

6125021:Introduction to Combinatorics

-Deep Learning: An Overview

Dr. Zichen Xu (徐子晨)

Department of Computer Science and Technology

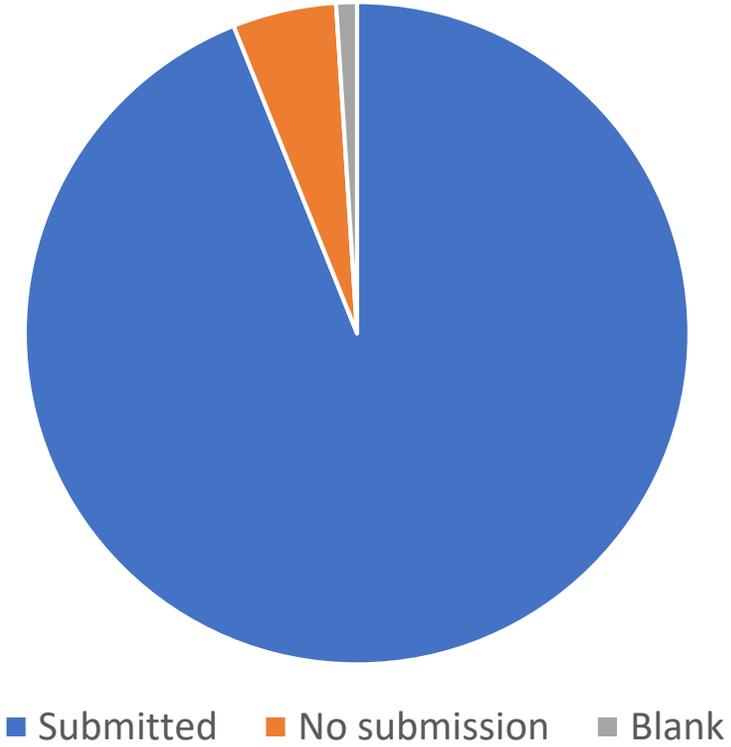
College of Information Science and Technology

Notes From Last Lecture

- We have discussed a list of following topics
 - Complexity Theory
 - Set Problems
 - Solution space and scale of the problem
 - Polynomial Time Complexity

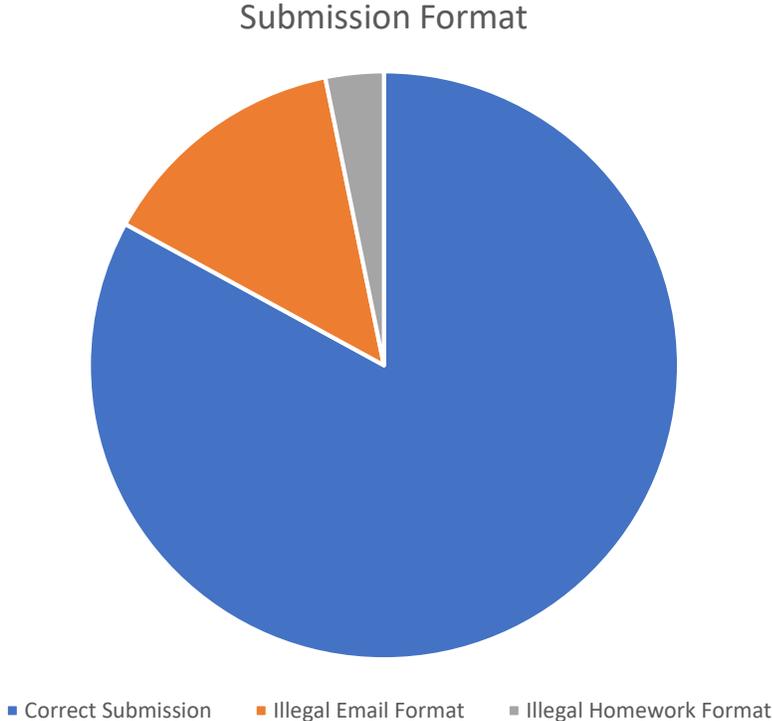
The Homework I Statistics

Homework Submission



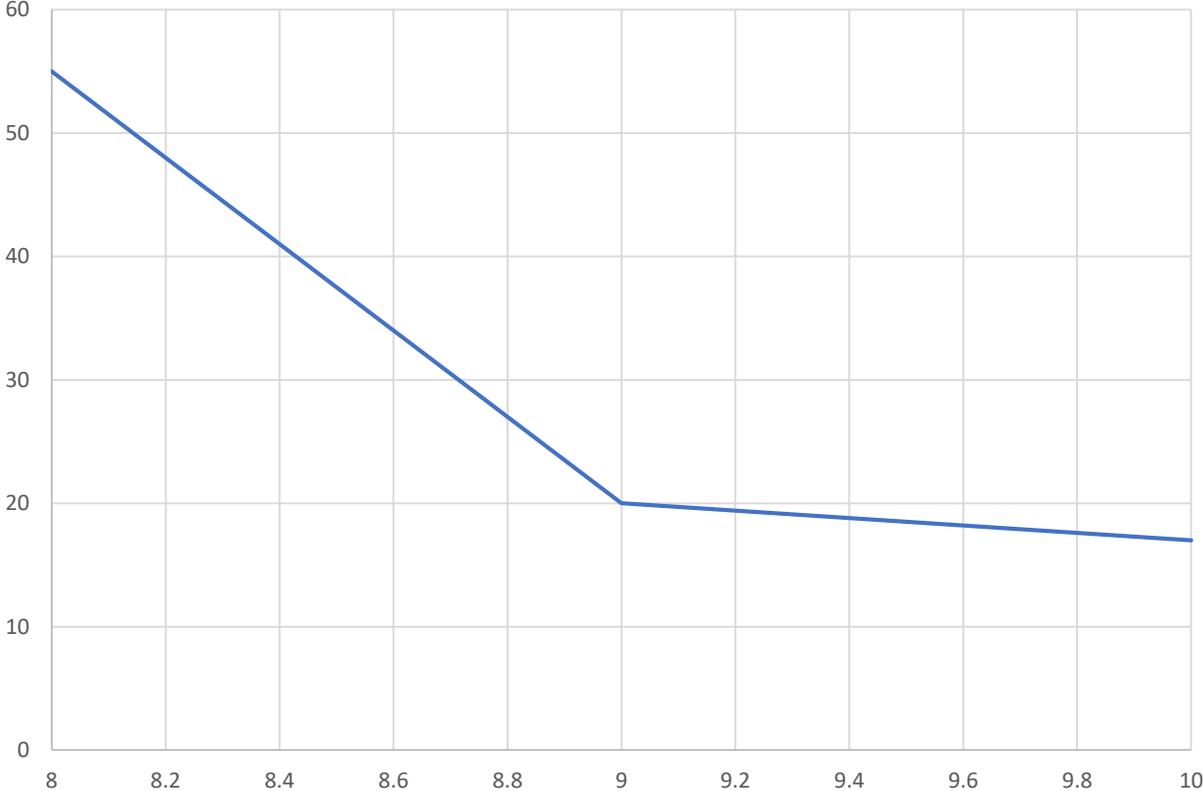
The Homework I Statistics

Earliest Submission Time: Sep. 15 4:54PM
Latest Submission Time: Sep. 23, 8:41PM



The Homework I Statistics

Grade Distribution



Acknowledgements

- Many of the pictures, results, and other materials are taken from:
 - Aarti Singh, Carnegie Mellon University
 - Andrew Ng, Stanford University
 - Barnabas Póczos, Carnegie Mellon University
 - Christopher Manning, Stanford University
 - Geoffrey Hinton, Google & University of Toronto
 - Richard Socher, MetaMind
 - Richard Turner, University of Cambridge
 - Yann LeCun, New York University
 - Yoshua Bengio, Université de Montréal
 - Weifeng Li, Victor Benjamin, Xiao Liu, and Hsinchun Chen, ASU

Outline

- Introduction
 - Motivation
 - Deep Learning Intuition
- Neural Network
 - Neuron, Feedforward algorithm, and Backpropagation algorithm
 - Limitations
- Deep Learning
 - Deep Belief Network: Autoencoder & Restricted Boltzman Machine
 - Convolutional Neural Network
- Deep Learning for Text Mining
 - Word Embedding
 - Recurrent Neural Network
 - Recursive Neural Network
- Conclusions

Outline

- Introduction
 - Motivation
 - Deep Learning Intuition
- Neural Network
 - Neuron, Feedforward algorithm, and Backpropagation algorithm
 - Limitations
- Deep Learning
 - Deep Belief Network: Autoencoder & Restricted Boltzman Machine
 - Convolutional Neural Network
- Deep Learning for Text Mining
 - Word Embedding
 - Recurrent Neural Network
 - Recursive Neural Network
- Conclusions

nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

At last — a computer program that
can beat a champion Go player **PAGE 484**

ALL SYSTEMS GO

CONSERVATION

SONGBIRDS À LA CARTE

Illegal harvest of millions
of Mediterranean birds

PAGE 452

RESEARCH ETHICS

SAFEGUARD TRANSPARENCY

Don't let openness backfire
on individuals

PAGE 459

POPULAR SCIENCE

WHEN GENES GOT 'SELFISH'

Dawkins's calling
card forty years on

PAGE 462

NATURE.COM/NATURE

28 January 2016 £10

Vol. 529, No. 7587



In “Nature” 27 January 2016:

- “DeepMind’s program AlphaGo beat Fan Hui, the European Go champion, five times out of five in tournament conditions...”
- “AlphaGo was not preprogrammed to play Go: rather, it learned using a general-purpose algorithm that allowed it to interpret the game’s patterns.”
- “...AlphaGo program applied **deep learning** in neural networks (convolutional NN) — brain-inspired programs in which connections between layers of simulated neurons are strengthened through examples and experience.”

10 BREAKTHROUGH TECHNOLOGIES 2013

Deep Learning

With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart. →

Temporary Social Media

Messages that quickly self-destruct could enhance the privacy of online communications and make people freer to be spontaneous. →

Prenatal DNA Sequencing

Reading the DNA of fetuses will be the next frontier of the genomic revolution. But do you really want to know about the genetic problems or musical aptitude of your unborn child? →

Add Man

Ske
prin
wor
mar
the
tech
jet p

Memory Implants

A maverick neuroscientist believes he has deciphered the code by which the brain

Smart Watches

Ultra-Efficient Solar Power

Doubling the efficiency of a solar cell would completely

Big Pho

Coll
ana
fron
pho

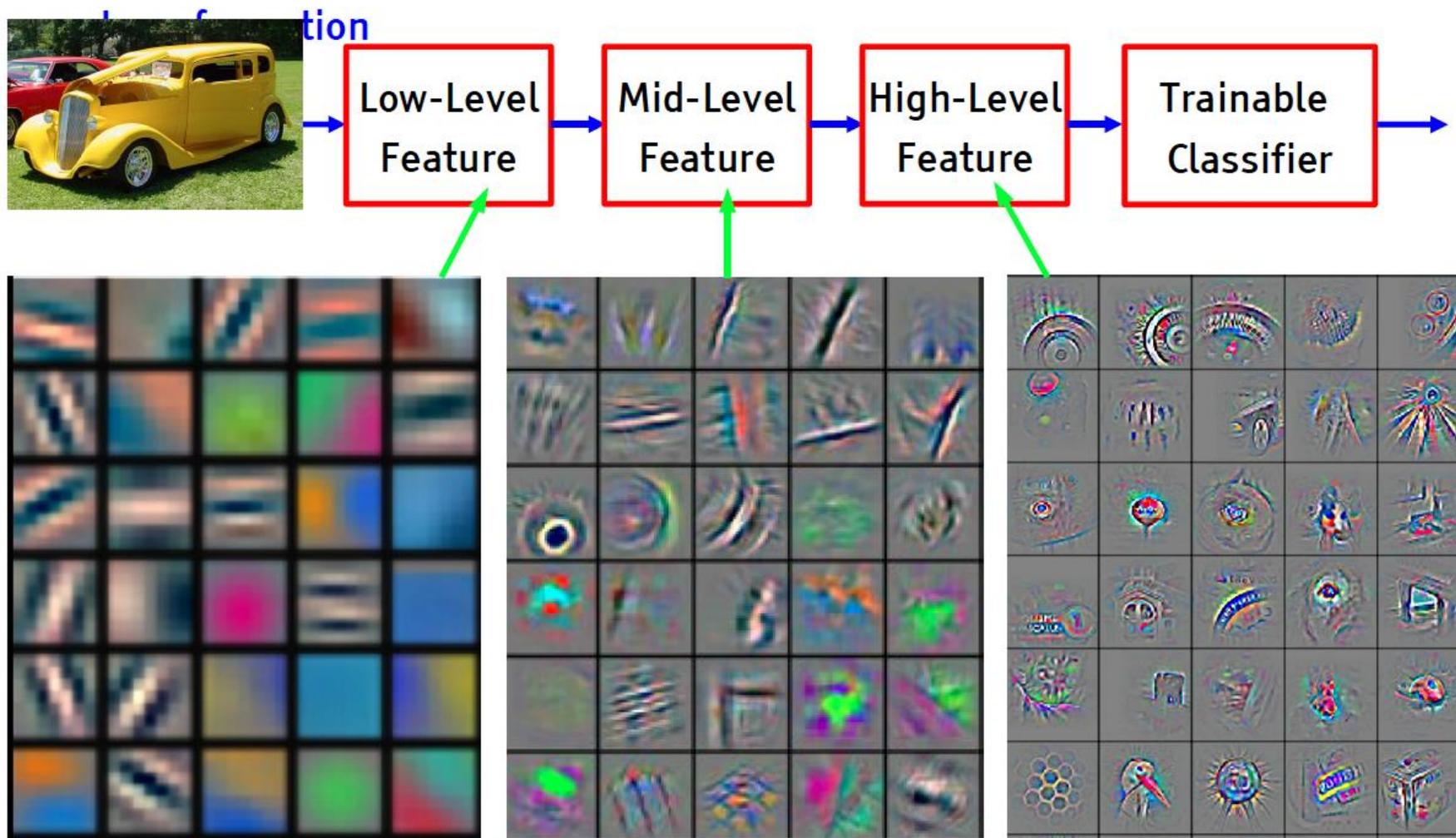
Deep Learning Today

- Advancement in speech recognition in the last 2 years
 - A few long-standing performance records were broken with deep learning methods
 - Microsoft and Google have both deployed DL-based speech recognition systems in their products
- Advancement in Computer Vision
 - Feature engineering is the bread-and-butter of a large portion of the CV community, which creates some resistance to feature learning
 - But the record holders on ImageNet and Semantic Segmentation are convolutional nets
- Advancement in Natural Language Processing
 - Fine-grained sentiment analysis, syntactic parsing
 - Language model, machine translation, question answering

Motivations for Deep Architectures

- Insufficient depth can hurt
 - With shallow architecture (SVM, NB, KNN, etc.), the required number of nodes in the graph (i.e. computations, and also number of parameters, when we try to learn the function) may grow very large.
 - Many functions that can be represented efficiently with a deep architecture cannot be represented efficiently with a shallow one.
- The brain has a deep architecture
 - The visual cortex shows a sequence of areas each of which contains a representation of the input, and signals flow from one to the next.
 - Note that representations in the brain are in between dense distributed and purely local: they are **sparse**: about 1% of neurons are active simultaneously in the brain.
- Cognitive processes seem deep
 - Humans organize their ideas and concepts hierarchically.
 - Humans first learn simpler concepts and then compose them to represent more abstract ones.
 - Engineers break-up solutions into multiple levels of abstraction and processing

Deep Learning = Learning Hierarchical Representations



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

Learning Representations: a challenge for ML, CV, AI, Neuroscience, Cognitive Science...

■ **How do we learn representations of the perceptual world?**

- ▶ How can a perceptual system build itself by looking at the world?
- ▶ How much prior structure is necessary

■ **ML/AI: how do we learn features or feature hierarchies?**

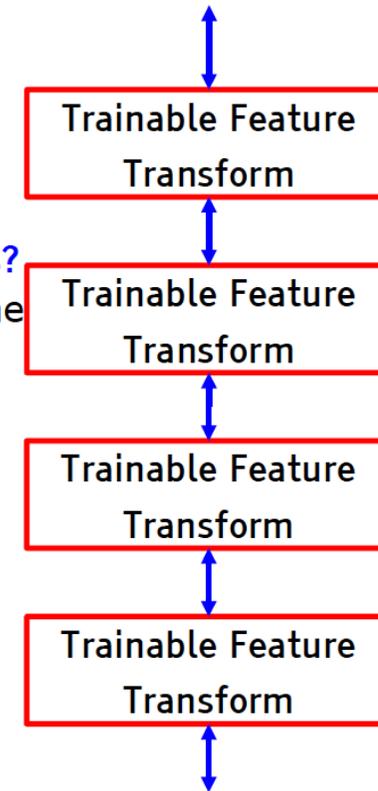
- ▶ What is the fundamental principle? What is the learning algorithm? What is the architecture?

