# Outline

# Outline

# Multi-layer perceptron a.k.a. feedforward neural network

# Outline

Preliminary

11

# Back propagation



until convergence:
- do a forward pass
- compute the cost/error
- adjust weights ← how??

Adjust every weight $w_{i,j}$ by:

$$\Delta w_{i,j} = -\alpha \frac{\partial cost}{\partial w_{i,j}}$$

$\alpha$ is the learning rate.

# Back propagation



$$\Delta w_{i,j} = -\alpha \ \frac{\partial cost}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial w_{i,j}} \ \leftarrow \text{chain rule}$$

$$cost(\hat{y}, y) = \frac{1}{2}(y - \hat{y})^2$$

$$\hat{y}_j = x_{i,j} = \phi(o_{i,j})$$

# Back propagation



$$cost(\hat{y}, y) = \frac{1}{2}(y - \hat{y})^2$$

$$\hat{y}_j = x_{i,j} = \phi(o_{i,j}), \text{e.g. } \sigma(o_{i,j})$$

$$x_{i,j} = \sigma(o) = \frac{1}{1 + e^{-o}}$$

$$o_{i,j} = \sum_{k=1}^{K} w_{i,k} \cdot x_{i-1,k}$$

$$\Delta w_{i,j} = -\alpha \ \frac{\partial cost}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial o_{i,j}} \ \frac{\partial o_{i,j}}{\partial w_{i,j}} \quad \leftarrow \text{chain rule}$$

# Back propagation



$$cost(\hat{y}, y) = \frac{1}{2}(y - \hat{y})^2$$

$$\hat{y}_j = x_{i,j} = \phi(o_{i,j}), \text{e.g. } \sigma(o_{i,j})$$

$$x_{i,j} = \sigma(o) = \frac{1}{1 + e^{-o}}$$

$$o_{i,j} = \sum_{k=1}^{K} w_{i,k} \cdot x_{i-1,k}$$

$$\Delta w_{i,j} = -\alpha \ \frac{\partial cost}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial o_{i,j}} \ \frac{\partial o_{i,j}}{\partial w_{i,j}}$$

# Back propagation


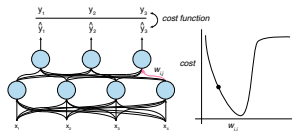
$$cost(\hat{y}, y) = \frac{1}{2}(y - \hat{y})^2$$

$$\hat{y}_j = x_{i,j} = \phi(o_{i,j}), \text{e.g. } \sigma(o_{i,j})$$

$$x_{i,j} = \sigma(o) = \frac{1}{1 + e^{-o}}$$

$$o_{i,j} = \sum_{k=1}^{K} w_{i,k} \cdot x_{i-1,k}$$

$$\Delta w_{i,j} = -\alpha \ \frac{\partial cost}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial o_{i,j}} \ \frac{\partial o_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial o_{i,j}} \ x_{i-1,j}$$

# Back propagation



$$cost(\hat{y}, y) = \frac{1}{2}(y - \hat{y})^2$$

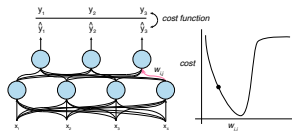$$\hat{y}_j = x_{i,j} = \phi(o_{i,j}), \text{e.g. } \sigma(o_{i,j})$$

$$x_{i,j} = \sigma(o) = \frac{1}{1 + e^{-o}}$$

$$\sigma'(o) = \sigma(o)(1 - \sigma(o))$$

$$o_{i,j} = \sum_{k=1}^{K} w_{i,k} \cdot x_{i-1,k}$$

$$\Delta w_{i,j} = -\alpha \ \frac{\partial cost}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial o_{i,j}} \ \frac{\partial o_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ x_{i,j}(1 - x_{i,j}) \ x_{i-1,j}$$

## Back propagation



$$cost(\hat{y}, y) = \frac{1}{2}(y - \hat{y})^2$$
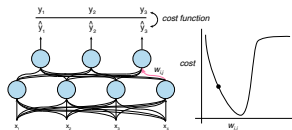
$$\hat{y}_j = x_{i,j} = \phi(o_{i,j}), \text{e.g. } \sigma(o_{i,j})$$

$$x_{i,j} = \sigma(o) = \frac{1}{1 + e^{-o}}$$

$$\sigma'(o) = \sigma(o)(1 - \sigma(o))$$

$$o_{i,j} = \sum_{k=1}^{K} w_{i,k} \cdot x_{i-1,k}$$

$$\Delta w_{i,j} = -\alpha \ \frac{\partial cost}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial o_{i,j}} \ \frac{\partial o_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \ y_j - x_{i,j} \ \ x_{i,j}(1 - x_{i,j}) \ \ x_{i-1,j}$$

