)}) p(\mathbf{h}^{(2)} | \mathbf{h}^{(3)}) p(\mathbf{h}^{(3)})$

 examples: variational autoencoders (VAE), generative adversarial networks (GAN), sparse coding, helmholtz machines



Properties

pros: easy to sample from (ancestral sampling)

• cons: $p(\mathbf{x})$ is intractable, so hard to train

https://ift6135h18.wordpress.com/author/aaroncourville/page/1/

Intern © Siemens AG 2017

Seite 72 May 2017

(details in Lecture on "Generative Models")



Family of Generative Models

Undirected graphical models

define a joint energy function

 $E(\mathbf{x}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = -\mathbf{x}\mathbf{W}^{(1)}\mathbf{h}^{(1)} - \mathbf{h}^{(2)}\mathbf{W}^{(2)}\mathbf{h}^{(3)} - \mathbf{h}^{(3)}\mathbf{W}^{(3)}\mathbf{h}^{(4)}$

exponentiate and normalize

$$p(\mathbf{x}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = \exp\left(-E(\mathbf{x}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)})\right)/Z$$

- examples: deep Boltzmann machines (DBM), deep energy models
- Properties
 - pros: can compute $p(\mathbf{x})$ up to a multiplicative factor (true for RBMs not general BMs)
 - cons: hard to sample from (MCMC), $p(\mathbf{x})$ is intractable, so hard to train

https://ift6135h18.wordpress.com/author/aaroncourville/page/1/

Intern © Siemens AG 2017






Family of Generative Models

Autoregressive generative models

- choose an ordering of the dimensions in x
- define the conditionals in the product rule expression of $p(\mathbf{x})$

$$p(\mathbf{x}) = \prod_{k=1}^{D} p(x_k | \mathbf{x}_{< k})$$

 examples: masked autoencoder distribution estimator (MADE), pixelCNN neural autoregressive distribution estimator (NADE) spatial LSTM, pixelRNN

Properties

pros: $p(\mathbf{x})$ is tractable, so easy to train, easy to sample (though slower)

cons: doesn't have a natural latent representation

https://ift6135h18.wordpress.com/author/aaroncourville/page/1/

Intern © Siemens AG 2017





Undirected (Generative) Probabilistic Graphical Models

Intern © Siemens AG 2017

Seite 75 May 2017

Corporate Technology



- → undirected *probabilistic* graphical model
- → unsupervised *stochastic extractor* of **binary** features (**h**)
- \rightarrow trained using *reconstruction* objective
- \rightarrow transform data into latent feature space and then reconstruct to learn data distribution



Further reading: <u>https://ift6266h15.files.wordpress.com/2015/03/chapter21.pdf</u>

Intern © Siemens AG 2017

Seite 76 May 2017



 \rightarrow Two layers: binary visible (v) and binary hidden (h) layer

\rightarrow energy-based models,

Therefore, joint probability distribution is given by its energy function:

 $P(\boldsymbol{v}, \boldsymbol{h}) = \frac{1}{Z} \exp \left\{-E(\boldsymbol{v}, \boldsymbol{h})\right\}.$ partition function energy binary visible binary hidden features \rightarrow energy function that parameterizes the relationship between the visible
and hidden variables

$$E(\boldsymbol{v},\boldsymbol{h}) = -\boldsymbol{b}^{\top}\boldsymbol{v} - \boldsymbol{c}^{\top}\boldsymbol{h} - \boldsymbol{v}^{\top}\boldsymbol{W}\boldsymbol{h}$$

 \rightarrow normalizing constant known as the **partition function**

$$Z = \sum_{\boldsymbol{v}} \sum_{\boldsymbol{h}} \exp\left\{-E(\boldsymbol{v}, \boldsymbol{h})\right\}$$

 $(h_1)(h_2)(h_3)(h_4)$ sible $(v_1)(v_2)(v_3)$

E.g., a document of 2000 unique words,

 $2^{\dim(v)} 2^{\dim(h)} \rightarrow 2^{2000} 2^{50} \rightarrow \text{intractable}$

summing over all states > summing exponential number of terms > computationally intractable P(v,h)

Intern © Siemens AG 2017

Seite 77 May 2017

Corporate Technology



\rightarrow energy-based models,

Therefore, *joint probability distribution* is given by its energy function: $P(v, h) = \frac{1}{Z} \exp \{-E(v, h)\}$.

 \rightarrow energy function, $E(\boldsymbol{v}, \boldsymbol{h}) = -\boldsymbol{b}^{\top} \boldsymbol{v} - \boldsymbol{c}^{\top} \boldsymbol{h} - \boldsymbol{v}^{\top} \boldsymbol{W} \boldsymbol{h}$

 \rightarrow partition function

$$Z = \sum_{\boldsymbol{v}} \sum_{\boldsymbol{h}} \exp\left\{-E(\boldsymbol{v}, \boldsymbol{h})\right\}$$

→ conditional distributions from the joint distribution

Full conditional over the **hidden layer** as the factorial distribution:

$$p(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^{n} p(h_j|\mathbf{v}) = \prod_{j=1}^{n} \operatorname{sigmoid}(c_j + \mathbf{v}^T W_{:,j})$$

Full conditional over the **visible layer** as the factorial distribution:

$$p(\mathbf{v}|\mathbf{h}) = \prod_{i=1}^{d} p(v_i|\mathbf{h}) = \prod_{i=1}^{d} \operatorname{sigmoid}(b_i + W_{i,:}\mathbf{h})$$

Intern © Siemens AG 2017

Seite 78 May 2017



\rightarrow energy-based models,

Therefore, *joint probability distribution* is given by its energy function: $P(v, h) = \frac{1}{Z} \exp \{-E(v, h)\}$.