■ **Neuroscience: how does the cortex learn perception?**

- ▶ Does the cortex "run" a single, general learning algorithm? (or a small number of them)

■ **CogSci: how does the mind learn abstract concepts on top of less abstract ones?**

■ **Deep Learning addresses the problem of learning hierarchical representations with a single algorithm**



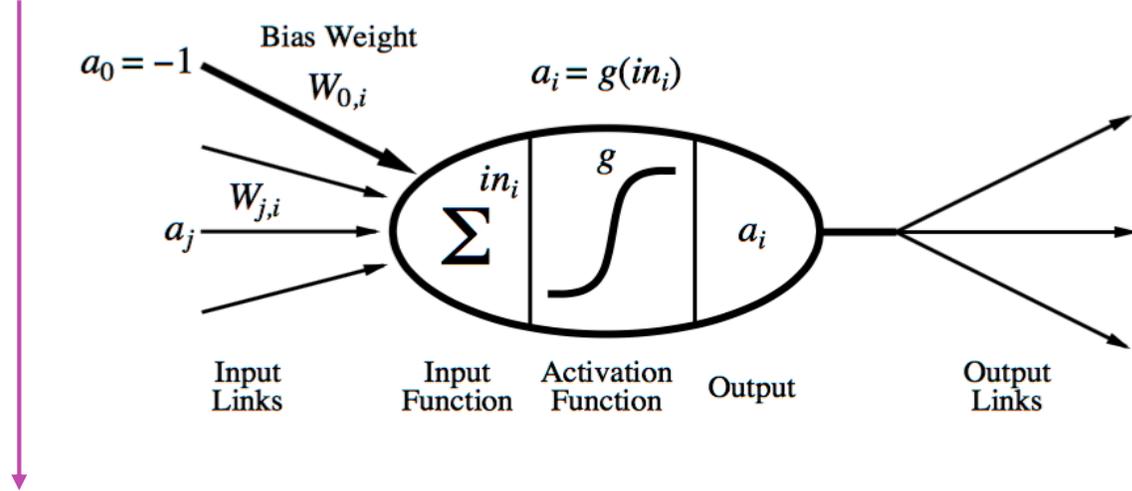
Outline

- Introduction
 - Motivation
 - Deep Learning Intuition
- **Neural Network**
 - Neuron, Feedforward algorithm, and Backpropagation algorithm
 - Limitations
- Deep Learning
 - Deep Belief Network: Autoencoder & Restricted Boltzman Machine
 - Convolutional Neural Network
- Deep Learning for Text Mining
 - Word Embedding
 - Recurrent Neural Network
 - Recursive Neural Network
- Conclusions

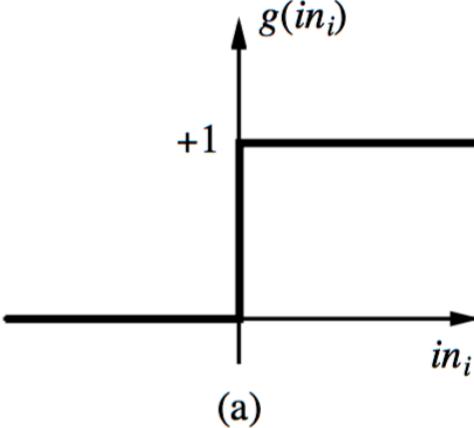
Neural Network: A Neuron

- A neuron is a computational unit in the neural network that exchanges messages with each other.

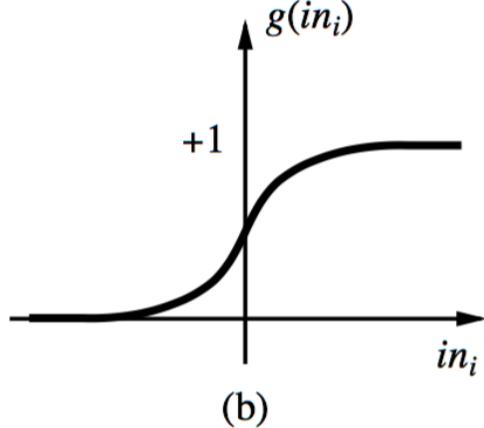
$$a_i \leftarrow g(in_i) = g(\sum_j W_{j,i} a_j)$$



Possible activation functions:

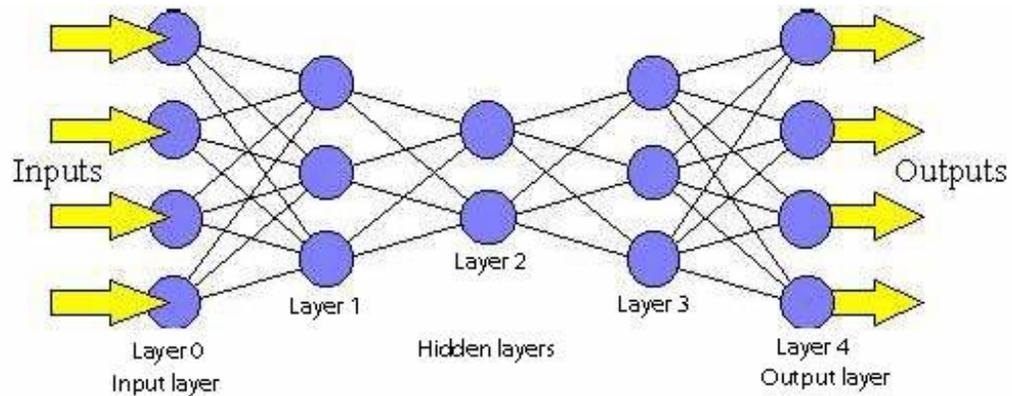


(a) Step function/threshold function



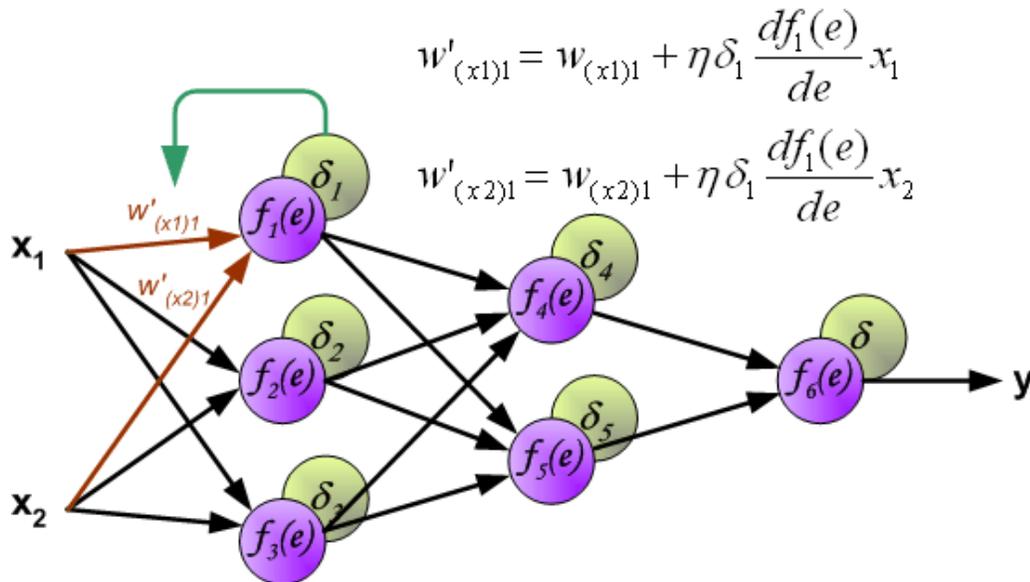
(b) Sigmoid function (a.k.a, logistic function)

Feed forward/Backpropagation Neural Network



Feed forward algorithm:

- Activate the neurons from the bottom to the top.



Backpropagation:

- Randomly initialize the parameters
- Calculate total error at the top, $f_6(e)$
- Then calculate contributions to error, δ_n , at each step going backwards.

Limitations of Neural Networks

Random initialization + densely connected networks lead to:

- High cost
 - Each neuron in the neural network can be considered as a logistic regression.
 - Training the entire neural network is to train all the interconnected logistic regressions.
- Difficult to train as the number of hidden layers increases
 - Recall that logistic regression is trained by gradient descent.
 - In backpropagation, gradient is progressively getting more dilute. That is, below top layers, the correction signal δ_n is minimal.
- Stuck in local optima
 - The objective function of the neural network is usually not convex.
 - The random initialization does not guarantee starting from the proximity of global optima.
- Solution:
 - Deep Learning/Learning multiple levels of representation

Outline

- Introduction
 - Motivation
 - Deep Learning Intuition
- Neural Network
 - Neuron, Feedforward algorithm, and Backpropagation algorithm
 - Limitations
- **Deep Learning**
 - **Deep Belief Network: Autoencoder & Restricted Boltzman Machine**
 - **Convolutional Neural Network**
- Deep Learning for Text Mining
 - Word Embedding
 - Recurrent Neural Network
 - Recursive Neural Network
- Conclusions

Deep Learning

We will cover two major deep learning models:

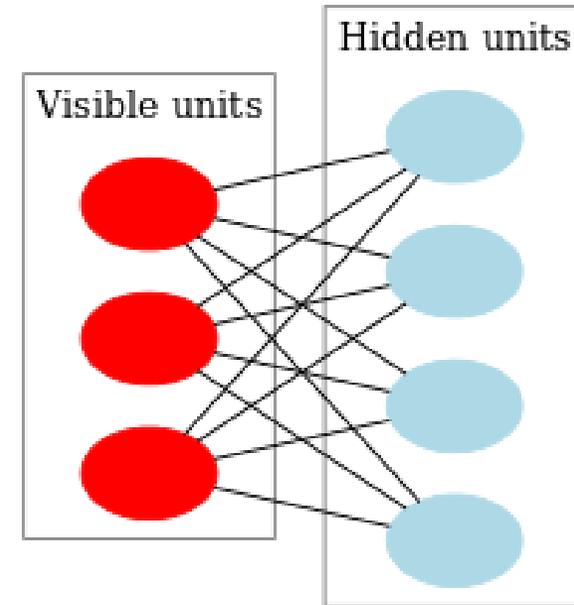
- ***Deep Belief Networks*** and ***Autoencoders*** employs layer-wise unsupervised learning to initialize each layer and capture multiple levels of representation simultaneously.
 - Hinton, G. E, Osindero, S., and Teh, Y. W. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527-1554.
 - Bengio, Y., Lamblin, P., Popovici, P., Larochelle, H. (2007). Greedy Layer-Wise Training of Deep Networks, *Advances in Neural Information Processing Systems* 19
- ***Convolutional Neural Network*** organizes neurons based on animal's visual cortex system, which allows for learning patterns at both local level and global level.
 - Y. LeCun, L. Bottou, Y. Bengio and P. Haffner: Gradient-Based Learning Applied to Document Recognition, *Proceedings of the IEEE*, 86(11):2278-2324, November 1998

Deep Belief Networks

- A *deep belief network* (DBN) is a probabilistic, generative model made up of multiple layers of hidden units.
 - A composition of simple learning modules that make up each layer
- A DBN can be used to generatively pre-train a DNN by using the learned DBN weights as the initial DNN weights.
 - Back-propagation or other discriminative algorithms can then be applied for fine-tuning of these weights.
- Advantages:
 - Particularly helpful when limited training data are available
 - These pre-trained weights are closer to the optimal weights than are randomly chosen initial weights.

Restricted Boltzmann Machine

- A DBN can be efficiently trained in an unsupervised, layer-by-layer manner, where the layers are typically made of *restricted Boltzmann machines* (RBM).
- RBM: undirected, generative energy-based model with a "visible" input layer and a hidden layer, and connections between the layers but not within layers.
- The training method for RBMs proposed by Geoffrey Hinton for use with training "Product of Expert" models is called *contrastive divergence* (CD).



An RBM with fully connected visible and hidden units. Note there are no hidden-hidden or visible-visible connections.

Contrastive Divergence (CD)

- CD provides an approximation to the maximum likelihood method that would ideally be applied for learning the weights of the RBM with gradient ascent:

$$\Delta w_{ij}(t + 1) = w_{ij}(t) + \eta \frac{\partial \log(p(v))}{\partial w_{ij}}$$

- $p(v)$: the probability of a visible vector, given by $p(v) = \frac{\sum_h \exp(-E(v,h))}{Z}$.
- $E(v, h)$: the energy function assigned to the state of the network,
 - Given by $E(v, h) = -v^T \mathbf{W}h$
 - A lower energy indicates the network is in a more “desirable” configuration.

Contrastive Divergence (CD)

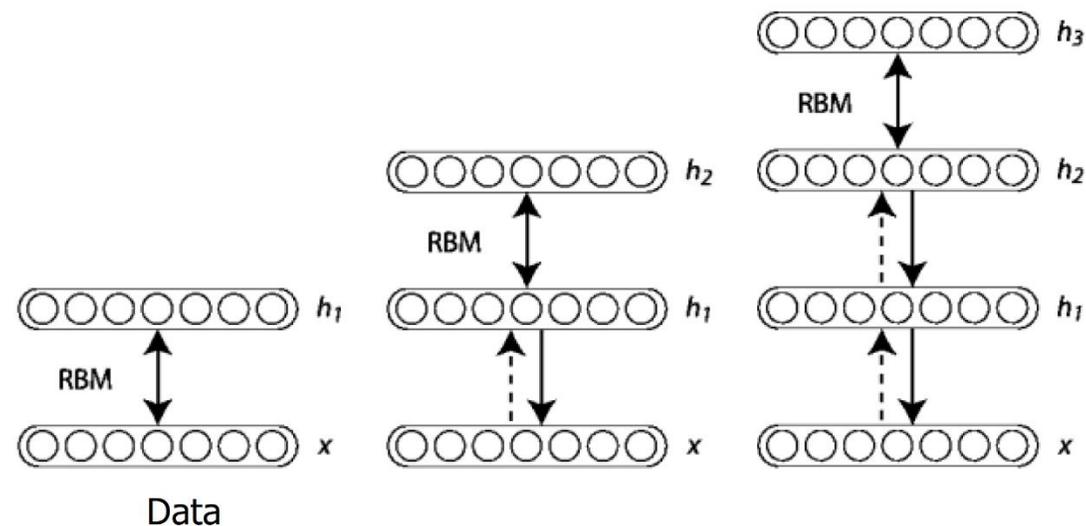
- $\frac{\partial \log(p(v))}{\partial w_{ij}}$ has the simple form $\langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model}$.
 - $\langle \dots \rangle_p$ represent averages with respect to distribution p .
 - Intuition: fitting the data (reducing error)
- Challenge: sampling $\langle v_i h_j \rangle_{model}$ requires alternating Gibbs sampling for a long time.
- CD replaces this step by running alternating Gibbs sampling for n steps. After n steps, the data are sampled and that sample is used in place of $\langle v_i h_j \rangle_{model}$.

The CD procedure

1. Initialize the visible units to a training vector.
2. Update the hidden units in parallel given the visible units:
$$p(h_j = 1|\mathbf{V}) = \sigma(b_j + \sum_i v_i w_{ij})$$
 - $\sigma(\cdot)$: sigmoid function; b_j : the bias of h_j
3. (Reconstruction) Update the visible units in parallel given the hidden units:
$$p(v_i = 1|\mathbf{H}) = \sigma(a_i + \sum_j h_j w_{ij})$$
 - a_i : the bias of v_i
4. Re-update the hidden units in parallel given the reconstructed visible units using the same equation as in step 2.
5. Perform the weight update:
$$\Delta w_{ij} \propto \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{reconstruction}$$

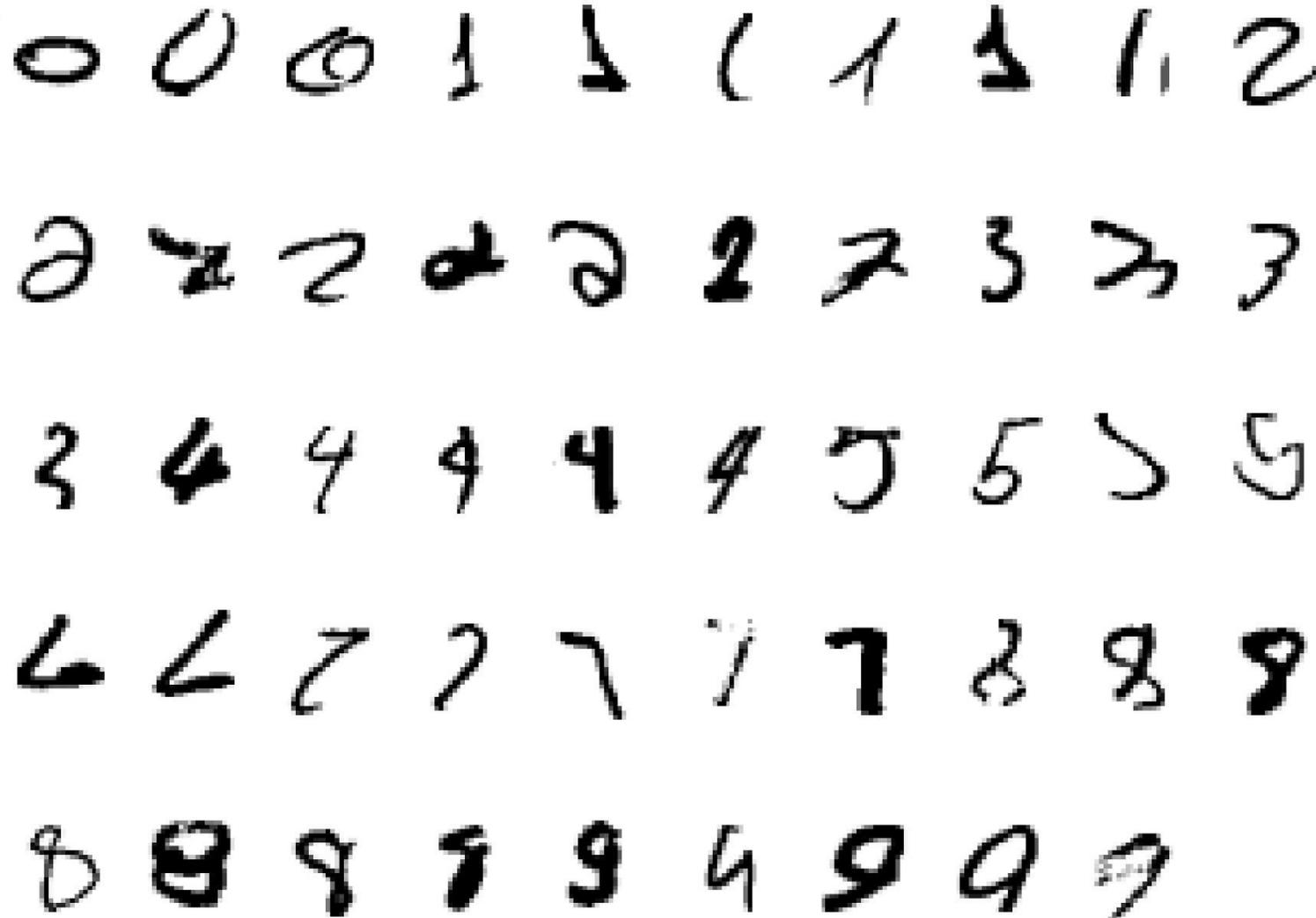
Deep Belief Network Architecture

- Once an RBM is trained, another RBM is "stacked" atop it, taking its input from the final already-trained layer.
- The new visible layer is initialized to a training vector, and values for the units in the already-trained layers are assigned using the current weights and biases.
- The new RBM is then trained with the CD procedure.
- This whole process is repeated until some desired stopping criterion is met.