18

# Back propagation



$$cost(\hat{y}, y) = \frac{1}{2}(y - \hat{y})^2$$

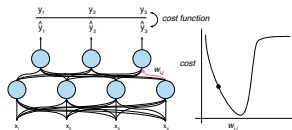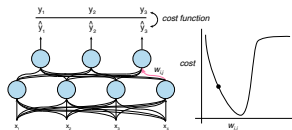$$\hat{y}_j = x_{i,j} = \phi(o_{i,j}), \text{e.g. } \sigma(o_{i,j})$$

$$x_{i,j} = \sigma(o) = \frac{1}{1 + e^{-o}}$$

$$\sigma'(o) = \sigma(o)(1 - \sigma(o))$$

$$o_{i,j} = \sum_{k=1}^{K} w_{i,k} \cdot x_{i-1,k}$$

$$\Delta w_{i,j} = -\alpha \ \frac{\partial cost}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ \frac{\partial cost}{\partial x_{i,j}} \ \frac{\partial x_{i,j}}{\partial o_{i,j}} \ \frac{\partial o_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \ y_j - x_{i,j} \ x_{i,j}(1 - x_{i,j}) \ x_{i-1,j}$$

$$= \text{l.rate} \ cost \ activation \ input$$

$\sigma(o) \quad \sigma'(o)$

## Back propagation

$$\Delta w_{i,j} = -\alpha \quad \frac{\partial cost}{\partial w_{i,j}}$$

$$= -\alpha \quad \frac{\partial cost}{\partial x_{i,j}} \quad \frac{\partial x_{i,j}}{\partial o_{i,j}} \quad \frac{\partial o_{i,j}}{\partial w_{i,j}}$$

$$= \mathsf{l.rate} \quad cost \quad activation \quad input$$

$$= -\alpha \quad \frac{\partial cost}{\partial x_{i,j}} \quad \frac{\partial x_{i,j}}{\partial o_{i,j}} \quad \frac{\partial o_{i,j}}{\partial w_{i,j}}$$

$$= -\alpha \quad \delta \quad x_{i-1,j}$$



$$\delta_{output} = (y_j - x_{i,j}) \quad x_{i,j}(1 - x_{i,j}) \qquad \leftarrow \mathsf{previous\ slide}$$

$$\delta_{hidden} = \left( \sum_{n \in nodes} \delta_n w_{n,j} \right) \quad x_{i,j}(1 - x_{i,j})$$

## Network representation



$\vec{y}$

$\vec{x}_3$

activation: $\quad \vec{x}_3 = \sigma(\vec{o}_2) = [\ 1 \times 3\ ]$

$\vec{x}^2 \quad * \quad W_2 \quad = \quad \vec{o}_2$
$[\ 1 \times 4\ ]\ [\ 4 \times 3\ ] = [\ 1 \times 3\ ]$

$\vec{x}_2$

activation: $\quad \vec{x}_2 = \sigma(\vec{o}_1) = [\ 1 \times 4\ ]$

$\vec{x}_1 \quad * \quad W_1 \quad = \quad \vec{o}_1$
$[\ 1 \times 4\ ]\ [\ 4 \times 4\ ] = [\ 1 \times 4\ ]$

$\vec{x}_1$ | $x_1[1]$ | $x_1[2]$ | $x_1[3]$ | $x_1[4]$ |

$y_1 \qquad y_2 \qquad y_3$

$\hat{y}_1 \qquad \hat{y}_2 \qquad \hat{y}_3$

$\leftrightarrow$

$x_1 \qquad x_2 \qquad x_3 \qquad x_4$

# Outline

Preliminary

22

# Distributed representations

- ▶ Represent units, e.g., words, as vectors
- ▶ Goal: words that are similar, e.g., in terms of meaning, should get similar embeddings

Cosine similarity to determine how similar two vectors are:

$$cosine(\vec{v}, \vec{w}) = \frac{\vec{v}^\top \cdot \vec{w}}{||\vec{v}||_2 ||\vec{w}||_2}$$

$$= \frac{\sum_{i=1}^{|v|} v_i * w_i}{\sqrt{\sum_{i=1}^{|v|} v_i^2} \sqrt{\sum_{i=1}^{|w|} w_i^2}}$$

$\overline{newspaper}$ = <0.08, 0.31, 0.41>

$\overline{magazine}$ = <0.09, 0.35, 0.36>

$\overline{biking}$ = <0.59, 0.25, 0.01>

# Distributed representations

How do we get these vectors?