 \rightarrow energy function, $E(\boldsymbol{v}, \boldsymbol{h}) = -\boldsymbol{b}^{\top} \boldsymbol{v} - \boldsymbol{c}^{\top} \boldsymbol{h} - \boldsymbol{v}^{\top} \boldsymbol{W} \boldsymbol{h}$

 \rightarrow partition function

$$Z = \sum_{\boldsymbol{v}} \sum_{\boldsymbol{h}} \exp\left\{-E(\boldsymbol{v}, \boldsymbol{h})\right\}$$

 \rightarrow conditional distributions from the joint distribution

Full conditional over the hidden layer as the factorial distribution:

$$p(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^{n} p(h_j|\mathbf{v}) = \prod_{j=1}^{n} \operatorname{sigmoid}(c_j + \mathbf{v}^T W_{:,j})$$

encoding

 $P(h_i =$

$$(h_1 \ h_2 \ h_3 \ h_4)$$

$$(j) \ v_1 \ v_2 \ v_3$$

$$j = 1$$

$$(c_j + v^\top W_{:,j})$$

Seite 79 May 2017



\rightarrow energy-based models,

Therefore, *joint probability distribution* is given by its energy function: $P(v, h) = \frac{1}{Z} \exp \{-E(v, h)\}$.

 \rightarrow energy function, $E(\boldsymbol{v}, \boldsymbol{h}) = -\boldsymbol{b}^{\top} \boldsymbol{v} - \boldsymbol{c}^{\top} \boldsymbol{h} - \boldsymbol{v}^{\top} \boldsymbol{W} \boldsymbol{h}$

 \rightarrow partition function

$$Z = \sum_{\boldsymbol{v}} \sum_{\boldsymbol{h}} \exp\left\{-E(\boldsymbol{v}, \boldsymbol{h})\right\}$$

 \rightarrow conditional distributions from the joint distribution

Full conditional over the hidden layer as the factorial distribution:

$$p(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^{n} p(h_j|\mathbf{v}) = \prod_{j=1}^{n} \operatorname{sigmoid}(c_j + \mathbf{v}^T W_{:,j})$$

encoding

 $P(h_i)$

$$\begin{array}{c|c} h_1 & h_2 & h_3 & h_4 \\ \hline W_{:,j} \end{pmatrix} & & & & & & \\ \hline \mathbf{y}_{:,j} \end{pmatrix} \\ \hline \mathbf{y}_1 & \mathbf{y}_2 & \mathbf{y}_3 \\ \hline \mathbf{y}_1 & \mathbf{y}_1 & \mathbf{y}_2$$

Seite 80 May 2017



\rightarrow energy-based models,

Therefore, *joint probability distribution* is given by its energy function:

- \rightarrow energy function, $E(\boldsymbol{v}, \boldsymbol{h}) = -\boldsymbol{b}^{\top} \boldsymbol{v} \boldsymbol{c}^{\top} \boldsymbol{h} \boldsymbol{v}^{\top} \boldsymbol{W} \boldsymbol{h}$
- \rightarrow partition function

$$Z = \sum_{\boldsymbol{v}} \sum_{\boldsymbol{h}} \exp\left\{-E(\boldsymbol{v}, \boldsymbol{h})\right\}$$

 \rightarrow conditional distributions from the joint distribution

Full conditional over the hidden layer as the factorial distribution:

$$p(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^{n} p(h_j|\mathbf{v}) = \prod_{j=1}^{n} \operatorname{sigmoid}(c_j + \mathbf{v}^T W_{:,j})$$

encoding

Full conditional over the visible layer as the factorial distribution:

$$p(\mathbf{v}|\mathbf{h}) = \prod_{i=1}^{d} p(v_i|\mathbf{h}) = \prod_{i=1}^{d} \operatorname{sigmoid}(b_i + W_{i,:}\mathbf{h})$$

Decoding/reconstruction

Intern © Siemens AG 2017

Seite 81 May 2017

n:
$$P(v, h) = \frac{1}{Z} \exp \{-E(v, h)\}.$$

Same W matrix in encoding and decoding





Restricted Boltzmann Machines (RBMs): Illustration



https://jamesmccaffrey.wordpress.com/2017/06/02/restricted-boltzmann-machines-using-c/ Intern © Siemens AG 2017

Seite 82 May 2017



Training RBMs

Cost: maximize log-likelihood of the data, v

Let us consider that we have a batch (or minibatch) of n examples taken from an i.i.d dataset (independently and identically distributed examples) $\{v^{(1)}, \ldots, v^{(t)}, \ldots, v^{(n)}\}$. The log likelihood under the RBM with parameters \boldsymbol{b} (visible unit biases), \boldsymbol{c} (hidden unit biases) and \boldsymbol{W} (interaction weights) is given by:



Intern © Siemens AG 2017

https://ift6266h15.files.wordpress.com/2015/03/chapter21.pdf

Seite 83 May 2017



Training RBMs

Cost: maximize log-likelihood of the data, v

Let us consider that we have a batch (or minibatch) of n examples taken from an i.i.d dataset (independently and identically distributed examples) $\{v^{(1)}, \ldots, v^{(t)}, \ldots, v^{(n)}\}$. The log likelihood under the RBM with parameters b (visible unit biases), c (hidden unit biases) and W (interaction weights) is given by:



Corporate Technology



Stacked RBMs: Deep Boltzmann Machine (DBM)



Detailed lecture: https://www.youtube.com/watch?v=MnGXXDjGNd0

Intern © Siemens AG 2017

Seite 85 May 2017



Applications of RBMs

- \rightarrow Unsupervised pre-training and transfer learning
- \rightarrow Once trained, can use W and **biases** as initial values for a neural net!