DBN Example: Hand-written Digit Recognition

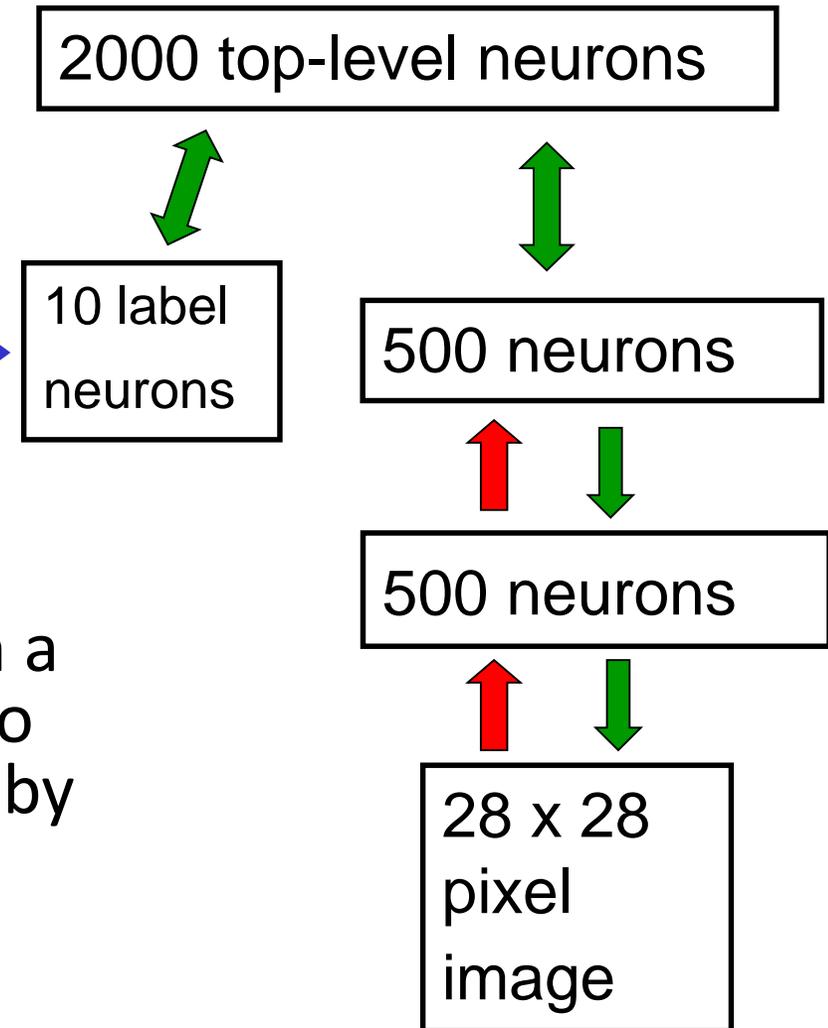
- Input:



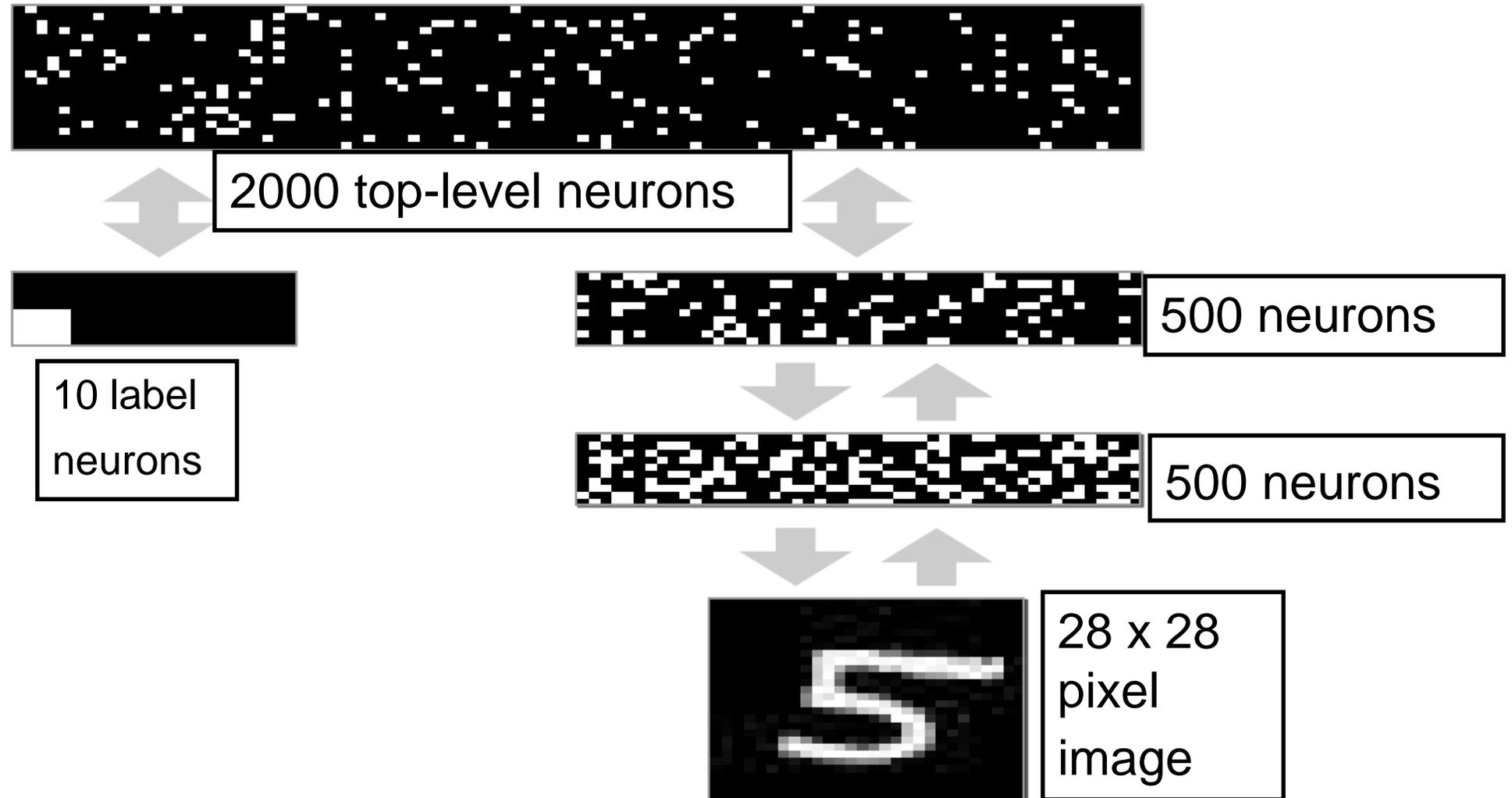
DBN Example: Hand-written Digit Recognition

The top two layers form an associative memory whose energy landscape models the low dimensional manifolds of the digits.

- The model learns to generate combinations of labels and images.
- To perform recognition we start with a neutral state of the label units and do an up-pass from the image followed by a few iterations of the top-level associative memory.



DBN Example: Hand-written Digit Recognition



DBN Example: Hand-written Digit Recognition

How well does it discriminate on MNIST test set (hand-written image) with no extra information about geometric distortions?

- **Generative model based on RBM's** **1.25% (error rate)**
- Support Vector Machine (Decoste et. al.) 1.4%
- Backprop with 1000 hiddens (Platt) ~1.6%
- Backprop with 500 -->300 hiddens ~1.6%
- K-Nearest Neighbor ~ 3.3%
- See Le Cun et. al. 1998 for more results

- Its better than backprop and much more neurally plausible because the neurons only need to send one kind of signal, and the teacher can be another sensory input.

Auto Encoders

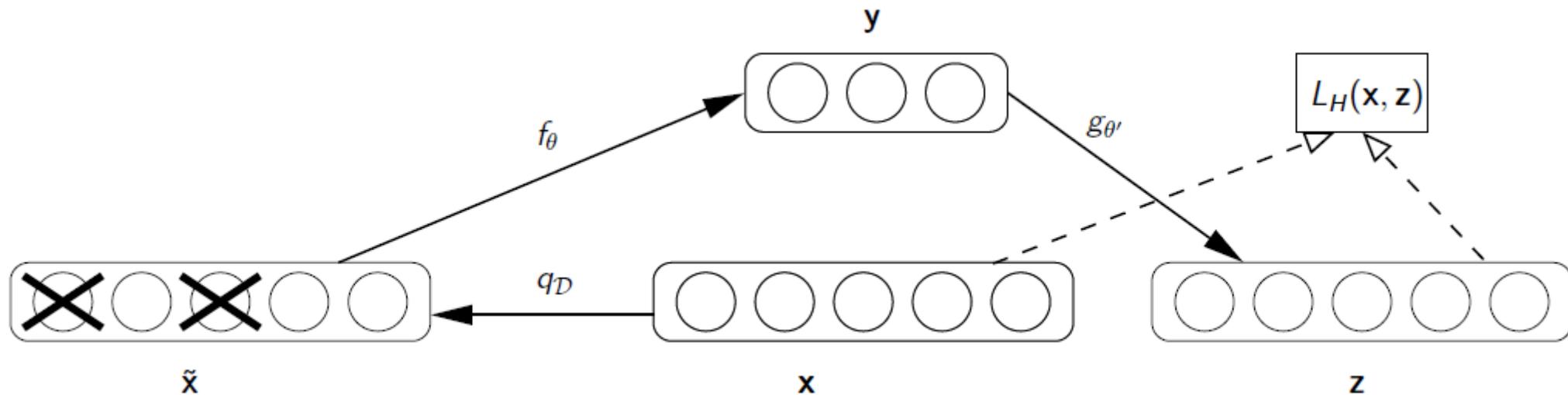
- The *auto encoder* idea is motivated by the concept of a good representation.
 - For example, for a classifier, a good representation can be defined as one that will yield a better performing classifier.
- An *encoder* is a deterministic mapping f_{θ} that transforms an input vector \mathbf{x} into hidden representation \mathbf{y}
 - $\theta = \{\mathbf{W}, b\}$, where \mathbf{W} is the weight matrix and b is bias (an offset vector)
- A *decoder* maps back the hidden representation \mathbf{y} to the reconstructed input \mathbf{z} via g_{θ} .
- Auto encoding: compare the reconstructed input \mathbf{z} to the original input \mathbf{x} and try to minimize this error to make \mathbf{z} as close as possible to \mathbf{x} .

De-noising Auto Encoders

- In Vincent et al. (2010), “*a **good representation** is one that can be obtained **robustly** from a **corrupted input** and that will be useful for **recovering** the corresponding **clean input**.*”
 - The higher level representations are relatively stable and robust to input corruption.
 - It is necessary to extract features that are useful for representation of the input distribution.
- In de-noising auto encoders, the partially *corrupted* output is cleaned (de-noised).

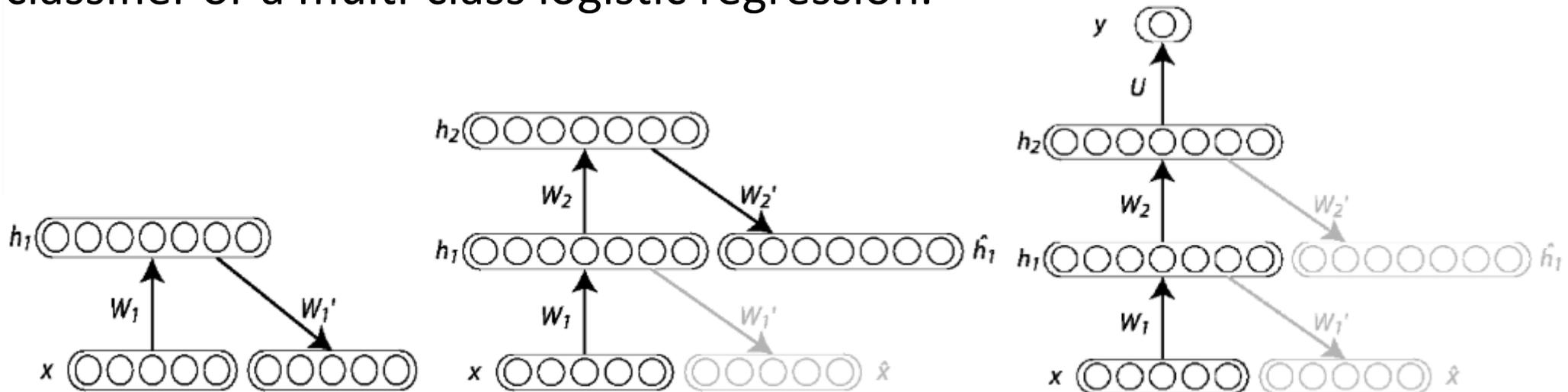
De-noising Auto Encoders Procedure

1. (Corrupting) Clean input \mathbf{x} is partially corrupted through a stochastic mapping $q_D(\tilde{\mathbf{x}}|\mathbf{x})$, yielding corrupted input $\tilde{\mathbf{x}} \sim q_D(\tilde{\mathbf{x}}|\mathbf{x})$.
2. The corrupted input $\tilde{\mathbf{x}}$ passes through a basic auto encoder and is mapped to a hidden representation $\mathbf{y} = f_\theta(\tilde{\mathbf{x}}) = s(\mathbf{W}\tilde{\mathbf{x}} + b)$.
3. From this hidden representation, we can reconstruct $\mathbf{z} = g_\theta(\mathbf{y})$.
4. Minimize the reconstruction error $L_H(\mathbf{x}, \mathbf{z})$.
 - $L_H(\mathbf{x}, \mathbf{z})$ choices: cross-entropy loss or squared error loss



Stacked De-noising Auto Encoder Architecture

- In order to make a deep architecture, auto encoders stack one on top of another.
- Once the encoding function f_θ of the first denoising auto encoder is learned and used to reconstruct the corrupted input, we can train the second level.
- Once the stacked auto encoder is trained, its output can be used as the input to a supervised learning algorithm such as support vector machine classifier or a multi-class logistic regression.