- You shall know a word by the company it keeps [Firth, 1957]
- The vector of a word should be similar to the vectors of the words surrounding it

$$\overrightarrow{all} \quad \overrightarrow{you} \quad \overrightarrow{need} \quad \overrightarrow{is} \quad \overrightarrow{love}$$

# Embedding methods



*target distribution* — all / answer / amtrak / is / need / you / what / zorro

*target distribution* `0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0` ... `0 0 0 0 0 0 0`

*vocabulary size probabitity distribution*

*turn this into a probability distribution*

*vocabulary size layer*

*embedding size × vocabulary size weight matrix*

*embedding size hidden layer*

*vocabulary size × embedding size weight matrix*

*vocabulary size inputs* `1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0` ... `0 1 0 0 0 0`

all / answer / amtrak / is / love / need / you / what / zorro

25

## Probability distributions

$\mathrm{softmax}$ = normalize the logits

$$= \frac{e^{logits[i]}}{\sum_{j=1}^{|logits|} e^{logits[j]}}$$

$cost$ = cross entropy loss

$$= -\sum_{x} p(x) \log \hat{p}(x)$$
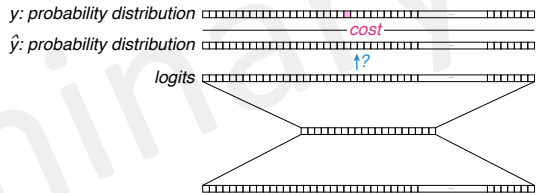
$$= -\sum_{i} p_{\mathsf{ground\ truth}}(word = vocabulary[i]) \log p_{\mathsf{predictions}}(word = vocabulary[i])$$

$$= -\sum_{i} y_i \log \hat{y}_i$$



y: probability distribution

ŷ: probability distribution

logits

cost

↕?

26

# Outline

Preliminary

# Recurrent neural networks

- Lots of information is sequential and requires a memory for successful processing
- Sequences as input, sequences as output



one to one    one to many    many to one    many to many    many to many

- Recurrent neural networks (RNNs) are called recurrent because they perform same task for every element of sequence, with output dependent on previous computations
- RNNs have memory that captures information about what has been computed so far
- RNNs can make use of information in arbitrarily long sequences – in practice they limited to looking back only few steps

Image credits: http://karpathy.github.io/assets/rnn/diags.jpeg

28

# Recurrent neural networks

▶ RNN being unrolled (or unfolded) into full network

▶ Unrolling: write out network for complete sequence



▶ Formulas governing computation:
  ▶ $x_t$ input at time step $t$
  ▶ $s_t$ hidden state at time step $t$ – memory of the network, calculated based on previous hidden state and input at the current step: $s_t = f(Ux_t + Ws_{t-1})$; $f$ usually nonlinearity, e.g., $\tanh$ or $\mathrm{ReLU}$; $s_{-1}$ typically initialized to all zeroes
  ▶ $o_t$ output at step $t$. E.g.,, if we want to predict next word in sentence, a vector of probabilities across vocabulary: $o_t = \mathrm{softmax}(Vs_t)$

Image credits: Nature

# Language modeling using RNNs

- **Language model** allows us to predict probability of observing sentence (in a given dataset) as: $P(w_1, \ldots, w_m) = \prod_{i=1}^{m} P(w_i \mid w_1, \ldots, w_{i-1})$

- In RNN, set $o_t = x_{t+1}$: we want output at step $t$ to be actual next word

- Input $x$ a sequence of words; each $x_t$ is a single word; we represent each word as a one-hot vector of size `vocabulary_size`

- Initialize parameters $U$, $V$, $W$ to small random values around $0$

# Language modeling using RNNs

- **Language model** allows us to predict probability of observing sentence (in a given dataset) as: $P(w_1, \ldots, w_m) = \prod_{i=1}^{m} P(w_i \mid w_1, \ldots, w_{i-1})$