Filter (W) from 1st layer: "pen-strokes"

\rightarrow W can be used to initialize neural networks

 \rightarrow latent vector as features into neural network

Intern © Siemens AG 2017

Seite 86 May 2017



Applications of DBMs

- \rightarrow Unsupervised pre-training and transfer learning
- → Once trained, can use W and biases as initial values for a neural net!



 \rightarrow W can be used to initialize neural networks

 \rightarrow latent vector as features into neural network



RBM variants: How to model Word Counts?

→ RBM and DBM model **binary** input or real-valued input using Gaussian-RBMs

How to model count data? Example: Text document i.e., word counts.

'a cat catches a mouse"
$$\mathbf{s} = \begin{pmatrix} 1\\1\\2\\1\\0\\0 \end{pmatrix}$$

Seite 88 May 2017



RBM variants: Overview of Replicated Softmax (RSM)



Gupta et al. 2018, Deep Temporal-Recurrent-Replicated-Softmax for Topical Trends over Time. Intern © Siemens AG 2017

Seite 89 May 2017



RBM variants: Overview of Replicated Softmax (RSM)



Gupta et al. 2018, Deep Temporal-Recurrent-Replicated-Softmax for Topical Trends over Time. Intern © Siemens AG 2017

Seite 90 May 2017



Idea: Tractable log-likelihood, p(v)

Intern © Siemens AG 2017

Seite 91 May 2017

Corporate Technology



- > **NADE:** Neural Autoregressive Distributional Estimator
- > Inspired from **RBM**,
- Generative model over binary observations v
- sampling each dimension one after another

given previous variables



Uria, Benigno, et al. "Neural Autoregressive Distribution Estimation"

Intern © Siemens AG 2017

Seite 92 May 2017



- > **NADE:** Neural Autoregressive Distributional Estimator
- Inspired from RBM Family therefore, energy based
- > sampling each dimension one after another

 $p(\mathbf{v}) = \prod_{i=1}^{D} p(v_i | \mathbf{v}_{< i})$ and computes all $p(v_i | \mathbf{v}_{< i})$ using the feed-forward architecture

 $\mathbf{h}_i(\mathbf{v}_{< i}) = \operatorname{sigm}\left(\mathbf{c} + \mathbf{W}_{:,<i}\mathbf{v}_{< i}\right)$

NADE is for binary data, $p(v_i = 1 | \mathbf{v}_{< i}) = \operatorname{sigm} (b_i + \mathbf{V}_{i,:} \mathbf{h}_i(\mathbf{v}_{< i}))$

$$\mathbf{h}_{i+1}(\mathbf{v}_{< i+1}) = \operatorname{sigm}(\mathbf{c} + \sum_{k < i+1} W_{:,v_k}) = \operatorname{sigm}(\mathbf{W}_{:,v_i} + \mathbf{c} + \sum_{k < i} \mathbf{W}_{:,v_k})$$



for $i \in \{1, ..., D\}$, where sigm $(x) = 1/(1 + \exp(-x))$, $\mathbf{W} \in \mathbb{R}^{H \times D}$ and $\mathbf{V} \in \mathbb{R}^{D \times H}$ are connection parameter matrices, $\mathbf{b} \in \mathbb{R}^{D}$ and $\mathbf{c} \in \mathbb{R}^{H}$ are bias parameter vectors, $\mathbf{v}_{< i}$ is the subvector $[v_1, ..., v_{i-1}]^{\top}$ and $\mathbf{W}_{:,<i}$ is a matrix made of the i - 1 first columns of \mathbf{W} .

corresponds to a neural network with several parallel $h_i(\mathbf{v}_{< i})$ hidden layers

Advantages: Tractable $p(v) \rightarrow$ easier to train

Intern © Siemens AG 2017

Seite 93 May 2017



- > **NADE:** Neural Autoregressive Distributional Estimator
- Inspired from RBM Family therefore, energy based
- ➤ sampling each dimension one after another

 $p(\mathbf{v}) = \prod_{i=1}^{D} p(v_i | \mathbf{v}_{< i})$ and computes all $p(v_i | \mathbf{v}_{< i})$ using the feed-forward architecture

 $\mathbf{h}_i(\mathbf{v}_{< i}) = \operatorname{sigm}\left(\mathbf{c} + \mathbf{W}_{:, < i} \mathbf{v}_{< i}\right)$

NADE is for binary data, $p(v_i = 1 | \mathbf{v}_{< i}) = \operatorname{sigm} (b_i + \mathbf{V}_{i,:} \mathbf{h}_i(\mathbf{v}_{< i}))$

$$\mathbf{h}_{i+1}(\mathbf{v}_{< i+1}) = \operatorname{sigm}(\mathbf{c} + \sum_{k < i+1} W_{:,v_k}) = \operatorname{sigm}(\mathbf{W}_{:,v_i} + \mathbf{c} + \sum_{k < i} \mathbf{W}_{:,v_k})$$

 $\frac{re}{W} = \frac{11}{W} = \frac{12}{W} = \frac{13}{W} = \frac{14}{W} = \frac{14}{W}$