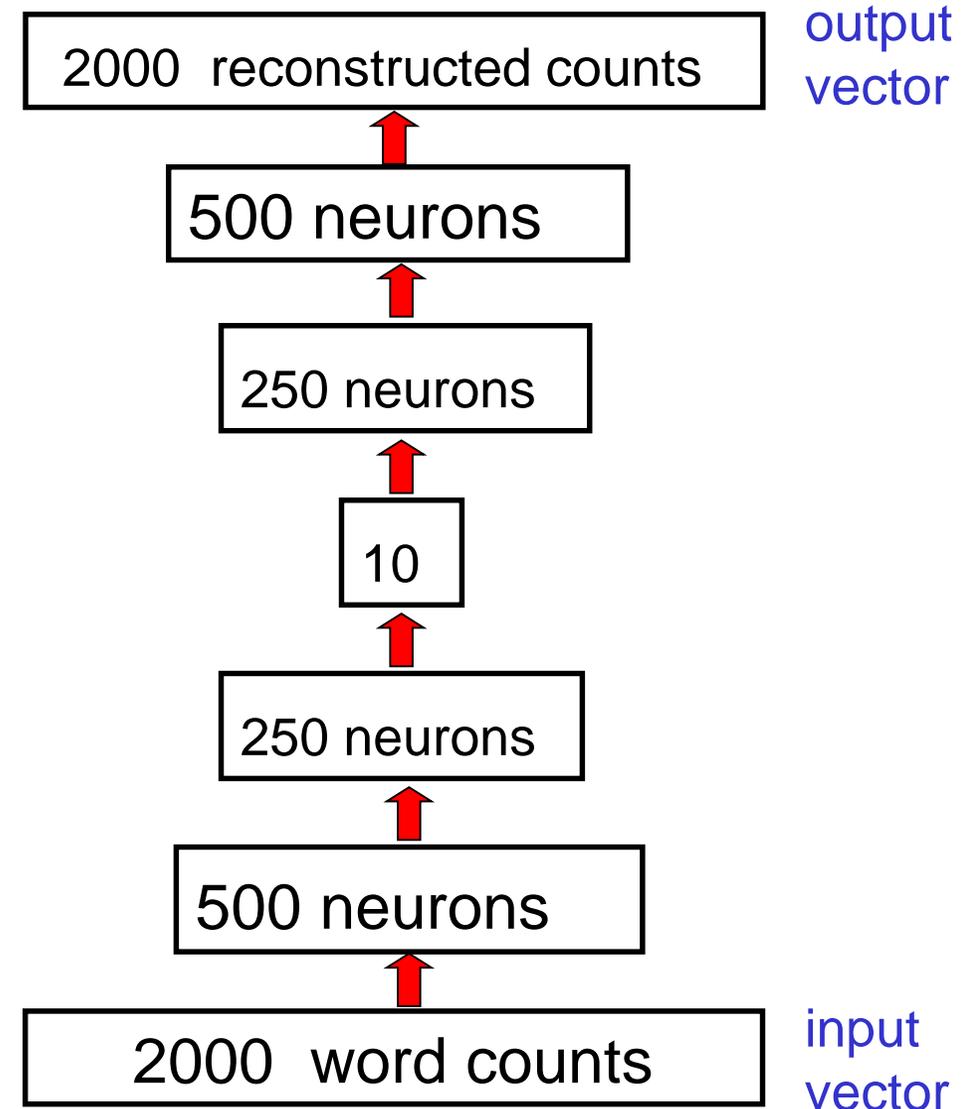


Example: Information Retrieval

- We can use an auto encoder to find low-dimensional codes for documents that allow fast and accurate retrieval of similar documents from a large set.
- We start by converting each document into a “bag of words”. This a 2000 dimensional vector that contains the counts for each of the 2000 commonest words.

Example: Information Retrieval

- We train the neural network to reproduce its input vector as its output
- This forces it to compress as much information as possible into the 10 numbers in the central bottleneck.
- These 10 numbers are then a good way to compare documents.

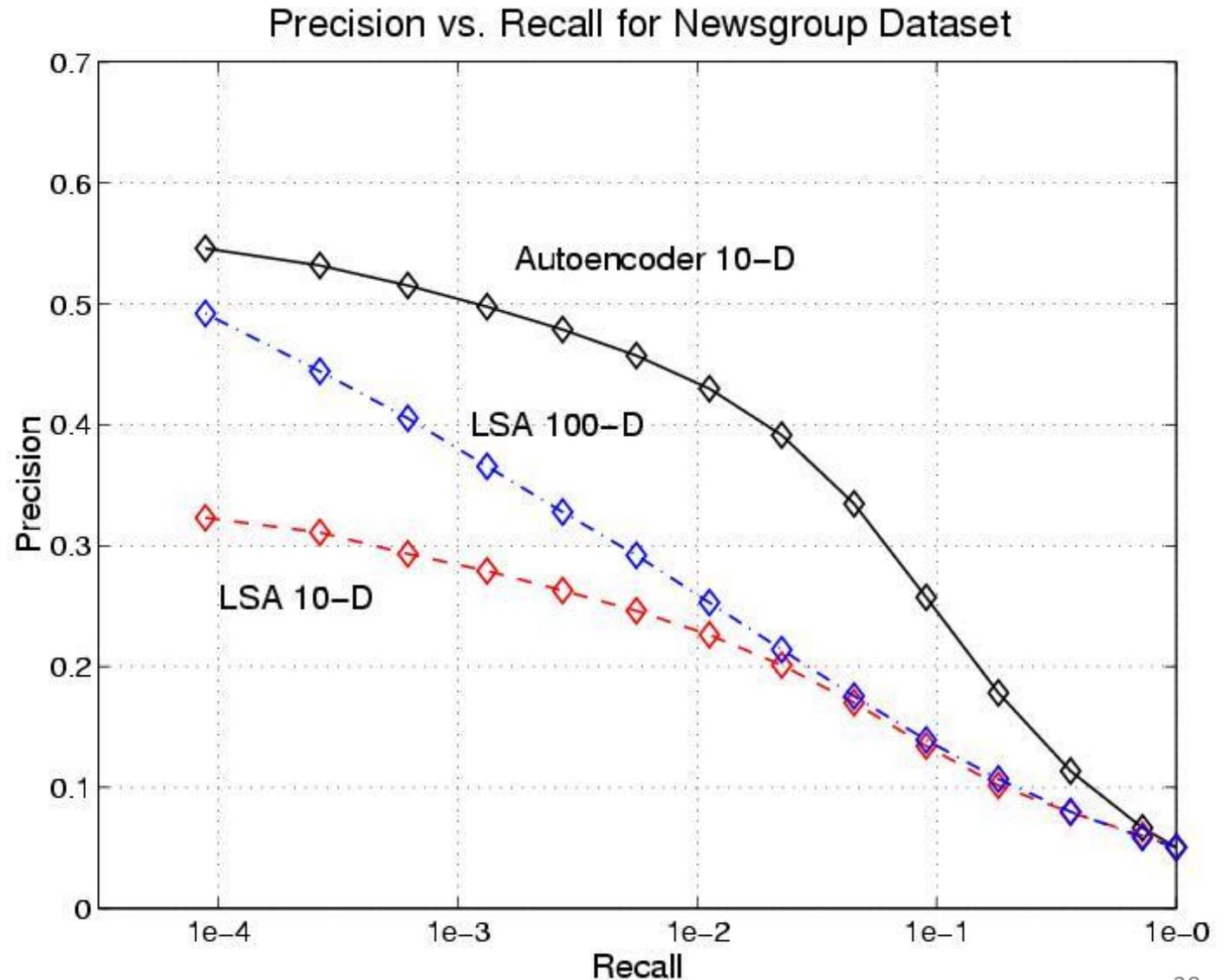


Example: Information Retrieval

- Train on bags of 2000 words for 400,000 training cases of business documents.
 - First train a stack of RBM's. Then fine-tune with backprop.
- Test on a separate 400,000 documents.
 - Pick one test document as a query. Rank order all the other test documents by using the cosine of the angle between codes.
 - Repeat this using each of the 400,000 test documents as the query (requires 0.16 trillion comparisons).
- Plot the number of retrieved documents against the proportion that are in the same hand-labeled class as the query document.

Example: Information Retrieval

- The shortlist found using binary codes actually improves the precision-recall curves of TF-IDF.
 - Locality sensitive hashing (the fastest other method) is 50 times slower and has worse precision-recall curves.
 - Vs. Latent Semantic Analysis (LSA)

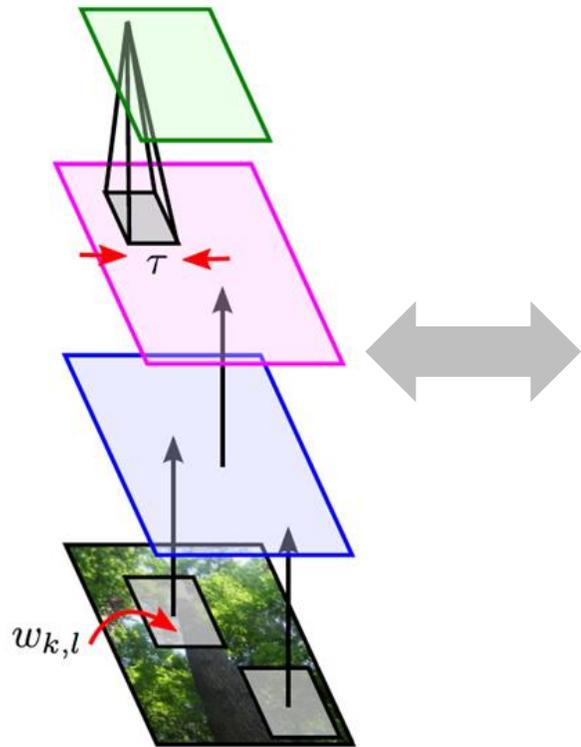


Convolutional Neural Network (CNN)

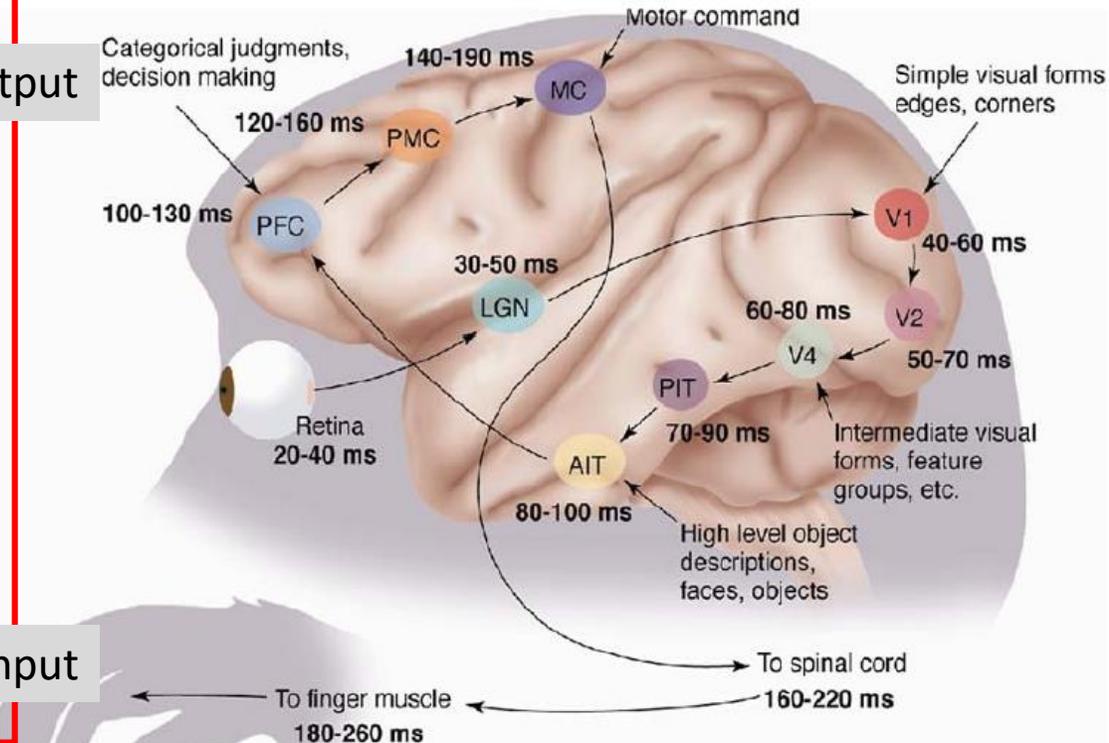
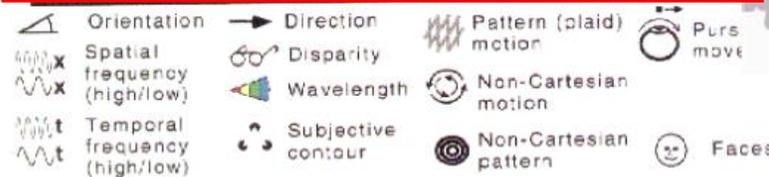
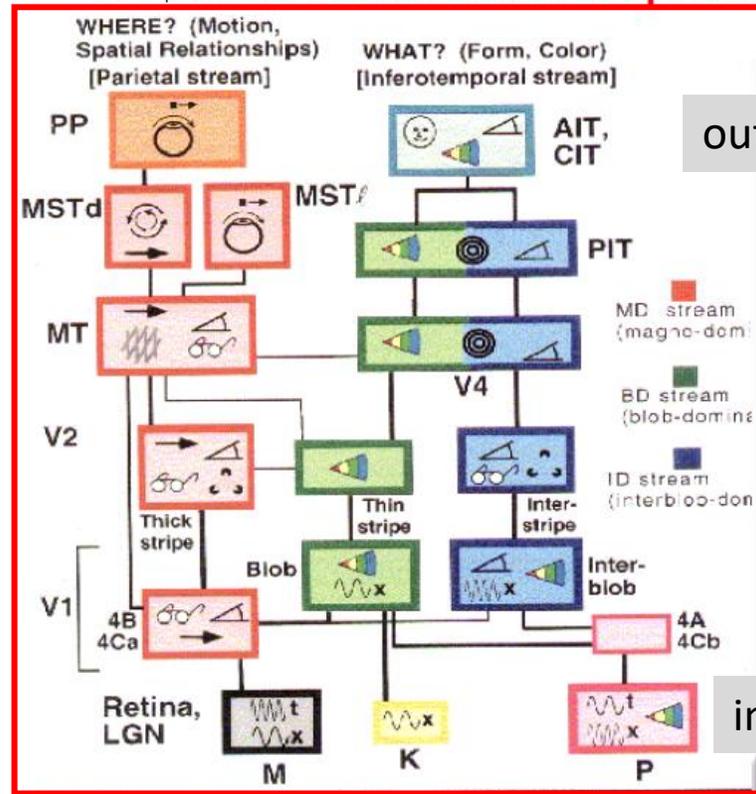
- Convolutional Neural Networks are inspired by mammalian visual cortex.
 - The visual cortex contains a complex arrangement of cells, which are sensitive to small sub-regions of the visual field, called a receptive field. These cells act as local filters over the input space and are well-suited to exploit the strong spatially local correlation present in natural images.
 - Two basic cell types:
 - Simple cells respond maximally to specific edge-like patterns within their receptive field.
 - Complex cells have larger receptive fields and are locally invariant to the exact position of the pattern.

The Mammalian Visual Cortex Inspires CNN

Convolutional Neural Net



- The ventral (recognition) pathway in the visual cortex has multiple stages
- Retina - LGN - V1 - V2 - V4 - PIT - AIT
- Lots of intermediate representations



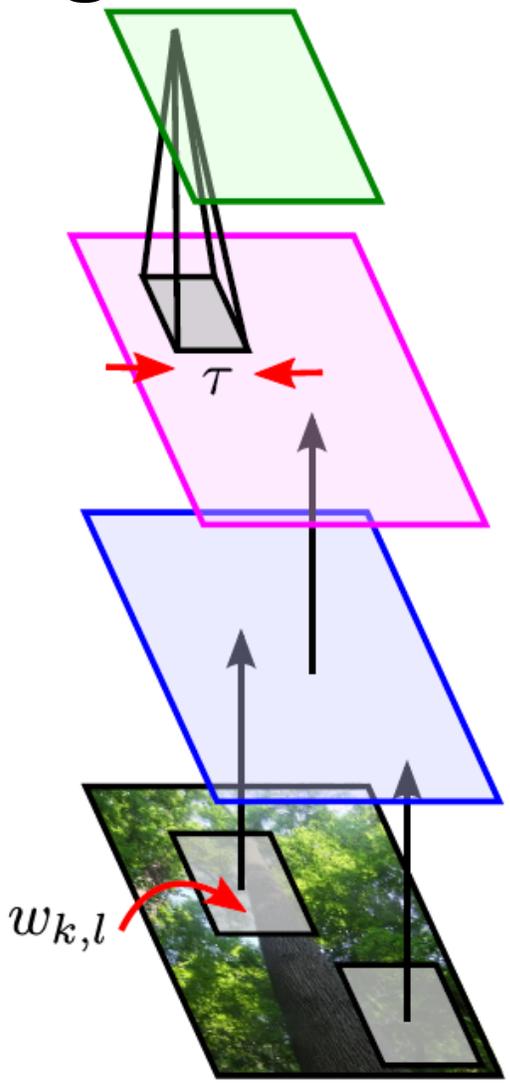
[picture from Simon Thorpe]

[Gallant & Van Essen]

CNN Architecture

- Intuition: Neural network with specialized connectivity structure,
 - Stacking multiple layers of feature extractors
 - Low-level layers extract local features.
 - High-level layers extract learn global patterns.
- A CNN is a list of layers that transform the input data into an output class/prediction.
- There are a few distinct types of layers:
 - Convolutional layer
 - Non-linear layer
 - Pooling layer

Building-blocks for CNN's



$$x_{i,j} = \max_{|k| < \tau, |l| < \tau} y_{i-k, j-l}$$

mean or subsample also used

pooling stage

Feature maps of a larger region are combined.

$$y_{i,j} = f(a_{i,j})$$

e.g. $f(a) = [a]_+$
 $f(a) = \text{sigmoid}(a)$

non-linear stage

Feature maps are trained with neurons.

$$a_{i,j} = \sum_{k,l} w_{k,l} z_{i-k, j-l}$$

Shared weights

convolutional stage

Each sub-region yields a feature map, representing its feature.