- In RNN, set $o_t = x_{t+1}$: we want output at step $t$ to be actual next word

- Input $x$ a sequence of words; each $x_t$ is a single word; we represent each word as a one-hot vector of size `vocabulary_size`

- Initialize parameters $U$, $V$, $W$ to small random values around $0$

- Cross-entropy loss as loss function

- For $N$ training examples (words in text) and $C$ classes (the size of our vocabulary), loss with respect to predictions $o$ and true labels $y$ is: $\mathcal{L}(y, o) = -\frac{1}{N} \sum_{n \in N} y_n \log o_n$

- Training RNN similar to training a traditional NN: backpropagation algorithm, but with small twist

- Parameters shared by all time steps, so gradient at each output depends on calculations of previous time steps: Backpropagation Through Time

# Vanishing and exploding gradients

- For training RNNs, calculate gradients for $U$, $V$, $W$ – ok for $V$ but for $W$ and $U$ ...

- Gradients for $W$:

$$\frac{\partial \mathcal{L}_3}{\partial W} = \frac{\partial \mathcal{L}_3}{\partial o_3} \frac{\partial o_3}{\partial s_3} \frac{\partial s_3}{\partial W} \quad = \sum_{k=0}^{3} \frac{\partial \mathcal{L}_3}{\partial o_3} \frac{\partial o_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial W}$$



- More generally: $\frac{\partial \mathcal{L}}{\partial s_t} = \frac{\partial \mathcal{L}}{\partial s_m} \cdot \underbrace{\frac{\partial s_m}{\partial s_{m-1}}}_{< 1} \cdot \underbrace{\frac{\partial s_{m-1}}{\partial s_{m-2}}}_{< 1} \cdot \ldots \cdot \underbrace{\frac{\partial s_{t+1}}{\partial s_t}}_{< 1} \Rightarrow \; \ll 1$

- Gradient contributions from far away steps become zero: state at those steps doesn't contribute to what you are learning

Image credits: http://www.wildml.com/2015/10/

recurrent-neural-networks-tutorial-part-3-backpropagation-through-time-and-vanishing-grad
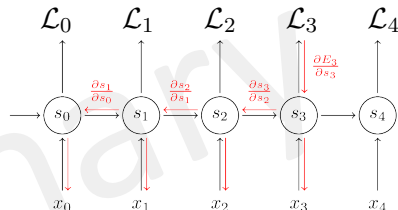
# Long Short Term Memory [Hochreiter and Schmidhuber, 1997]

LSTMs designed to combat vanishing gradients through gating mechanism

▶ How LSTM calculates hidden state $s_t$

$$i = \sigma(x_t U^i + s_{t-1} W^i)$$
$$f = \sigma(x_t U^f + s_{t-1} W^f)$$
$$o = \sigma(x_t U^o + s_{t-1} W^o)$$
$$g = \tanh(x_t U^g + s_{t-1} W^g)$$
$$c_t = c_{t-1} \circ f + g \circ i$$
$$s_t = \tanh(c_t) \circ o$$

($\circ$ is elementwise multiplication)

▶ RNN computes hidden state as
$s_t = \tanh(U x_t + W s_{t-1})$ – an LSTM
unit does exact same thing

# Long Short Term Memory [Hochreiter and Schmidhuber, 1997]

LSTMs designed to combat vanishing gradients through gating mechanism

- How LSTM calculates hidden state $s_t$

$$i = \sigma(x_t U^i + s_{t-1} W^i)$$
$$f = \sigma(x_t U^f + s_{t-1} W^f)$$
$$o = \sigma(x_t U^o + s_{t-1} W^o)$$
$$g = \tanh(x_t U^g + s_{t-1} W^g)$$
$$c_t = c_{t-1} \circ f + g \circ i$$
$$s_t = \tanh(c_t) \circ o$$

  ($\circ$ is elementwise multiplication)

- RNN computes hidden state as
  $s_t = \tanh(U x_t + W s_{t-1})$ – an LSTM
  unit does exact same thing

- $i$, $f$, $o$: input, forget and output gates
- Gates optionally let information through: composed out of sigmoid neural net layer and pointwise multiplication operation
- $g$ is a candidate hidden state computed based on current input and previous hidden state
- $c_t$ is internal memory of LSTM unit: combines previous memory $c_{t-1}$ multiplied by forget gate, and newly computed hidden state $g$, multiplied by input gate