V

for $i \in \{1, ..., D\}$, where sigm $(x) = 1/(1 + \exp(-x))$, $\mathbf{W} \in \mathbb{R}^{H \times D}$ and $\mathbf{V} \in \mathbb{R}^{D \times H}$ are connection parameter matrices, $\mathbf{b} \in \mathbb{R}^{D}$ and $\mathbf{c} \in \mathbb{R}^{H}$ are bias parameter vectors, $\mathbf{v}_{< i}$ is the subvector $[v_1, ..., v_{i-1}]^{\top}$ and $\mathbf{W}_{:,<i}$ is a matrix made of the i - 1 first columns of \mathbf{W} .

corresponds to a neural network with several parallel $h_i(v_{< i})$ hidden layers

Advantages: Tractable $p(v) \rightarrow$ easier to train

Intern © Siemens AG 2017

Seite 94 May 2017



- > **NADE:** Neural Autoregressive Distributional Estimator
- Inspired from RBM Family therefore, energy based
- ➤ sampling each dimension one after another

 $p(\mathbf{v}) = \prod_{i=1}^{D} p(v_i | \mathbf{v}_{< i})$ and computes all $p(v_i | \mathbf{v}_{< i})$ using the feed-forward architecture

$$\mathbf{h}_{i}(\mathbf{v}_{< i}) = \operatorname{sigm}\left(\mathbf{c} + \mathbf{W}_{:,< i}\mathbf{v}_{< i}\right)$$
NADE is for binary data,
$$p(v_{i} = 1 | \mathbf{v}_{< i}) = \operatorname{sigm}\left(b_{i} + \mathbf{V}_{i,:}\mathbf{h}_{i}(\mathbf{v}_{< i})\right)$$

$$\mathbf{h}_{i+1}(\mathbf{v}_{< i+1}) = \operatorname{sigm}(\mathbf{c} + \sum_{k < i+1} W_{:,v_k}) = \operatorname{sigm}(\mathbf{W}_{:,v_i} + \mathbf{c} + \sum_{k < i} \mathbf{W}_{:,v_k})$$



for $i \in \{1, ..., D\}$, where sigm $(x) = 1/(1 + \exp(-x))$, $\mathbf{W} \in \mathbb{R}^{H \times D}$ and $\mathbf{V} \in \mathbb{R}^{D \times H}$ are connection parameter matrices, $\mathbf{b} \in \mathbb{R}^{D}$ and $\mathbf{c} \in \mathbb{R}^{H}$ are bias parameter vectors, $\mathbf{v}_{< i}$ is the subvector $[v_1, ..., v_{i-1}]^{\top}$ and $\mathbf{W}_{:,<i}$ is a matrix made of the i - 1 first columns of \mathbf{W} .

corresponds to a neural network with several parallel $h_i(v_{< i})$ hidden layers

Advantages: Tractable $p(v) \rightarrow$ easier to train

Intern © Siemens AG 2017

Seite 95 May 2017



Training NADE

→ Ground truth values of the pixels are used for conditioning when predicting subsequent values

- → *cost*: maximize log-likelihood (logL)
- → optimize to maximize the logL by *stochastic gradient descent* (SGD)

$$\mathcal{L}^{DocNADE}(\mathbf{v}) = \sum_{i=1}^{D} \log p(v_i | \mathbf{v}_{< i})$$

where, D is the number of words in document, \boldsymbol{v} and

autoregressive conditional is given by:

$$p(v_i = 1 | \mathbf{v}_{< i}) = \operatorname{sigm} (b_i + \mathbf{V}_{i,:} \mathbf{h}_i(\mathbf{v}_{< i}))$$

 $p(v_i = 1 | \mathbf{v}_{< i}) = \operatorname{sigm} (b_i + \mathbf{V}_{i,:} \mathbf{h}_i(\mathbf{v}_{< i}))$





where, $\mathbf{h}_i(\mathbf{v}_{< i}) = \operatorname{sigm} (\mathbf{c} + \mathbf{W}_{:, < i} \mathbf{v}_{< i})$

Intern © Siemens AG 2017

Seite 96 May 2017



Training NADE

 \rightarrow Ground truth values of the pixels are used for conditioning

NADE: Alternative to RBMs with tractable p(v)

→ optimize to maximize the logL by *stochastic gradient descent* (SGD)

$$\mathcal{L}^{DocNADE}(\mathbf{v}) = \sum_{i=1}^{D} \log p(v_i | \mathbf{v}_{< i})$$

where, D is the number of words in document, \mathbf{v} and

an autoregressive conditional is given by:

$$p(v_i = 1 | \mathbf{v}_{< i}) = \operatorname{sigm} (b_i + \mathbf{V}_{i,:} \mathbf{h}_i(\mathbf{v}_{< i}))$$

 $p(v_i = 1 | \mathbf{v}_{< i}) = \operatorname{sigm} (b_i + \mathbf{V}_{i,:} \mathbf{h}_i(\mathbf{v}_{< i}))$





where, $\mathbf{h}_i(\mathbf{v}_{< i}) = \operatorname{sigm} (\mathbf{c} + \mathbf{W}_{:,< i} \mathbf{v}_{< i})$