$z_{i,j}$ Images are segmented into sub-regions.

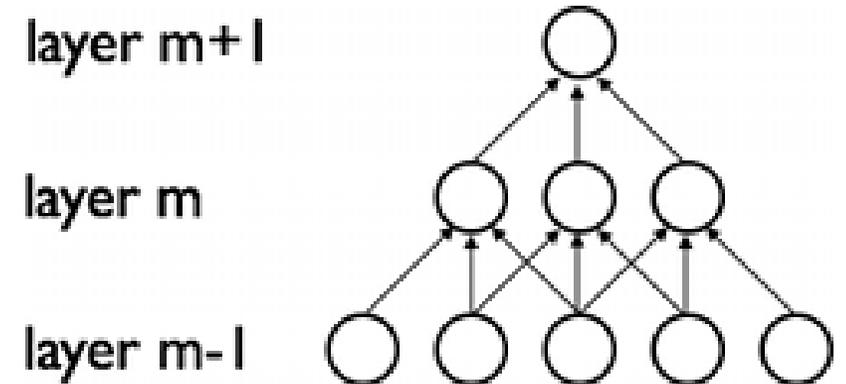
input image

CNN Architecture: Convolutional Layer

- The core layer of CNNs
- The convolutional layer consists of a set of filters.
 - Each filter covers a spatially small portion of the input data.
- Each filter is convolved across the dimensions of the input data, producing a multidimensional feature map.
 - As we convolve the filter, we are computing the dot product between the parameters of the filter and the input.
- Intuition: the network will learn filters that activate when they see some specific type of feature at some spatial position in the input.
- The key architectural characteristics of the convolutional layer is local connectivity and shared weights.

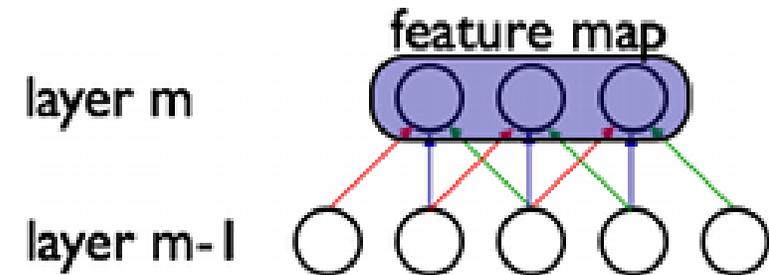
CNN Convolutional Layer: Local Connectivity

- Neurons in layer **m** are only connected to 3 adjacent neurons in the **m-1** layer.
- Neurons in layer **m+1** have a similar connectivity with the layer below.
- Each neuron is unresponsive to variations outside of its *receptive field* with respect to the input.
 - Receptive field: small neuron collections which process portions of the input data
- The architecture thus ensures that the learnt feature extractors produce the strongest response to a spatially local input pattern.



CNN Convolutional Layer: Shared Weights

- We show 3 hidden neurons belonging to the same feature map (the layer right above the input layer).
- Weights of the same color are shared—constrained to be identical.
- Gradient descent can still be used to learn such shared parameters, with only a small change to the original algorithm.
- The gradient of a shared weight is simply the sum of the gradients of the parameters being shared.
- Replicating neurons in this way allows for features to be detected regardless of their position in the input.
- Additionally, weight sharing increases learning efficiency by greatly reducing the number of free parameters being learnt.

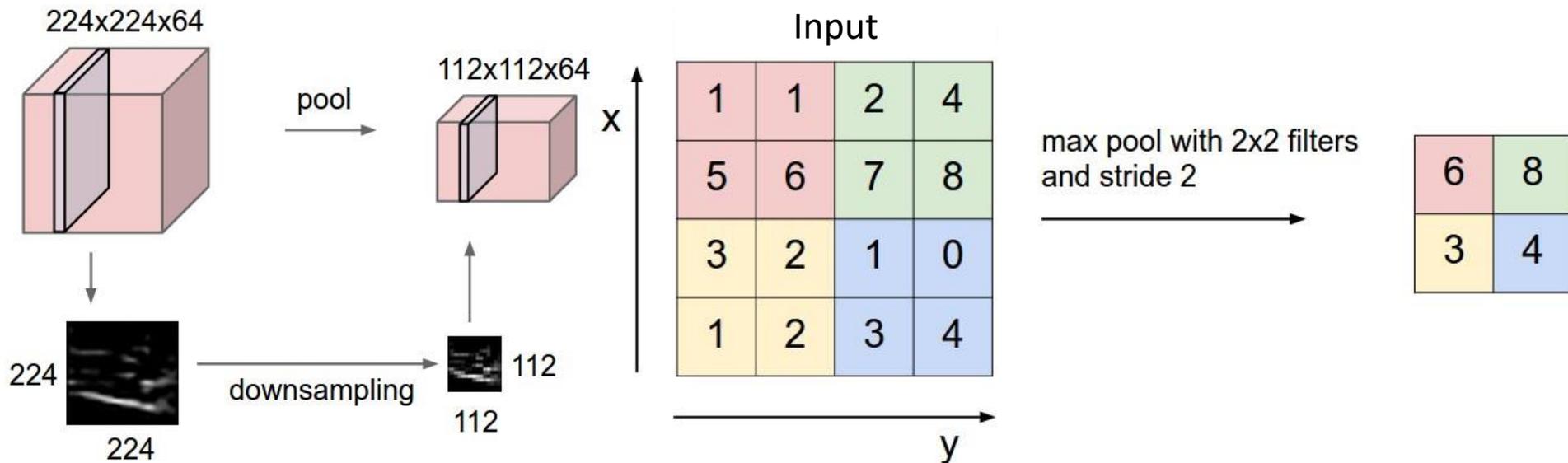


CNN Architecture: Non-linear Layer

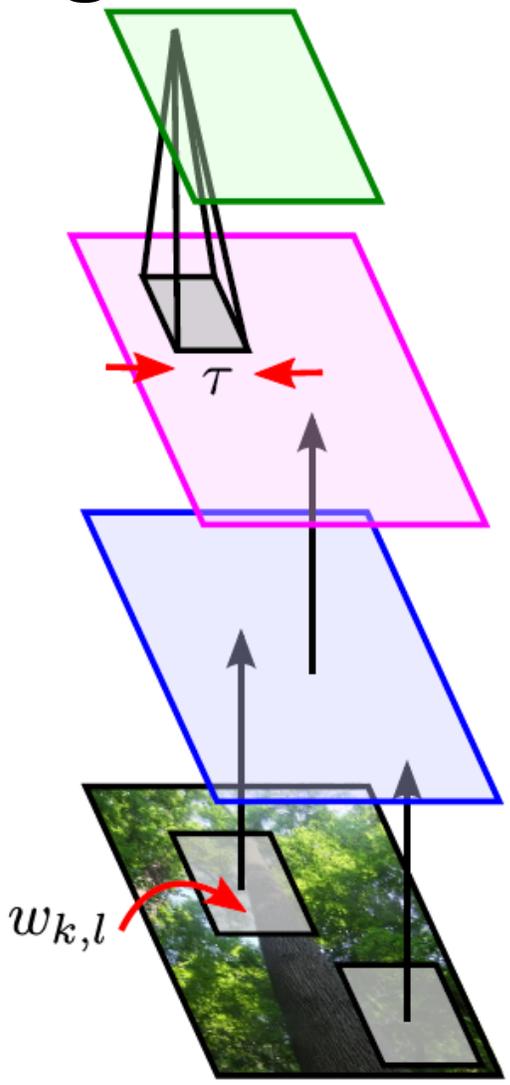
- Intuition: Increase the nonlinearity of the entire architecture without affecting the receptive fields of the convolution layer
- A layer of neurons that applies the non-linear activation function, such as,
 - $f(x) = \max(0, x)$
 - $f(x) = \tanh x$
 - $f(x) = |\tanh x|$
 - $f(x) = (1 + e^{-x})^{-1}$

CNN Architecture: Pooling Layer

- Intuition: to progressively reduce the spatial size of the representation to reduce the amount of parameters and computation in the network, and hence to also control overfitting
- Pooling partitions the input image into a set of non-overlapping rectangles and, for each such sub-region, outputs the maximum value of the features in that region.



Building-blocks for CNN's



$$x_{i,j} = \max_{|k| < \tau, |l| < \tau} y_{i-k, j-l}$$

mean or subsample also used

pooling stage

Feature maps of a larger region are combined.

$$y_{i,j} = f(a_{i,j})$$

e.g. $f(a) = [a]_+$
 $f(a) = \text{sigmoid}(a)$

non-linear stage

Feature maps are trained with neurons.

$$a_{i,j} = \sum_{k,l} w_{k,l} z_{i-k, j-l}$$

Shared weights

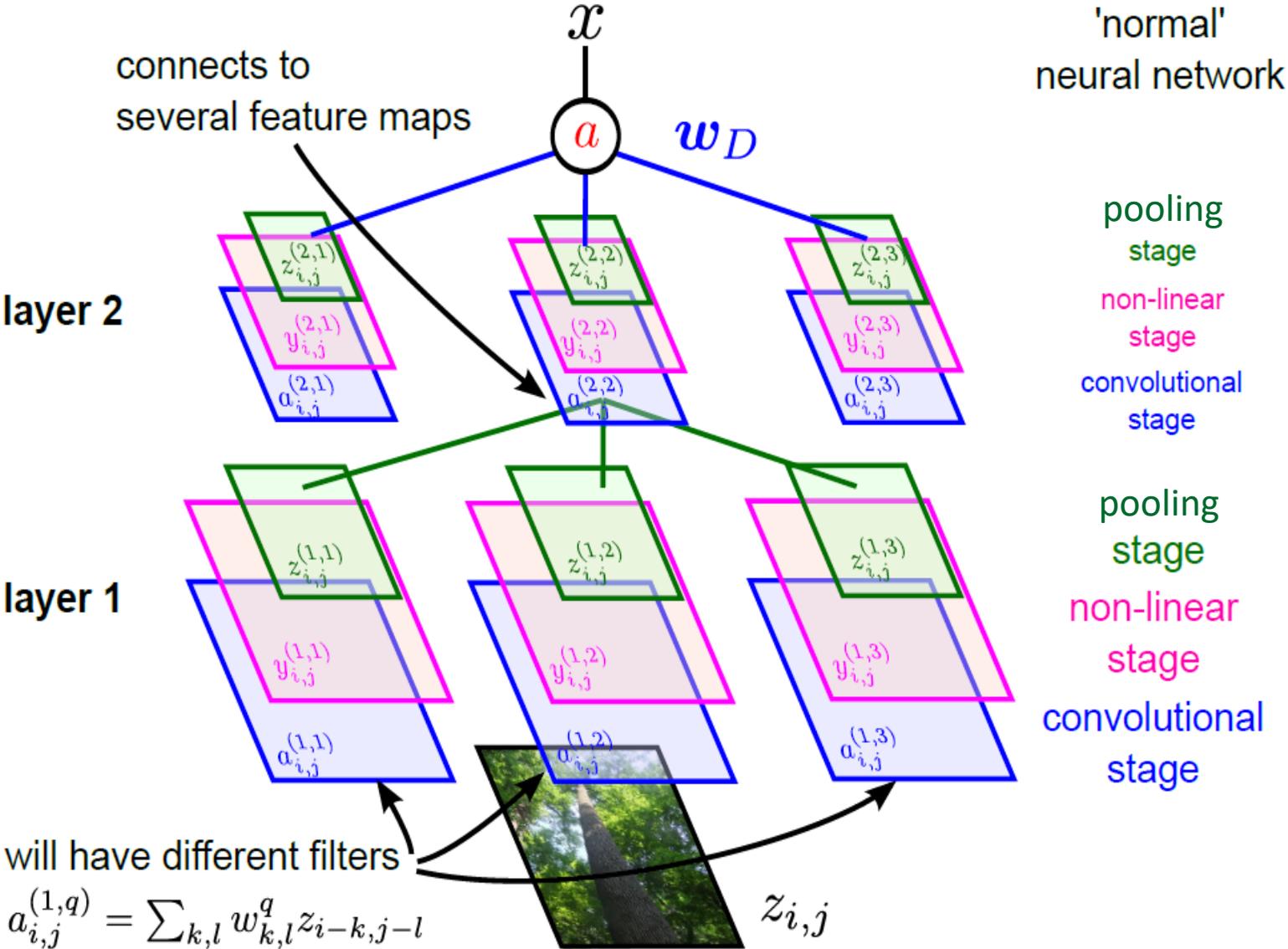
convolutional stage

Each sub-region yields a feature map, representing its feature.

$z_{i,j}$ Images are segmented into sub-regions.

input image

Full CNN



Outline

- Introduction
 - Motivation
 - Deep Learning Intuition
- Neural Network
 - Neuron, Feedforward algorithm, and Backpropagation algorithm
 - Limitations
- Deep Learning
 - Deep Belief Network: Autoencoder & Restricted Boltzman Machine
 - Convolutional Neural Network
- **Deep Learning for Text Mining**
 - Recurrent Neural Network
 - Recursive Neural Network
 - Word Embedding
- Conclusions

Deep Learning for Text Mining

Building deep learning models on textual data requires:

- Representation of the basic text unit, word.
- Neural network structure that can hierarchically capture the sequential nature of text.

Deep learning models for text mining use:

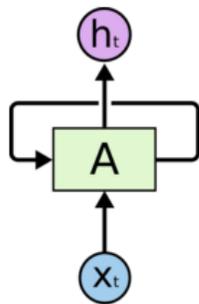
- Vector representation of words (i.e., word embeddings)
- Neural network structures
 - *Recurrent* Neural Network
 - *Recursive* Neural Network

Recurrent Neural Networks

- The applications of standard Neural Networks (and also Convolutional Networks) are limited due to:
 - They only accepted a fixed-size vector as input (e.g., an image) and produce a fixed-size vector as output (e.g., probabilities of different classes).
 - These models use a fixed amount of computational steps (e.g. the number of layers in the model).
- Recurrent Neural Networks are unique as they allow us to operate over sequences of vectors.
 - Sequences in the input, the output, or in the most general case both

Recurrent Neural Networks

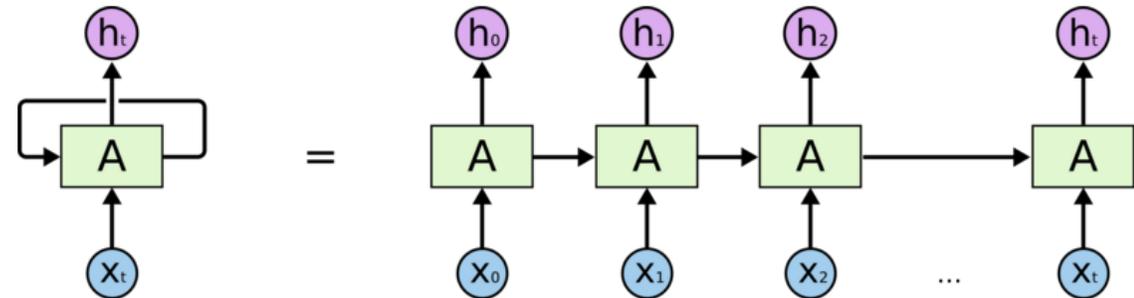
- Recurrent Neural Networks are networks with loops in them, allowing information to persist.



Recurrent Neural Networks have loops.

In the above diagram, a chunk of neural network, A , looks at some input x_t and outputs a value h_t .

A loop allows information to be passed from one step of the network to the next.



An unrolled recurrent neural network.

A recurrent neural network can be thought of as multiple copies of the same network, each passing a message to a successor.

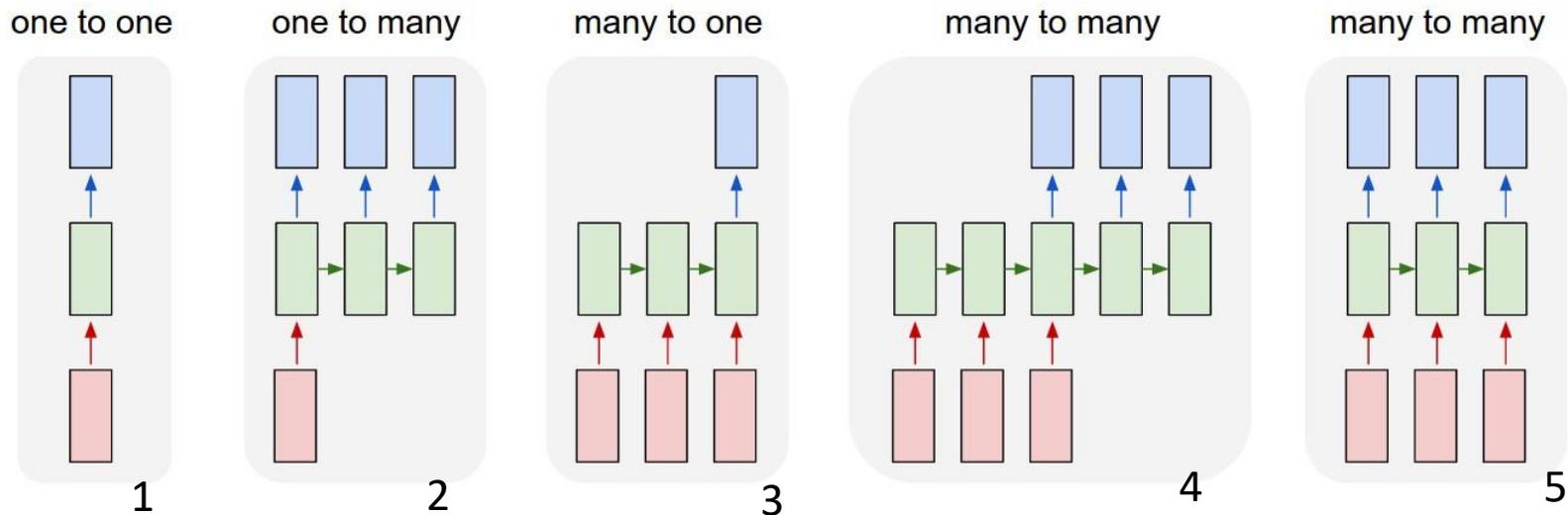
The diagram above shows what happens if we unroll the loop.