32

# Long Short Term Memory [Hochreiter and Schmidhuber, 1997]

LSTMs designed to combat vanishing gradients through gating mechanism

- How LSTM calculates hidden state $s_t$

$$i = \sigma(x_t U^i + s_{t-1} W^i)$$
$$f = \sigma(x_t U^f + s_{t-1} W^f)$$
$$o = \sigma(x_t U^o + s_{t-1} W^o)$$
$$g = \tanh(x_t U^g + s_{t-1} W^g)$$
$$c_t = c_{t-1} \circ f + g \circ i$$
$$s_t = \tanh(c_t) \circ o$$

  ($\circ$ is elementwise multiplication)

- RNN computes hidden state as $s_t = \tanh(U x_t + W s_{t-1})$ – an LSTM unit does exact same thing

- $\cdots$
- Compute output hidden state $s_t$ by multiplying memory with output gate
- Plain RNNs a special case of LSTMs:
  - Fix input gate to all 1's
  - Fix forget gate to all 0's (always forget the previous memory)
  - Fix output gate to all 1's (expose the whole memory)
  - Additional $\tanh$ squashes output
- Gating mechanism allows LSTMs to model long-term dependencies
- Learn parameters for gates, to learn how memory should behave

33

## Gated Recurrent Units

▶ GRU layer quite similar to that of LSTM layer, as are the equations:

$$z = \sigma(x_t U^z + s_{t-1} W^z)$$
$$r = \sigma(x_t U^r + s_{t-1} W^r)$$
$$h = \tanh(x_t U^h + (s_{t-1} \circ r) W^h)$$
$$s_t = (1 - z) \circ h + z \circ s_{t-1}$$

▶ GRU has two gates: reset gate $r$ and update gate $z$.
  ▶ Reset gate determines how to combine new input with previous memory; update gate defines how much of the previous memory to keep around
  ▶ Set reset to all 1's and update gate to all 0's to get plain RNN model
▶ On many tasks, LSTMs and GRUs perform similarly

# Bidirectional RNNs

- Bidirectional RNNs based on idea that output at time $t$ may depend on previous and future elements in sequence
    - Example: predict missing word in a sequence
- Bidirectional RNNs are two RNNs stacked on top of each other
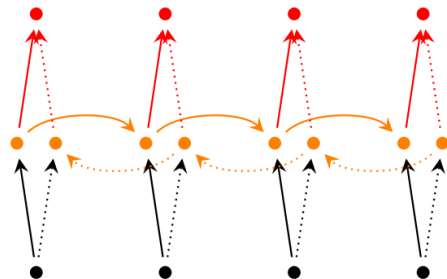- Output is computed based on hidden state of both RNNs



Image credits: http://www.wildml.com/2015/09/
recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/

35

## Attention

- Attention mechanisms come from visual analysis — suppose you want to locate a specific object in a large image; like people, focus on specific areas
- First applied to text and NLP in 2015
- Basic mechanism behind every attention mechanism
    1. Read operator: read a "patch" from the input
    2. Glimpse sensor: extract information from "patch"
    3. Locator: predict the next location of read operator
    4. RNN: combine the previous and current responses from glimpse sensor

Assume we are at time step $t$. From $t - 1$ we get next location to which we should pay attention (produced by locator). Move sensor there and extract information, which the RNN combines with previous outputs. After several iterations, we produce final response, e.g., classification or label.

# Attention

### Hard attention

- Read operator: fixed size, but there may be several of them
- Glimpse sensor can be any NN
- Locator predicts $x$ and $y$ location of sensor (images) or words ahead/previous (text)

Not differentiable, so Reinforcement Learning is used

### Soft attention

- Read operator: it only has fixed aspect ratio; it can zoom into the image and blur if needed
- Glimpse sensor can be any NN
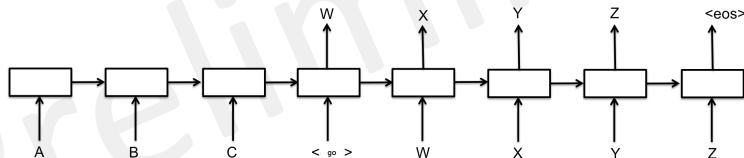- Locator predict more than $x$ and $y$ parameter (like sigma for blur .etc)

Differentiable

# Outline

## Sequence-to-sequence models

Increasingly important: not just retrieval but also generation

- Snippets, summaries, small screen versions of search results, spoken results, chatbots, conversational interfaces, . . . , but also query suggestion, query correction, . . .