Intern © Siemens AG 2017

Seite 97 May 2017



7	7	3	į	1	4	3	1	ê.	2	7	7	3		ł	3	3	1	2	2	開き		-		感				No.	1
3	4	2	3	6	Z	i.	ł	8	¥	3	6	2	3	6	3	41	l	8	4	新	1	100				-		20	14
4	č	9	3	ŝ,	*ų,	1	5	3	6	4	6	9	3	5	1	1	3	3	6		the state	ANA ANA	1	and the second s		家	S.	4	Section of the
5	1	ę.	2	Ģ	3	7	7	ч	D	5	1	5	2	9	3	7	7	ч	D	No.	St.	S.	Harris Contraction	Contraction of the second			X	1	
B	8	9	5	5	14	÷	Ĵ	\$	7	B	8	9	5	δ	4	4	3	3	7	4	1		44			and the second s	藏	1	
9	3	\$	Ó	•	q	3	1	ŕ	1	9	3	4	6	3	9	3	1	7	Y	No.	the second	ST.	藩						
8	3	3	9	5	X	2	1	5	14	8	3	4	9	5	1	2	1	5	4			•	1	The second		A STATE	語		100
5	11.	2	.*	14.	2	5	7	2	5	6	120	9		1.1	2	8	H	2	3	-		The second	1		a state		100	and and	1/2
1	6	1	4	0	1	6	4	3	60	1	0	1	4	0	1	6	\$		60	3.	19	and the second	2	「ない	No.		and the second	100 M	ALS.
0	2	6	\$	42	2	9	\$	Ş	4	0	2	6	5	23	2	9	8	8	9	1		•	175	and the second		-	100	J.	1

(Left): samples from NADE trained on a binary version of MNIST (Middle): probabilities from which pixel was samples (**Right**): Visualization of some of the rows

https://ift6135h18.files.wordpress.com/2018/04/autoregressive_gen.pdf Intern © Siemens AG 2017

Seite 98 May 2017





(Left): samples from NADE trained on a binary version of MNIST (Middle): probabilities from which pixel was samples (**Right**): Visualization of some of the rows

https://ift6135h18.files.wordpress.com/2018/04/autoregressive_gen.pdf Intern © Siemens AG 2017

May 2017 Seite 99



NADE models **binary** input or *real-valued* input using realNADE

How to model count data in NADE architecture? Example: Text document i.e., word counts.



Seite 100 May 2017





Intern © Siemens AG 2017

Seite 101 May 2017

Corporate Technology



→ NADE model **binary** input or real-valued input using realNADE

How to model count data in NADE architecture? Example: Text document i.e., word counts.



Seite 102 May 2017



Probabilistic graphical model that *learns topics over sequences of words*

- learn topic-word distribution based on word co-occurrences
- learn distributed word representations
- > compute joint distribution via autoregressive conditionals
- compute joint distribution or log-likelihood for a document, v in *language modeling fashion*
- > interpreted as a neural network with several parallel hidden layers
- \succ predict the word v_i, given the sequence of preceding words v_i

Modeling documents in NADE





Probabilistic graphical model that *learns topics over sequences of words*

- learn topic-word distribution based on word co-occurrences
- learn distributed word representations
- > compute joint distribution via autoregressive conditionals
- compute joint distribution or log-likelihood for a document, v in *language modeling fashion*
- > interpreted as a neural network with several parallel hidden layers
- > predict the word v_i, given the sequence of preceding words v_i Limitations:
- \succ does not take into account the following words $\mathbf{v}_{>i}$ in the sequence
- poor in modeling short-text documents

(i.e., does not use pre-trained word embeddings)

Intern © Siemens AG 2017





Probabilistic graphical model that *learns topics over sequences of words*

- learn topic-word distribution based on word co-occurrences
- learn distributed word representations
- > compute joint distribution via autoregressive conditionals
- compute joint distribution or log-likelihood for a document, v in *language modeling fashion*
- ➤ interpreted as a neural network with several parallel hidden layers
- > predict the word v_i, given the sequence of preceding words v_i Limitations:
- \succ does not take into account the following words $\mathbf{v}_{>i}$ in the sequence
- > poor in modeling short-text documents due to limited context i.e., co-occurrences



DocNADE

Intern © Siemens AG 2017

Seite 105 May 2017



DocNADE Formulation

- ➢ inspired by RBM, RSM and NADE models
- \succ models the joint distribution of all words v_{i}

 $v_i \in \{1, ..., K\}$

- i.e., the index of the *ith* word in the dictionary of vocabulary size K
- > a document v of size D is represented as, $v = [v_1, ..., v_D]$
- \succ joint distribution p(v) computed via each autoregressive conditionals,

 $p(\mathbf{v}) = \prod_{i=1}^{D} p(v_i | \mathbf{v}_{< i})$



Intern © Siemens AG 2017

Seite 106 May 2017



DocNADE Formulation

- ➢ inspired by RBM, RSM and NADE models
- > models the joint distribution of all words v_i

 $v_i \in \{1, ..., K\}$

- i.e., the index of the *ith* word in the dictionary of vocabulary size K
- > a document v of size D is represented as, $v = [v_1, ..., v_D]$
- \succ joint distribution p(v) computed via each autoregressive conditionals,

 $p(\mathbf{v}) = \prod_{i=1}^{D} p(v_i | \mathbf{v}_{< i}) = \prod_{i=1}^{D} \frac{\exp(b_w + \mathbf{U}_{w,:} \overrightarrow{\mathbf{h}}_i(\mathbf{v}_{< i}))}{\sum_{w'} \exp(b_{w'} + \mathbf{U}_{w',:} \overrightarrow{\mathbf{h}}_i(\mathbf{v}_{< i}))}$ via a feed-forward neural network, where $\mathbf{v}_{< i} \in \{v_1, ..., v_{i-1}\}$

autoregressive conditional, $p(\mathbf{v}_{i=3} | \mathbf{v}_{<3})$ h h₁ h, h, w (VD) v₂ v_3

DocNADE

Intern © Siemens AG 2017

Seite 107 May 2017



DocNADE Formulation

- ➢ inspired by RBM, RSM and NADE models
- > models the joint distribution of all words v_i