Recurrent Neural Networks

- Intuition of Recurrent Neural Networks
 - Human thoughts have persistence; humans don't start their thinking from scratch every second.
 - As you read this sentence, you understand each word based on your understanding of previous words.
 - One of the appeals of RNNs is the idea that they are able to connect previous information to the present task
 - E.g., using previous video frames to inform the understanding of the present frame.
 - E.g., a language model tries to predict the next word based on the previous ones.

Recurrent Neural Networks

- Examples of Recurrent Neural Networks



- Each rectangle is a vector and arrows represent functions (e.g. matrix multiply).
- Input vectors are in red, output vectors are in blue and green vectors hold the RNN's state

(1) Standard mode of processing without RNN, from fixed-sized input to fixed-sized output (e.g. image classification).

(2) Sequence output (e.g. image captioning takes an image and outputs a sentence of words).

(3) Sequence input (e.g. sentiment analysis where a given sentence is classified as expressing positive or negative sentiment).

(4) Sequence input and sequence output (e.g. Machine Translation: an RNN reads a sentence in English and then outputs a sentence in French).

(5) Synced sequence input and output (e.g. video classification where we wish to label each frame of the video).

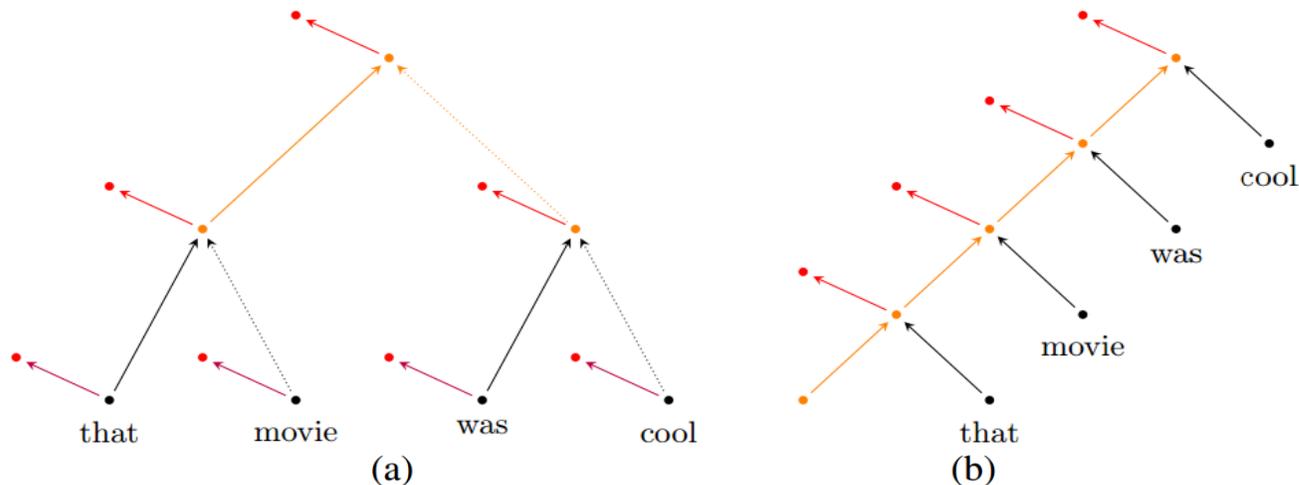
Recurrent Neural Networks

- Incredible success applying RNNs to language modeling and sequence learning problems

Task	Input Sequence	Output Sequence
Machine translation (Sutskever et al. 2014)	English	French
Question answering (Bordes et al. 2014)	Question	Answer
Speech recognition (Graves et al. 2013)	Voice	Text
Handwriting prediction (Graves 2013)	Handwriting	Text
Opinion mining (Irsoy et al. 2014)	Text	Opinion expression

Recursive Neural Networks

- A recursive NN can be seen as a generalization of the recurrent NN
 - A recurrent neural network is in fact a recursive neural network with the structure of a linear chain
 - Recursive NNs operate on hierarchical structure, Recurrent NN operate on progression of time



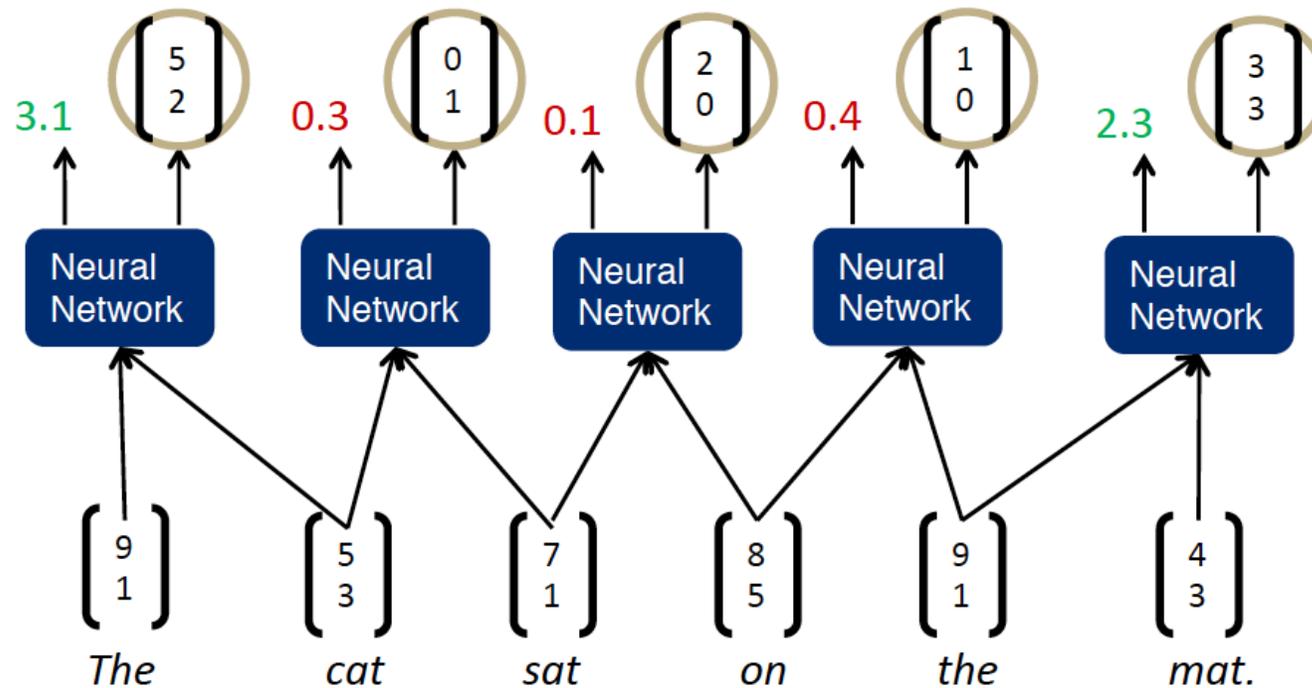
- Operation of a recursive net (a), and a recurrent net (b) on an example sentence. Note the linear chain in (b)
- Black, orange and red dots represent input, hidden and output layers, respectively.
- Directed edges having the same color-style combination denote shared connections.

- Applied to parsing, sentence-level sentiment analysis, and paraphrase detection.

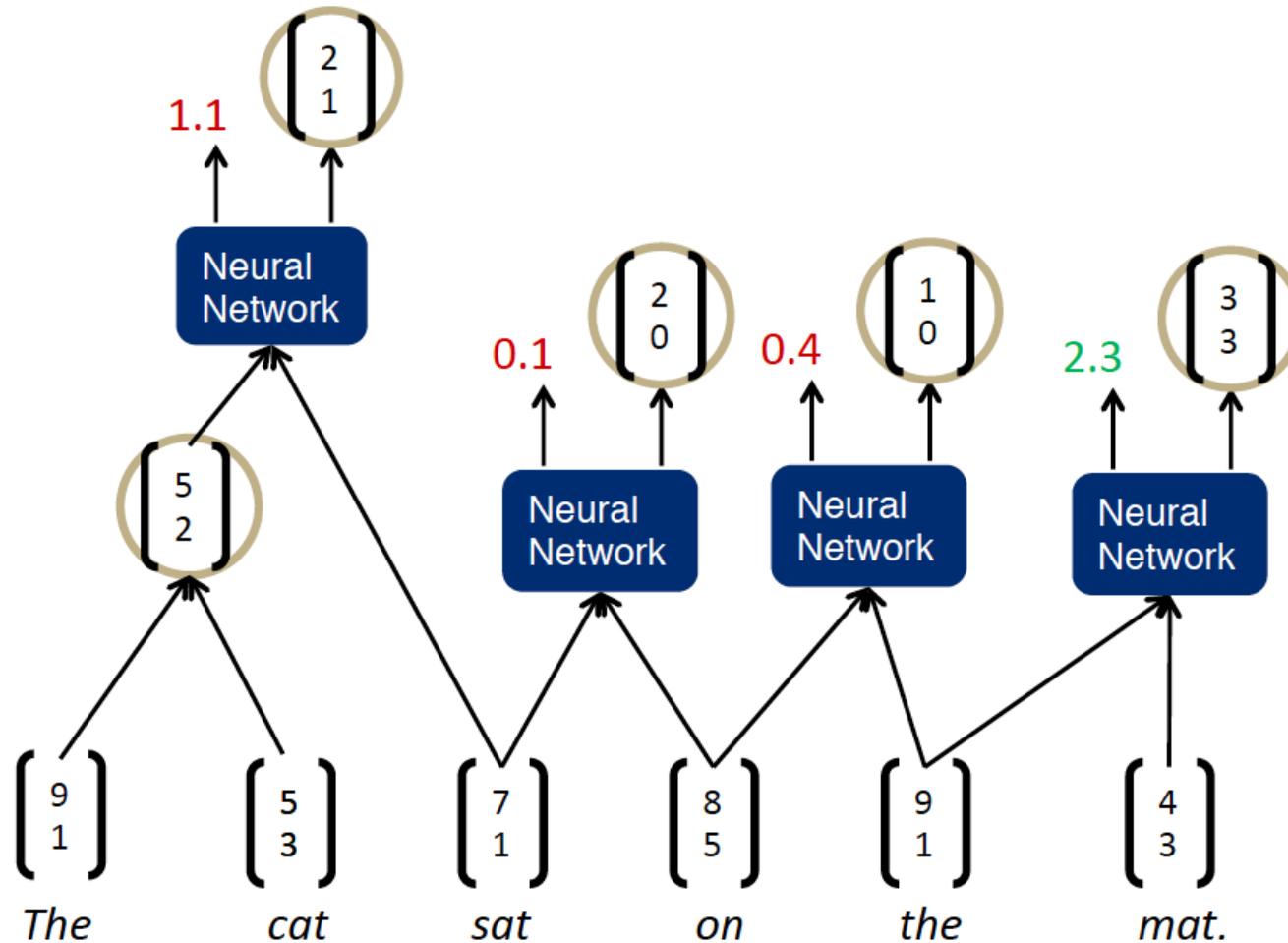
Recursive Neural Networks

- Given the structural representation of a sentence, e.g. a parse tree:
 - Recursive Neural Networks recursively generate parent representations in a bottom-up fashion,
 - By combining tokens to produce representations for phrases, eventually producing the whole sentence.
 - The sentence-level representation (or, alternatively, its phrases) can then be used to make a final classification for a given input sentence — e.g. whether it conveys a positive or a negative sentiment.

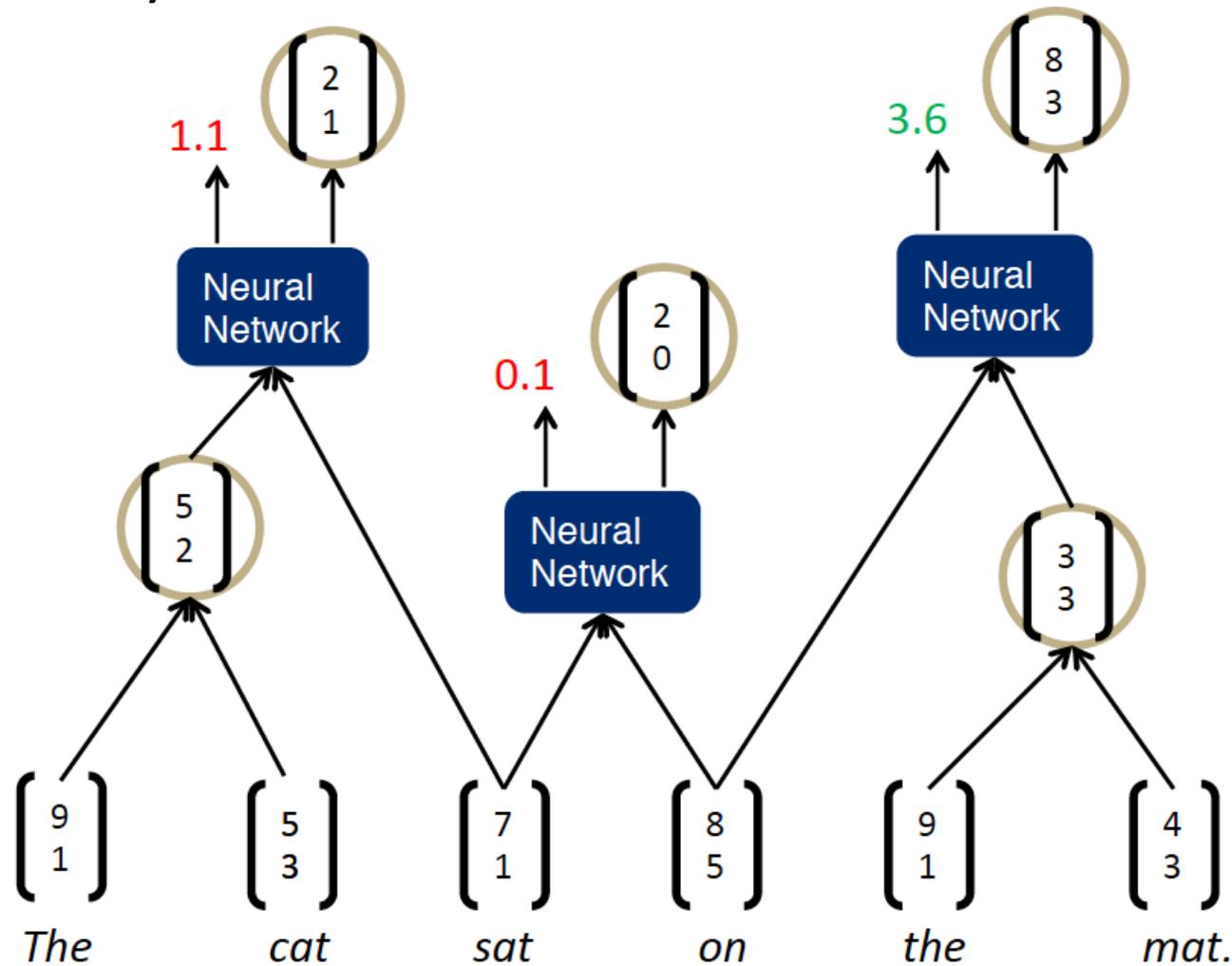
Recursive Neural Network Illustration (Animation)



Recursive Neural Network Illustration (Animation)



Recursive Neural Network Illustration (Animation)



Recursive Neural Networks

- SemEval 2014: Predicting semantic relatedness of two sentences (Tai et al. 2015)

Ranking by mean word vector cosine similarity	Score	Ranking by Dependency Tree LSTM model	Score
a woman is slicing potatoes		a woman is slicing potatoes	
a woman is cutting potatoes	0.96	a woman is cutting potatoes	4.82
a woman is slicing herbs	0.92	potatoes are being sliced by a woman	4.70
a woman is slicing tofu	0.92	tofu is being sliced by a woman	4.39
a boy is waving at some young runners from the ocean		a boy is waving at some young runners from the ocean	
a man and a boy are standing at the bottom of some stairs , which are outdoors	0.92	a group of men is playing with a ball on the beach	3.79
a group of children in uniforms is standing at a gate and one is kissing the mother	0.90	a young boy wearing a red swimsuit is jumping out of a blue kiddies pool	3.37
a group of children in uniforms is standing at a gate and there is no one kissing the mother	0.90	the man is tossing a kid into the swimming pool that is near the ocean	3.19
two men are playing guitar		two men are playing guitar	
some men are playing rugby	0.88	the man is singing and playing the guitar	4.08
two men are talking	0.87	the man is opening the guitar for donations and plays with the case	4.01
two dogs are playing with each other	0.87	two men are dancing and singing in front of a crowd	4.00

Recursive NN based model can pick up more subtle semantic relatedness in sentences compared to word vector model, e.g., ocean and beach.