Basic sequence-to-sequence (seq2seq) model consists of two RNNs: an encoder that processes input and a decoder that generates output:



Each box represents cell of RNN (often GRU cell or LSTM cell). Encoder and decoder can share weights or, as is more common, use a different set of parameters

## Sequence-to-sequence models

- seq2seq models build on top of language models
  - Encoder step: a model converts input sequence into a fixed representation
  - Decoder step: a language model is trained on both the output sequence (e.g., translated sentence) as well as fixed representation from the encoder
  - Since decoder model sees encoded representation of input sequence as well as output sequence, it can make more intelligent predictions about future words based on current word
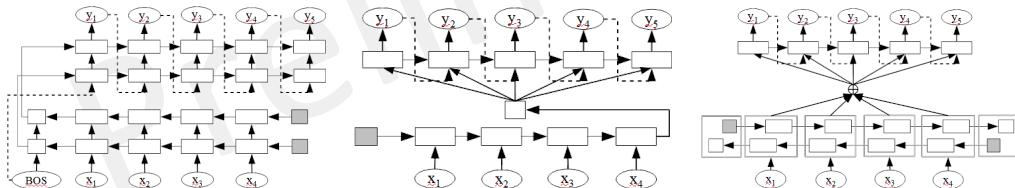


Image credits: [Sutskever et al., 2014; Cho et al., 2014; Bahdanou et al., 2014]

# Sequence-to-sequence models

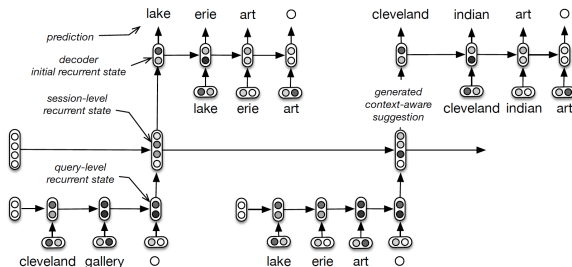Used for a "traditional information retrieval task"



Figure 3: The hierarchical recurrent encoder-decoder (HRED) for query suggestion. Each arrow is a non-linear transformation. The user types *cleveland gallery → lake erie art*. During training, the model encodes *cleveland gallery*, updates the session-level recurrent state and maximize the probability of seeing the following query *lake erie art*. The process is repeated for all queries in the session. During testing, a contextual suggestion is generated by encoding the previous queries, by updating the session-level recurrent states accordingly and by sampling a new query from the last obtained session-level recurrent state. In the example, the generated contextual suggestion is *cleveland indian art*.

Image credits: Sordoni et al. [2015a]
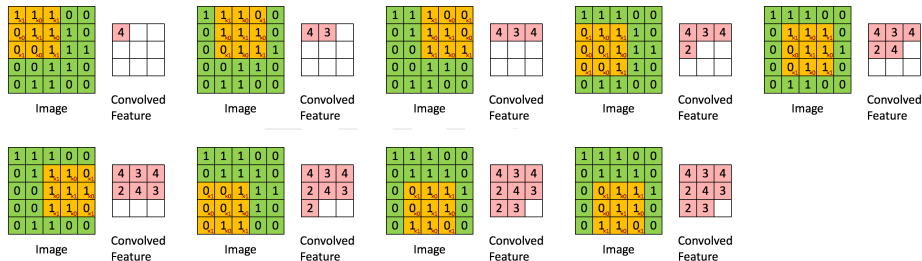
# Outline

Preliminary

# Convolutional neural networks

Major breakthroughs in image classification – at core of many computer visions systems

Some initial applications of CNNs to problems in text and information retrieval

What is a convolution? Intuition: sliding window function applied to a matrix

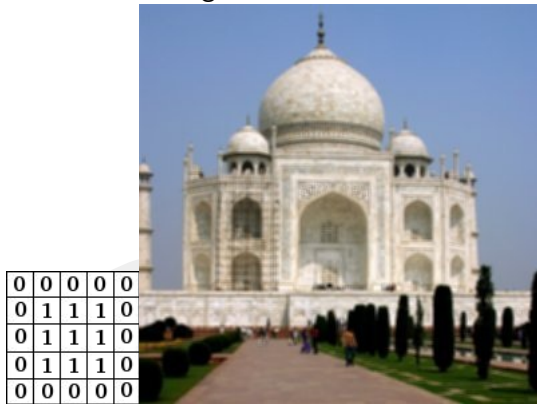Example: convolution with $3 \times 3$ filter



Multiply values element-wise with original matrix, then sum. Slide over whole matrix.