 $v_i \in \{1, ..., K\}$

- i.e., the index of the *ith* word in the dictionary of vocabulary size K
- > a document v of size D is represented as, $v = [v_1, ..., v_D]$
- \succ joint distribution p(v) computed via each autoregressive conditionals,

$$p(\mathbf{v}) = \prod_{i=1}^{D} p(v_i | \mathbf{v}_{< i}) = \prod_{i=1}^{D} \frac{\exp(b_w + \mathbf{U}_{w,:} \overrightarrow{\mathbf{h}}_i(\mathbf{v}_{< i}))}{\sum_{w'} \exp(b_{w'} + \mathbf{U}_{w',:} \overrightarrow{\mathbf{h}}_i(\mathbf{v}_{< i}))}$$

via a feed-forward neural network, where $\mathbf{v}_{< i} \in \{v_1, ..., v_{i-1}\}$

Topic matrix where,
$$\mathbf{W} \in \mathbb{R}^{H \times K}$$
 and $\mathbf{U} \in \mathbb{R}^{K \times H}$

autoregressive conditional, $p(\mathbf{v}_{i=3} | \mathbf{v}_{<3})$



DocNADE

Intern © Siemens AG 2017

Seite 108 May 2017





S1: Deal with stock index fallS2: Brace for market share drop

Intern © Siemens AG 2017

Seite 109 May 2017




S1: Deal with stock index fallS2: Brace for market share drop

Intern © Siemens AG 2017

Seite 110 May 2017





S1: Deal with stock index fallS2: Brace for market share drop

Seite 111 May 2017





Corporate Technology

Seite 112 May 2017





Seite 113 May 2017



DocNADE Formulation

Properties of weight matrix, W

- ➤ each column-vector W:,vi
 - \rightarrow a vector for the word v_i
- \succ each row-vector $\mathbf{W}_{j,:}$
 - → a distribution over vocabulary of size K, representing the *jth* topic
- > exploit column-vector property and introduce

additional matrix E, to incorporate pre-trained word embeddings

or distributional word representations









DocNADE Formulation

- ➢ inspired by RBM, RSM and NADE models
- > models the joint distribution of all words v_i

 $v_i \in \{1, ..., K\}$

- i.e., the index of the *ith* word in the dictionary of vocabulary size K
- > a document **v** of size D is represented as, $\mathbf{v} = [v_1, ..., v_D]$
- \succ joint distribution p(v) computed via each autoregressive conditionals,

$$p(\mathbf{v}) = \prod_{i=1}^{D} p(v_i | \mathbf{v}_{< i}) = \prod_{i=1}^{D} \frac{\exp(b_w + \mathbf{U}_{w,:} \overrightarrow{\mathbf{h}}_i(\mathbf{v}_{< i}))}{\sum_{w'} \exp(b_{w'} + \mathbf{U}_{w',:} \overrightarrow{\mathbf{h}}_i(\mathbf{v}_{< i}))}$$

via a feed-forward neural network, where $\mathbf{v}_{< i} \in \{v_1, ..., v_{i-1}\}$

Topic matrix where,
$$\mathbf{W} \in \mathbb{R}^{H \times K}$$
 and $\mathbf{U} \in \mathbb{R}^{K \times H}$

autoregressive conditional, $p(\mathbf{v}_{i=3} | \mathbf{v}_{<3})$



DocNADE

DOES NOT take into account the following words $v_{>i}$

Intern © Siemens AG 2017

Seite 115 May 2017

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior



DocNADE extensions......

Exploiting distributed word representations (word embeddings)

+ Full Context information

Checkout the backup slides, if interested !!!

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Seite 116 May 2017

Document Representation in iDocNADE variants

book			jesus			windows			gun		
neighbors	s_i	s_g	neighbors	s_i	s_g	neighbors	s_i	s_g	neighbors	s_i	s_g
books	.61	.84	christ	.86	.83	dos	.74	.34	guns	.72	.79
reference	.52	.51	god	.78	.63	files	.63	.36	firearms	.63	.63
published	.46	.74	christians	.74	.49	version	.59	.43	criminal	.63	.33
reading	.45	.54	faith	.71	.51	file	.59	.36	crime	.62	.42
author	.44	.77	bible	.71	.51	unix	.52	.47	police	.61	.43

Tasks:

→Information retrieval

- \rightarrow document representation
- \rightarrow word representation
- \rightarrow text classification
- \rightarrow text clustering, etc.

Table 6: 20NS dataset: The five nearest neighbors by iDoc-NADE. s_i : Cosine similarity between the word vectors from iDocNADE, for instance vectors of *jesus* and *god*. s_g : Cosine similarity in embedding vectors from glove.

Visualizing or interpreting filters i.e., W matrix. (column vectors \rightarrow word embeddings)

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Document Representation in iDocNADE variants

Checkout the backup slides!!!								
DocNADE	iDocNADE	DocNADEe	Tasks:					
beliefs, muslims,	scripture, atheists,	atheists, christianity,	\rightarrow document representation \rightarrow word representation					
forward, alt,	sin, religions,	belief, eternal,						
islam, towards,	christianity, lord,	atheism, catholic,	\rightarrow text classification					
atheism, christianity,	bible, msg,	bible, arguments,	\rightarrow text clustering, etc.					
hands, opinions	heaven, jesus	islam, religions						
0.44	0.46	0.52						
Topics (top 10 words) of 20NS with coherence								

Visualizing or interpreting filters i.e., W matrix. (row vectors → topic information)

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Seite 118 May 2017



Document Representation in iDocNADE variants

Limitations:

- \rightarrow Bag-of-word models
- \rightarrow missing word ordering
- \rightarrow missing local dynamics of the sequence

Extension(s):

- → Joint neural autoregressive topic (e.g., DocNADE) and neural language models (e.g., RNN or LSTM)
- → Introduce language concepts (e.g., word ordering, latent syntactic and semantic information) into DocNADE

Further reading: https://arxiv.org/abs/1810.03947

Intern © Siemens AG 2017

Tasks:

- →Information retrieval
- \rightarrow document representation
- \rightarrow word representation
- \rightarrow text classification
- \rightarrow text clustering, etc.