Table 4: Most similar sentences from a 1000-sentence sample drawn from the SICK test set. The Tree-LSTM model is able to pick up on more subtle relationships, such as that between “beach” and “ocean” in the second example.

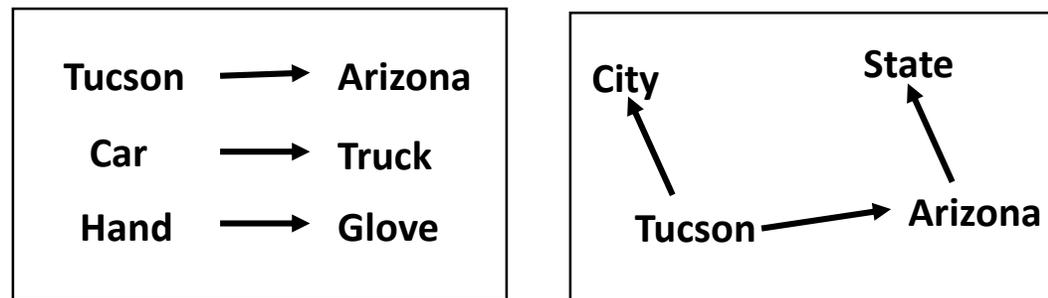
Vector Representation of Words

- Vector space models (VSMs) represent (embed) words in a continuous vector space
 - Theoretical foundation in Linguistics: Distributional Hypothesis
 - Words with similar meanings will occur with similar neighbors if enough text material is available (Rubenstein et al. 1967).
- Approaches that leverage VSMs can be divided into two categories

Approach	Example	Description
Count-based methods	Latent semantic analysis	Compute how often some word co-occurs with its neighbor words in a large text corpus, and then map these count-statistics down to a small, dense vector for each word
Predictive methods	Neural probabilistic language model	Directly predict a word from its neighbors in terms of learned small, dense embedding vectors (considered parameters of the model)

Word2vec –Vector Representation of Words (Mikolov et al. 2013)

- Word2vec: computationally-efficient, 2-layer predictive NN for learning word embeddings from raw text
 - Considered deep for its ability to digest expansive data sets quickly
- Can be used for unsupervised learning of words
 - Relationships between different words
 - Ability to abstract higher meaning between words (e.g., Tucson is a city in the state of Arizona)



- Useful for language modeling, sentiment analysis, and more

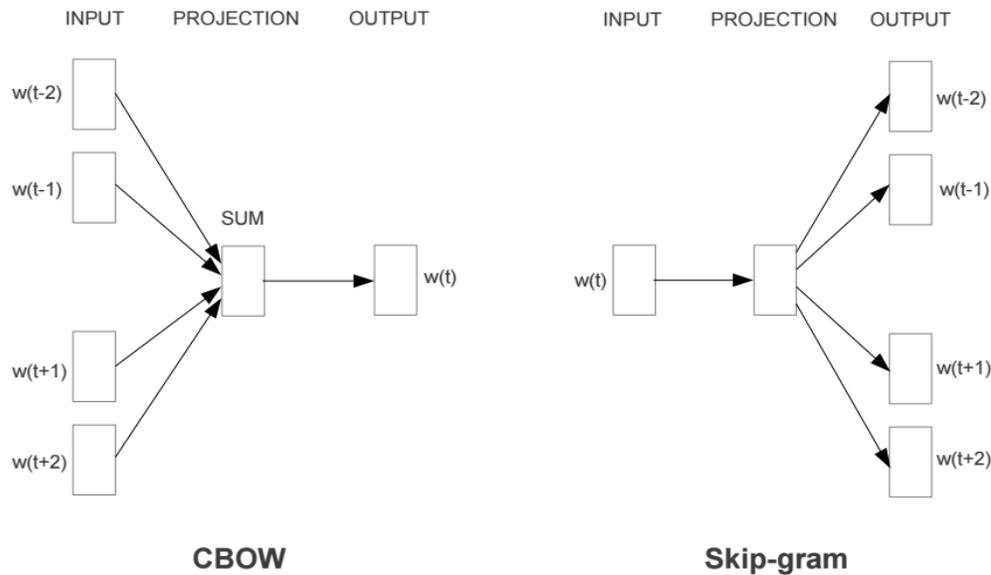
Word2vec –Vector Representation of Words (Mikolov et al. 2013)

- Its input is a text corpus and its output is a set of vectors or “embeddings” (feature vectors for words in that corpus)
 - Similarity between two embeddings represents conceptual similarity of words
- Example results: words associated with *Sweden*, in order of proximity:

Word	Cosine distance
norway	0.760124
denmark	0.715460
finland	0.620022
switzerland	0.588132
belgium	0.585835
netherlands	0.574631
iceland	0.562368
estonia	0.547621
slovenia	0.531408

Word2vec –Vector Representation of Words (Mikolov et al. 2013)

- Word2vec comes with two models:

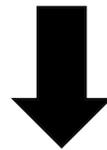


Model	Approach	Speed and Performance	Use case
Continuous Bag-of-Words model (CBOW)	The CBOW predicts the current word based on the context.	Faster to train than the skip-gram model	Predicts frequent words better
Skip-Gram model	Skip-gram predicts surrounding words given the current word.	Usually performs better than CBOW	Predicts rare words better

Word2vec – Vector Representation of Words (Mikolov et al. 2013)

- Skip-gram learning:
 - Given w_0 , predict w_{-2} , w_{-1} , w_1 , and w_2

w_{-2}	w_{-1}	w_0	w_1	w_2
?	?	Network	?	?



w_{-2}	w_{-1}	w_0	w_1	w_2
Recurrent	Neural		Language	Model

- Conversely, CBOW tries to predict w_0 when given w_{-2} , w_{-1} , w_1 , and w_2

Running Word2Vec

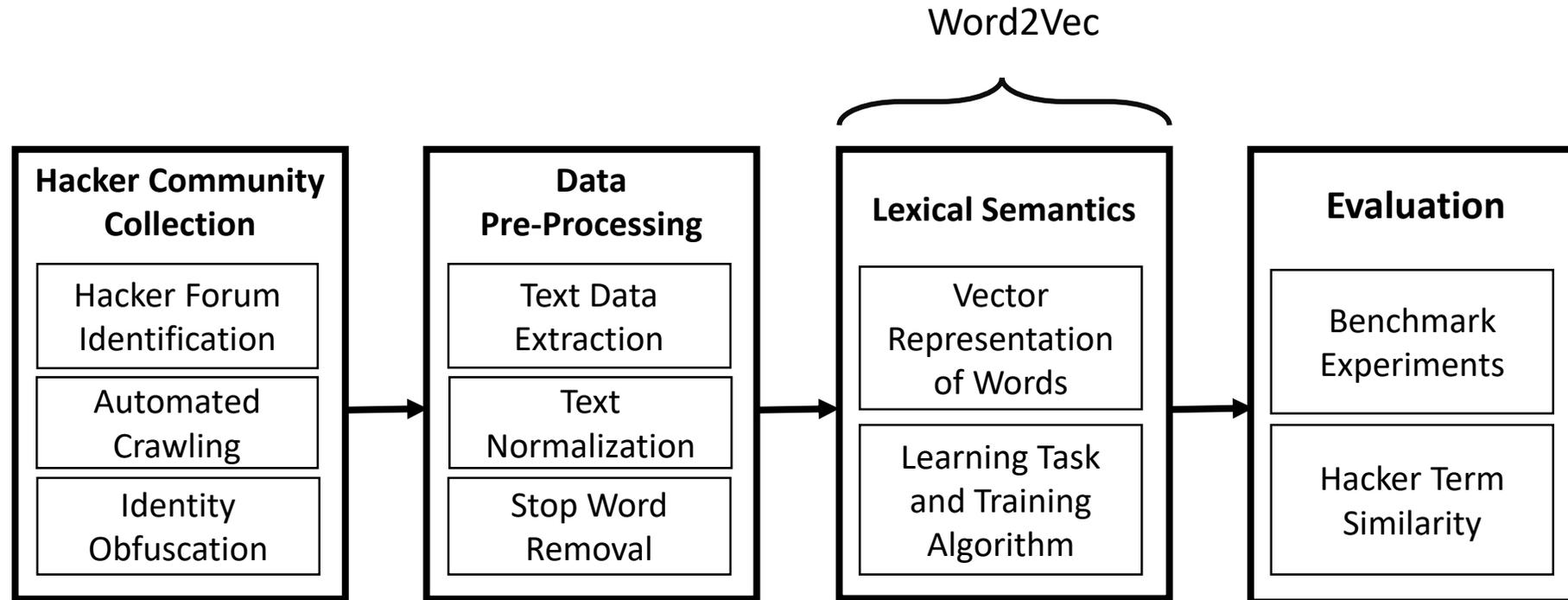
- Download: <https://code.google.com/archive/p/word2vec/>
- Word2Vec comes bundled with many files. Two important ones:
 - *Word2vec.c* - the actual Word2Vec program written in C; is executed in command line
 - *Demo-word.sh* - shell script containing example of how to run *Word2Vec.c* on test data
- To use Word2Vec, you need:
 1. A corpus (e.g., collection of tweets, news articles, product reviews)
 - Word2Vec expects a sequence of sentences as input
 - One input file containing many sentences, with one sentence per line
 2. A C programming language compiler
 - Unix environments are easiest - Linux generally ships with 'gcc' pre-installed, OSX can use Xcode
- Word2Vec is useful for language modeling tasks,

Word2Vec Example

- Example: *Zeus* refers to a botnet and not Greek mythology
- Word2Vec can provide automated understanding of unfamiliar terms and language
- We further explore this use-case as an illustrative example

The image shows a screenshot of a forum post from a user named 'kallysky'. The post title is 'Latest Zeus 2014 Botnet Tutorial Made Easy For Beginners'. The forum post includes a video player with a play button and a thumbnail showing a desktop interface with various icons. Below the video, there is text: 'Latest Zeus 2014 Botnet Tutorial Made Easy For Beginners (GRABS LATEST MOZILLA FIREFOX & IE BROWSER ETC.)'. There are two blue callout boxes with arrows pointing to the video and the text. The first callout box contains the text 'Video tutorial for configuring Zeus botnet'. The second callout box contains the text 'Message text: "Latest Zeus 2014 Botnet Tutorial Made Easy For Beginners"'. The forum post also includes a disclaimer: '(This Tutorial is for Educational Purpose Only. You agree to take sole responsibility for any damage caused by using this tutorial.)' and contact details: 'Contact Detail: FOR SETUP AND ASSISTANCE! | Yahoo Messenger: kallysky | Skype: kallysky1'.

Word2Vec Example



Benjamin, V., & Chen, H. (2015). Developing understanding of hacker language through the use of lexical semantics. In *2015 IEEE International Conference on Intelligence and Security Informatics (ISI)*, (pp. 79-84).

Word2Vec Example

- We identify and collect two hacker forums:

Forum	Members	Threads	Posts	Time Span
HackFive	947	1,108	5,334	1/24/2013 –12/30/2014
HackHound	633	507	3,622	12/10/2012 –12/30/2014

- **Text Data Extraction**

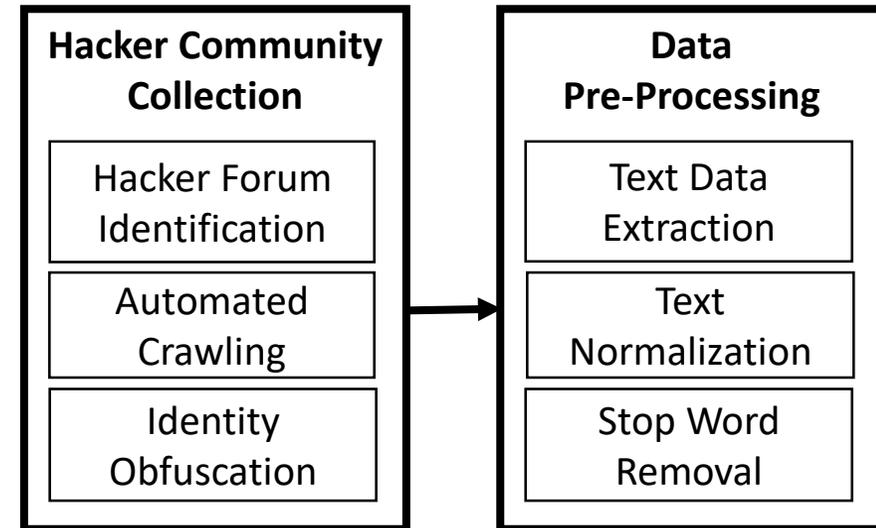
- Extract forum message data from raw webpages
- Thread titles, message bodies

- **Text Normalization**

- Strip punctuation from words
- Remove stop words

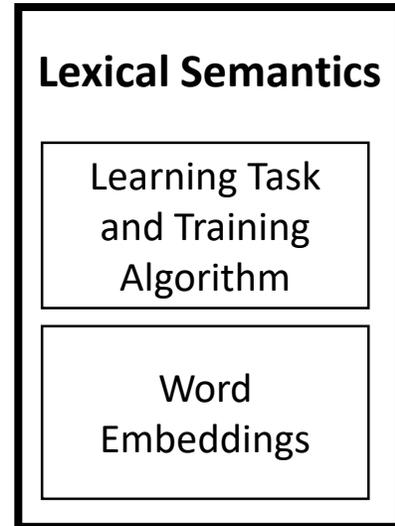
- **Scrubbed data is fed into Word2Vec one sentence at a time**

- Word2Vec's input is a single text file
- Each row of the input file represents one sentence



Word2Vec Example

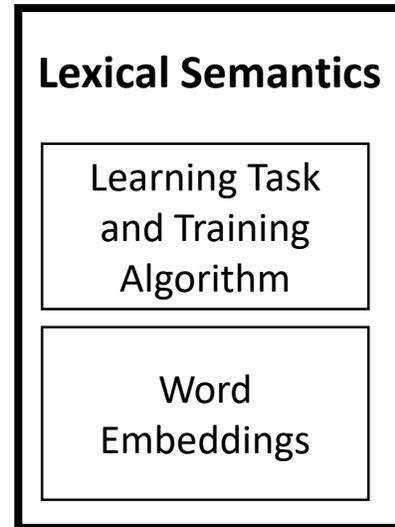
- Many parameters can be altered to change Word2Vec behavior
 - *Architecture*: Skip-gram (slower, better for infrequent words) vs CBOW (fast)
 - *Training algorithm*: Hierarchical Softmax (better for infrequent words) vs Negative Sampling (better for frequent words, better with low dimensional vectors)
 - *Dimensionality of the word vectors*: How many elements per word vector
 - *Context (window) size*: Size of sliding window used to iterate through words for training
- Skip-gram with negative sampling configuration generally provides best performance (Mikolov et al., 2013)
 - Additionally, Word2Vec comes default with recommended parameters
 - Default size for word vector dimensionality is 200 words
 - Default context window size of 10 words



Word2Vec Example



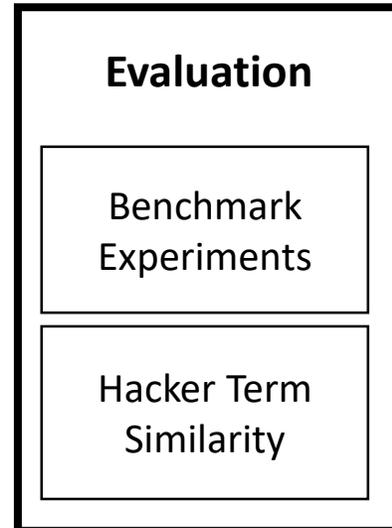
- The output of Word2Vec is a file called *Vectors.bin*
 - Can be opened and viewed in plaintext
 - Contains embeddings of each word in corpus
- Generated embeddings can be used in two ways:
 - Directly evaluated to better understand underlying corpus
 - Fed into other models and deep learning algorithms as features