Image credits:

http://deeplearning.stanford.edu/wiki/index.php/Feature_extraction_using_convolution

# Visual examples of CNNs

Averaging each pixel with neighboring values blurs image:

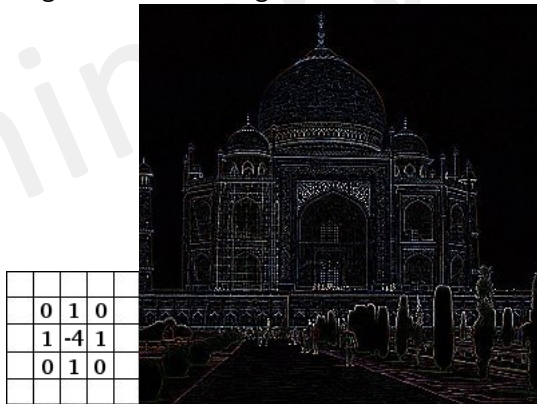Taking difference between pixel and its neighbors detects edges:



| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 |

|   | 0 | 1 | 0 |   |
|---|---|---|---|---|
|   | 1 | -4 | 1 |   |
|   | 0 | 1 | 0 |   |
|   |   |   |   |   |

Image credits: https://docs.gimp.org/en/plug-in-convmatrix.html
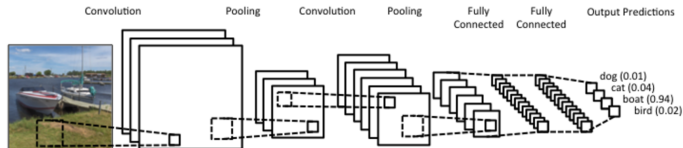
44

# Convolutional neural networks



- Use convolutions over input layer to compute output

- Yields local connections: each region of input connected to a neuron in output

- Each layer applies different filters and combines results

- Pooling (subsampling) layers

- During training, CNN learns values of filters

- Image classification a CNN may learn to detect edges from raw pixels in first layer

- Then use edges to detect simple shapes in second layer

- Then use shapes to detect higher-level features, such as facial shapes in higher layers

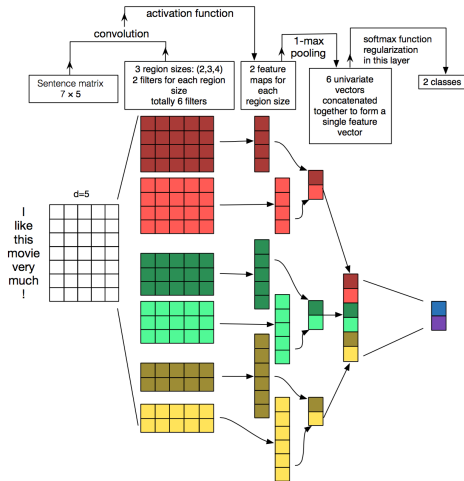- Last layer is then a classifier that uses high-level features

# CNNs in text

## Basic intution

- Instead of image pixels, input to most NLP tasks are sentences or documents represented as a matrix. Each row of matrix corresponds to one token, typically a word, but could be a character. That is, each row is vector that represents word.

- Typically, these vectors are word embeddings (low-dimensional representations) like word2vec or GloVe, but they could also be one-hot vectors that index the word into a vocabulary.

- For a 10 word sentence using a 100-dimensional embedding we would have a $10 \times 100$ matrix as our input.

- That's our "image"

- Typically use filters that slide over full rows of the matrix (words): the "width" of our filters is usually the same as the width of the input matrix. The height, or region size, may vary, but sliding windows over 2-5 words at a time is typical.

# CNNs in text

Example architecture (Zhang and Wallace, 2015; Sentence classification)

# CNNs in text

### Example uses in IR

- ► MSR: how to learn semantically meaningful representations of sentences that can be used for Information Retrieval
- ► Recommending potentially interesting documents to users based on what they are currently reading
- ► Sentence representations are trained based on search engine log data
- ► Gao et al. Modeling Interestingness with Deep Neural Networks. EMNLP 2014; Shen et al. A Latent Semantic Model with Convolutional-Pooling Structure for Information Retrieval. CIKM 2014.