Compositional Distributional Semantics

 $L_t = -\log P(x_t | x_{t-1}, x_{t-2}, \dots x_1)$

Compositional models in Neural Networks: RNN-LM or Recursive Neural Network or seq2seq (Lecture-05)



Seite 120 May 2017

Compositional Distributional Semantics

Compositional models in Neural Networks: RNN-LM or Recursive Neural Network or seq2seq (Lecture-05)

RNN-LM captures *local dynamics* of the sequence, i.e., word ordering, syntactic and semantic information from word co-occurrences in collocation/nearby patterns

→RNN-LMs (or LSTM-LM) lack in capturing global semantics, i.e. long-term dependencies

→ DocNADEs capture global semantics in form of topics



Generative Recurrent Neural Network (RNN)

$$L_t = -\log P(x_t | x_{t-1}, x_{t-2}, \dots x_1)$$

 $\begin{bmatrix} 9\\1 \end{bmatrix} \begin{bmatrix} 5\\3 \end{bmatrix} \begin{bmatrix} 7\\1 \end{bmatrix} \begin{bmatrix} 8\\5 \end{bmatrix} \begin{bmatrix} 9\\1 \end{bmatrix} \begin{bmatrix} 4\\3 \end{bmatrix}$ The cat sat on the mat.

Recursive Neural Network

Intern © Siemens AG 2017

Seite 121 May 2017

Distributional Representations: Local and Global Semantics

Combine or Joint training of DocNADE and LSTM-LM: textTOvec



Further reading: Gupta et al, 2018. *textTOvec*: <u>https://arxiv.org/abs/1810.03947</u>

Intern © Siemens AG 2017

Seite 122 May 2017

Metric Learning for Similarity (overview)



Aditya Thyagarajan and Jonas Mueller. 2016. Siamese Recurrent Architectures for Learning Sentence Similarity

Intern © Siemens AG 2017

Seite 123 May 2017



 \rightarrow unsupervised learning for distributed representations in neural networks

 \rightarrow pre-train and transfer learning to initialize neural networks for better convergence

 \rightarrow distributed word representations encode syntactic, semantic information into vectors

 \rightarrow metric learning to compute similarity over representations learned

Intern © Siemens AG 2017

Seite 124 May 2017



References, Resources and Further Reading

- https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/slides/L1_intro.pdf
- https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/
- Generative models: <u>https://ift6135h18.wordpress.com/author/aaroncourville/page/1/</u>
- Boltzmann Machines: <u>http://www.iro.umontreal.ca/~bengioy/ift6266/H12/html/contents.html#contents-en</u>
- Boltzmann Machines: <u>https://ift6266h15.wordpress.com/category/lectures/page/1/</u>
- http://www.iro.umontreal.ca/~bengioy/ift6266/H16/representation-learning.pdf
- https://ift6266h16.wordpress.com/2016/03/14/lecture-19-march-17th-2016-representation-learning/
- Generative models: <u>https://ift6266h15.files.wordpress.com/2015/03/chapter21.pdf</u>
- http://cse.iitkgp.ac.in/~sudeshna/courses/DL17/Autoencoder-15-Mar-17.pdf
- https://www.cl.cam.ac.uk/~pv273/slides/UCLSlides.pdf
- Autoregressive networks: <u>https://ift6135h18.files.wordpress.com/2018/04/autoregressive_gen.pdf</u>
- https://blog.acolyer.org/2016/04/21/the-amazing-power-of-word-vectors/
- Plotting Samples and Filters: <u>http://deeplearning.net/tutorial/utilities.html#how-to-plot</u>
- iDocNADEe: Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior
- ctx-DocANDEe (textTOvec): Gupta et al, 2018. textTOvec: <u>https://arxiv.org/abs/1810.03947</u>

Siamese LSTM: Gupta et al, 2018. Replicated Siamese LSTM in Ticketing System for Similarity Learning and Retrieval in Asymmetric Texts Intern © Siemens AG 2017



Thanks !!! Question ???

Write me, if interested in

firstname.lastname@siemens.com

@Linkedin: https://www.linkedin.com/in/pankaj-gupta-6b95bb17/

About my research contributions:

https://scholar.google.com/citations?user=_YjIJF0AAAAJ&hl=en

Intern © Siemens AG 2017

Seite 126 May 2017



Backup slides, if interested !!!

Intern © Siemens AG 2017

Seite 127 May 2017



Informed Document Autoregressive Topic Model with Word Embeddings

Source Text

Motivation1: Contextual Information

In biological brains, we study noisy **neurons** at cellular level \rightarrow

Like biological brains, study of noisy **neurons** in artificial neural networks

Sense/Topic

- "biological neural network"
- "artificial neural network" \rightarrow

Training Contexts (Preceding + Following) Like biological brains, study of noisy + in artificial neural networks \rightarrow "biological neural network"

Like biological brains, study of noisy + in artificial neural networks \rightarrow

- Sense/Topic of "neurons"
- "artificial neural network"

Context information around words helps in determining their actual meaning !!!