Word2Vec Example

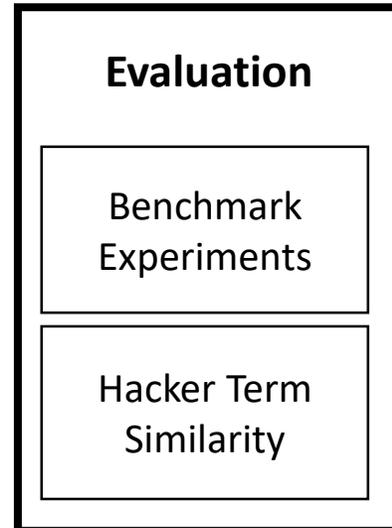
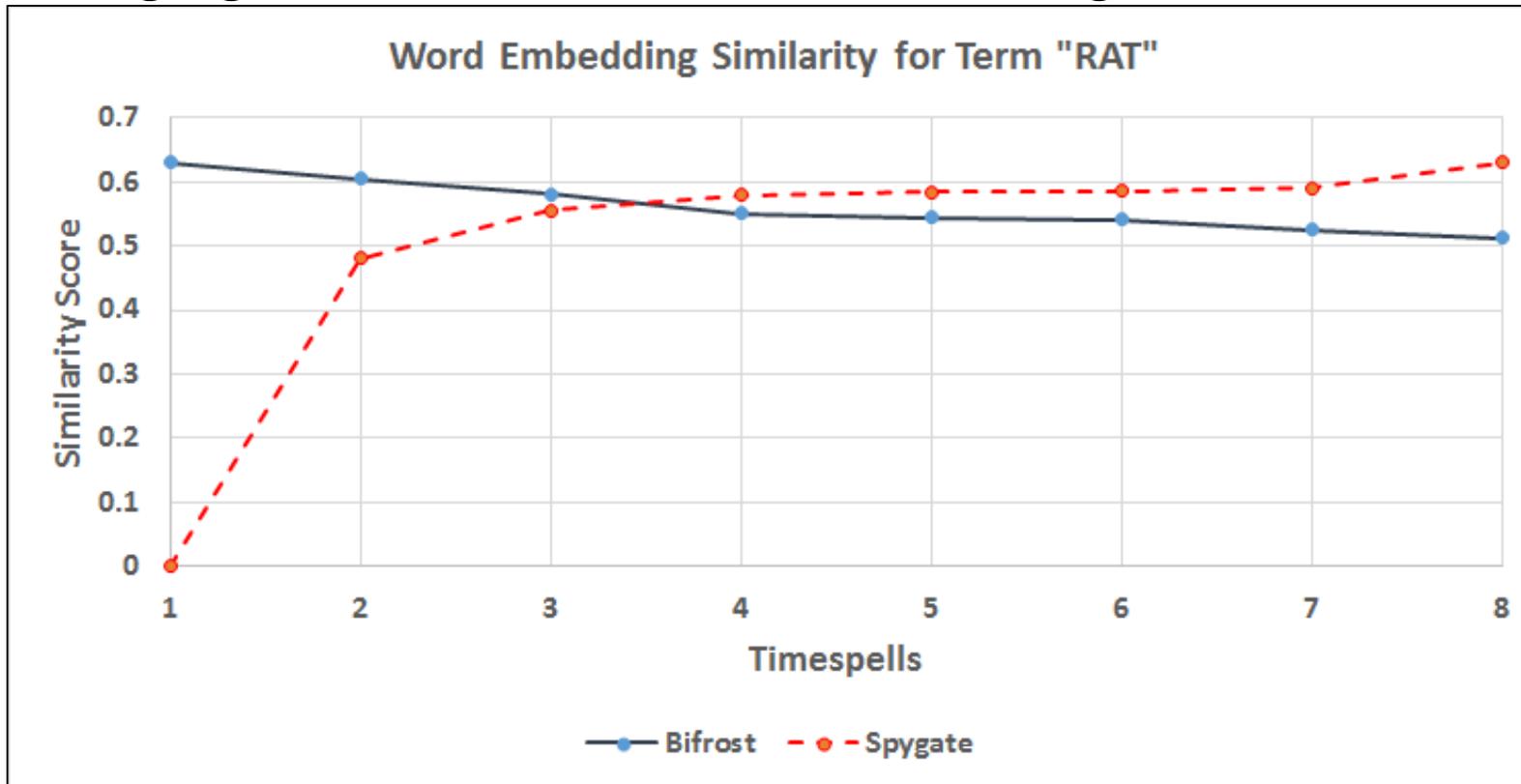
- We directly evaluate word embeddings in this study
 - Embeddings are vectors, can use cosine similarity to find similar words
 - Useful in hacker context to discover new hacker terms, tool names, etc.

	Most Similar Embeddings for "Botnet"	
	<i>Word</i>	<i>Similarity Score</i>
1	Citadel	0.561456
2	Zeus	0.554653
3	Partners	0.548900
4	Pandemiya	0.545221
5	Mailer	0.540075
6	Panel	0.524557
7	Linksys	0.498224
8	Cythosia	0.480465
9	Phase	0.464738
10	Spyeye	0.459695
<i>P@10</i>	70%	



Word2Vec Example

- *Bifrost* and *Spygate* are remote administration tools (RATs) that grant hackers backdoor access to victim computers
 - Can look at their similarity with word RAT over time to assess evolving significance in discussions concerning RATs

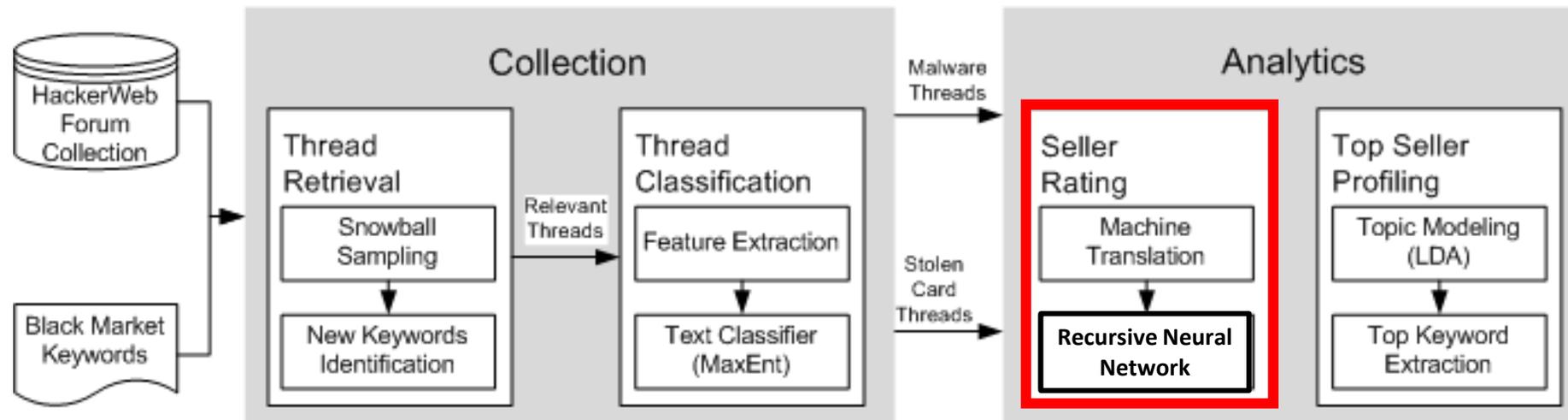


Recursive Neural Network Example – Identifying Key Carding Sellers

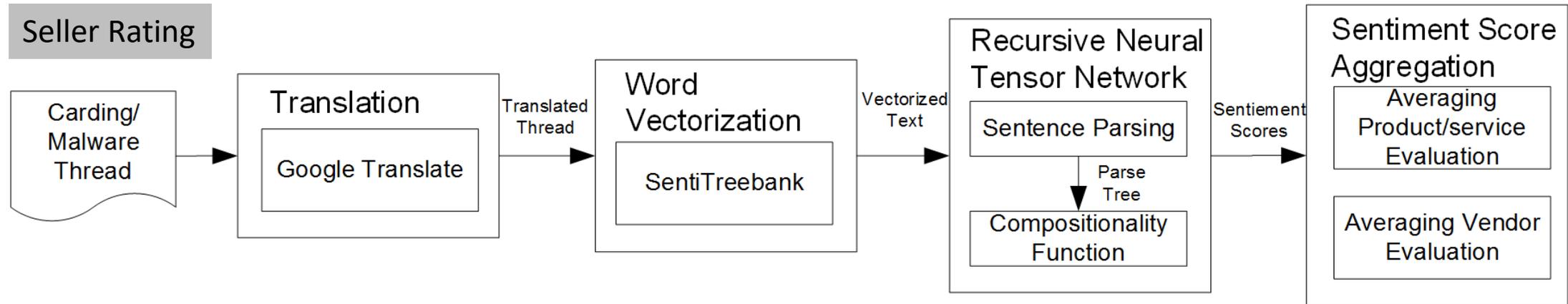
- The carding community rely heavily on the online black market for exchanging malwares and stolen data.
- Online black market for carders:
 - *Sellers* advertise their malwares and/or stolen data by posting a thread in the carding community
 - *Buyers* leave feedback commenting the their evaluation of the seller.
- **Objective:** to identify key carding sellers through evaluating buyers feedback.

Recursive Neural Network Example – Identifying Key Carding Sellers

- Input: Original threads from hacker forums
- Preprocessing:
 - Thread Retrieval: Identifying threads related to the underground economy by conducting snowball sampling-based keywords search
 - Thread Classification: Identifying advertisement threads using MaxEnt classifier
 - Focusing on malware advertisements and stolen card advertisement
 - Can be generalized to other advertisements.



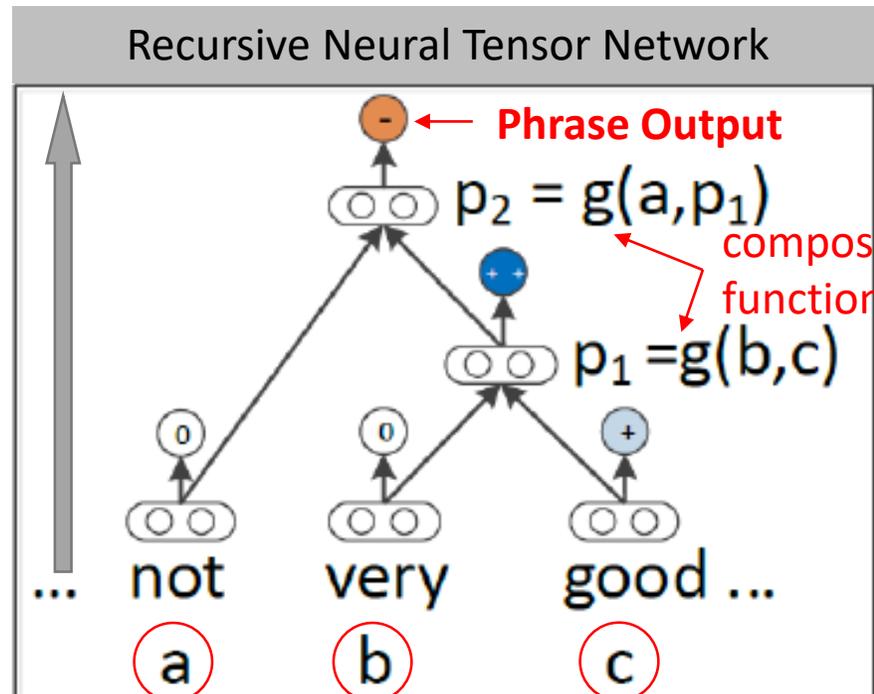
Recursive Neural Network Example – Identifying Key Carding Sellers



- *Translation* translates the thread content from the original language to English.
- *Word Vectorization* vectorizes each word in the text into a five dimensional vector using SentiTreeBank, a dictionary for customer feedback(Socher et al., 2013).
 - *SentiTreebank*: corpus with fully labeled parse trees with sentiment on each level
- *Recursive Neural Tensor Network* parses the sentence into a binary tree and aggregate semantics from constituents recursively.
- *Sentiment Score Aggregation* averages the sentiment measures for each advertisement and seller.

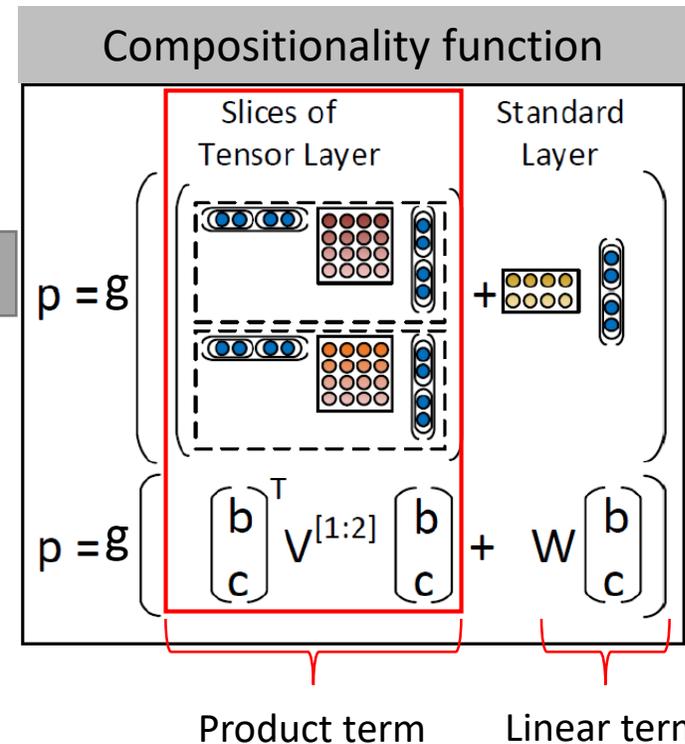
Recursive Neural Network Example – Identifying Key Carding Sellers

Example feedback: “*The dump is not very good*”



Word embeddings trained by SentiTreeBank

Output of the example feedback: **Negative**



The additional *tensor layer* makes the compositionality function more powerful by introducing the product of the constituent word embeddings.

Recursive Neural Network Example – Identifying Key Carding Sellers

Top 3 Best/Worst Malware and Stolen Data Sellers for Each Forum								
	Top 3 Best				Top 3 Worst			
	Malware		Stolen Data		Malware		Stolen Data	
Rank	User	Score	User	Score	User	Score	User	Score
Antichat								
1	L**G	5	i**o	3.6	N**g	1.8	l**s	2.3
2	V**U	4.5	a**s	3.5	k**a	2	p**A	2.3
3	g**l	4	D**R	3.4	D**i	2	s**8	2.4
CrdPro								
1	H**l	4	R**r	4.4	N**0	2	f**4	1.3
2	b**t	4	F**x	4	1**4	2	s**3	1.3
3	S**r	4	R**c	4	M**D	2	l**u	1.3
Zloy								
1	P**t	5	B**r	4	r**t	1.5	r**y	1
2	D**n	4	B**1	4	w**0	1.6	m**n	1
3	D**f	4	s**c	4	g**n	2	j**a	2

Outline

- Introduction
 - Motivation
 - Deep Learning Intuition
- Neural Network
 - Neuron, Feedforward algorithm, and Backpropagation algorithm
 - Limitations
- Deep Learning
 - Deep Belief Network: Autoencoder & Restricted Boltzman Machine
 - Convolutional Neural Network
- Deep Learning for Text Mining
 - Word Embedding
 - Recurrent Neural Network
 - Recursive Neural Network
- **Conclusions**

Conclusion

- Deep learning = Learning Hierarchical Representations
- Deep learning is thriving in big data analytics, including *image processing, speech recognition, and natural language processing*.
- Deep learning has matured and is very promising as an artificial intelligence method.
- Still has room for improvement:
 - Scaling computation
 - Optimization
 - Bypass intractable marginalization
 - More disentangled abstractions
 - Reasoning from incrementally added facts

Deep Learning Resources

Name	Language	Link	Note
Pylearn2	Python	http://deeplearning.net/software/pylearn2/	A machine learning library built on Theano
Theano	Python	http://deeplearning.net/software/theano/	A python deep learning library
Caffe	C++	http://caffe.berkeleyvision.org/	A deep learning framework by Berkeley
Torch	Lua	http://torch.ch/	An open source machine learning framework
Overfeat	Lua	http://cilvr.nyu.edu/doku.php?id=code:start	A convolutional network image processor
Deeplearning4j	Java	http://deeplearning4j.org/	A commercial grade deep learning library
Word2vec	C	https://code.google.com/p/word2vec/	Word embedding framework
GloVe	C	http://nlp.stanford.edu/projects/glove/	Word embedding framework
Doc2vec	C	https://radimrehurek.com/gensim/models/doc2vec.html	Language model for paragraphs and documents
StanfordNLP	Java	http://nlp.stanford.edu/	A deep learning-based NLP package
TensorFlow	Python	http://www.tensorflow.org	A deep learning based python library

References

- Bordes, A., Chopra, S., & Weston, J. (2014). Question answering with subgraph embeddings. *arXiv preprint arXiv:1406.3676*.
- Graves, A., Mohamed, A. R., & Hinton, G. (2013, May). Speech recognition with deep recurrent neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on* (pp. 6645-6649). IEEE.
- Graves, A. (2013). Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*.
- Irsoy, O., & Cardie, C. (2014, October). Opinion Mining with Deep Recurrent Neural Networks. In *EMNLP* (pp. 720-728).
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Rubenstein, H., & Goodenough, J. B. (1965). Contextual correlates of synonymy. *Communications of the ACM*, 8(10), 627-633.
- Socher, R., Perelygin, A., Wu, J. Y., Chuang, J., Manning, C. D., Ng, A. Y., & Potts, C. (2013, October). Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*(Vol. 1631, p. 1642).
- Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Advances in neural information processing systems* (pp. 3104-3112).
- Tai, K. S., Socher, R., & Manning, C. D. (2015). Improved semantic representations from tree-structured long short-term memory networks. *arXiv preprint arXiv:1503.00075*.

That is it for today

Thank You