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Seite 128 May 2017



Informed Document Autoregressive Topic Model with Word Embeddings

Motivation2: Distributional Semantics for Lack of Context
"Lack of Context" in short-text documents, e.g., headlines, tweets, etc.

> "Lack of Context" in a corpus of few documents



TO RESCUE: Use External/additional information, e.g., word embeddings (encodes semantic and syntactic relatedness in words)



Cosine similarity in Word Embedding space

Intern © Siemens AG 2017



DocNADE variants: Contextualized DocNADE (iDocNADE)



> incorporating full contextual information around words in a document (*preceding* and *following* words)

- > **boost the likelihood** of each word and subsequently the document
- improved representation learning

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Seite 130 May 2017



DocNADE variants: Contextualized DocNADE (iDocNADE)



Seite 131 May 2017



DocNADE variants: Contextualized DocNADE (iDocNADE)





DocNADE variants: **DocNADE** + Embedding Priors 'e' (DocNADEe)



$$\mathcal{L}^{DocNADE}(\mathbf{v}) = \sum_{i=1}^{D} \log p(v_i | \mathbf{v}_{< i})$$
$$\overrightarrow{\mathbf{h}}_i(\mathbf{v}_{< i}) = g(\mathbf{c} + \sum_{k < i} \mathbf{W}_{:, v_k})$$

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Seite 133 May 2017



DocNADE variants: DocNADE + Embedding Priors 'e' (DocNADEe)



 $\mathcal{L}^{DocNADE}(\mathbf{v}) = \sum_{i=1}^{D} \log p(v_i | \mathbf{v}_{< i})$ $\overrightarrow{\mathbf{h}}_i(\mathbf{v}_{< i}) = g(\mathbf{c} + \sum_{k < i} \mathbf{W}_{:,v_k})$

introduce weighted word embedding aggregation at each autoregressive step k

E as fixed prior

- topics with word embeddings
- generate a complementary

textual representation (duality)



Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Seite 134 May 2017



DocNADE variants: DocNADE + Embedding Priors 'e' (DocNADEe)



 introduce weighted word embedding aggregation at each autoregressive step k
 E as fixed prior
 topics with word embeddings
 generate a complementary textual representation (*duality*)
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w
 w

$$\mathcal{L}^{DocNADE}(\mathbf{v}) = \sum_{i=1}^{D} \log p(v_i | \mathbf{v}_{< i}) \qquad \qquad \mathbf{h}_i^{\vec{e}}(\mathbf{v}_{< i}) = g(\vec{\mathbf{c}} + \sum_{k < i} \mathbf{W}_{:,v_k} + \mathbf{\lambda} \sum_{k < i} \mathbf{E}_{:,v_k})$$
$$\qquad \qquad \mathbf{h}_i(\mathbf{v}_{< i}) = g(\mathbf{c} + \sum_{k < i} \mathbf{W}_{:,v_k})$$

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Seite 135 May 2017



Deep DocNADEs Variants with/without Embedding Priors

- > Deep, multiple hidden layer architectures
- > Adding new hidden layers as in a regular deep feed-forward neural network

$$\overrightarrow{\mathbf{h}}_{i}^{(d)}(\mathbf{v}_{< i}) = g(\overrightarrow{\mathbf{c}}^{(d)} + \mathbf{W}^{(d)} \cdot \overrightarrow{\mathbf{h}}_{i}^{(d-1)}(\mathbf{v}_{< i}))$$

introduce embedding prior \mathbf{E} in the first hidden layer, i.e.,

$$\overrightarrow{\mathbf{h}}_{i}^{e,(1)} = g(\overrightarrow{\mathbf{c}}^{(1)} + \sum \mathbf{W}_{:,v_{k}}^{(1)} + \lambda \sum \mathbf{E}_{:,v_{k}})$$

DeepVairant1	DeepVairant2	DeepVairant3	DeepVairant4
DeepDNE	iDeepDNE	DeepDNEe	iDeepDNEe

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior

Intern © Siemens AG 2017

Seite 136 May 2017

Document Representation in iDocNADE variants

With the trained *iDocNADEe* (or *DocNADE* variants), the representation ($\mathfrak{h}^e \in \mathbb{R}^H$) for a new document \mathbf{v}^* of size D^* is extracted by summing the hidden representations from the forward and backward networks to account for the context information around each word in the words' sequence, as

$$\overrightarrow{\mathbf{h}^{e}}(\mathbf{v}^{*}) = g(\overrightarrow{\mathbf{c}} + \sum_{k \leq D^{*}} \mathbf{W}_{:,v_{k}^{*}} + \lambda \sum_{k \leq D^{*}} \mathbf{E}_{:,v_{k}^{*}})$$

$$\overleftarrow{\mathbf{h}^{e}}(\mathbf{v}^{*}) = g(\overleftarrow{\mathbf{c}} + \sum_{k \geq 1} \mathbf{W}_{:,v_{k}^{*}} + \lambda \sum_{k \geq 1} \mathbf{E}_{:,v_{k}^{*}})$$

Therefore;
$$\overleftarrow{\mathbf{h}^{e}} = \overrightarrow{\mathbf{h}^{e}}(\mathbf{v}^{*}) + \overleftarrow{\mathbf{h}^{e}}(\mathbf{v}^{*})$$

Tasks:

- →Information retrieval
- \rightarrow document representation
- \rightarrow word representation
- \rightarrow text classification
- \rightarrow text clustering, etc.

Intern © Siemens AG 2017

Seite 137 May 2017

Gupta et al, 2018. Document Informed Neural Autoregressive Topic Models with Distributional Prior