

all other parameters. Word embeddings could be considered the first step towards deep learning (neural network) based NLP.

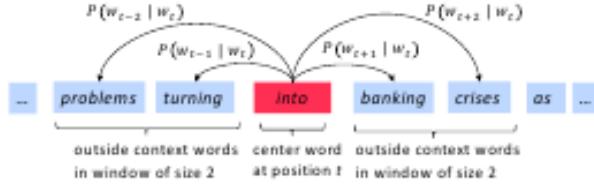
### 3.2.1 Word2Vec Overview

Word2Vec (Mikolov et al. 2013) is a framework for learning word vector models.

Suppose, we have a large corpus of text. We can construct a vocabulary from this corpus, including in the vocabulary e.g. all the unique words in it. We fix the vocabulary after this, i.e. we will not be adding any words there.

We represent each word in the vocabulary by a vector. Initially, this vector could be purely random, therefore we need somehow adjust it. We go through each position  $t$  in the corpus text; in each position  $t$  there is a center word  $c$  and a set of context (“outside”) words  $o$ . We use a similarity function to compare word vectors for  $c$  and  $o$  and to calculate the probability of  $o$  given  $c$  (or vice versa). Keep adjusting the word vectors to maximize this probability.

Let us have a look on more formal way of this process description.  $P(w_{t+j}|w_t)$  is a probability for a word  $w_{t+j}$  from set  $o$  to be present in context of central word  $w_t$ .



### 3.2.2 Word2Vec: objective function

First of all about the objective function itself. It could be known by several names, like objective, cost, or even loss function. But now we are discussing the word2vec loss function in details.

For each position  $t = 1, \dots, T$ , predict context words within a window of fixed size  $m$ , given center word  $w_t$ .

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j}|w_t; \theta),$$

where  $\theta$  refers to the whole set of optimizable model parameters.

The objective function  $J(\theta)$  is the (average) negative log likelihood.

Negative Log Likelihood Loss:

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j}|w_t; \theta)$$

is also called Cross-Entropy Loss.

Cross-entropy is a measure of the difference between two probability distributions for a given random variable or set of events.

Consider two probability distributions  $p$  and  $q$  defined on a set of events  $X = \{x_1, x_2, \dots, x_n\}$ , then cross-entropy between  $p$  and  $q$  is:

$$H(q, p) = - \sum_{x \in X} q(x) \log p(x)$$

Assume  $X$  is the vocabulary,  $p(x)$  is the model generated probabilities over the vocabulary,  $q(x)$  is the actual distribution of the content word at  $t+j$ :

$$q(x) = \begin{cases} 1, & \text{for } x = w_{t+j} \\ 0, & \text{otherwise} \end{cases}$$

then:

$$H(q, p) = -\log p(w_{t+j})$$

### 3.2.3 Word2vec prediction function

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

1. Dot product compares similarity of  $o$  and  $c$ . Larger dot product = larger probability
2. Exponentiation makes anything positive
3. Normalize over entire vocabulary gives "probability" distribution

This is an example of the **softmax function**  $\mathbb{R}^n \Rightarrow (0, 1)^n$ :

$$\text{softmax}(x_i) = \frac{\exp x_i}{\sum_{j=1}^n \exp x_j}$$

The softmax function maps arbitrary values  $x_i$  to a "probability" distribution  $p_i$ .

- "max" because amplifies "probability" of largest  $x_i$
- "soft" because still assigns some "probability" to smaller  $x_i$
- frequently used in Deep Learning

In mathematics, the *softmax* function, also known as *softmax* or *normalized exponential function*, is a function that takes as input a vector of K real numbers, and normalizes it into a probability distribution consisting of K probabilities proportional to the exponentials of

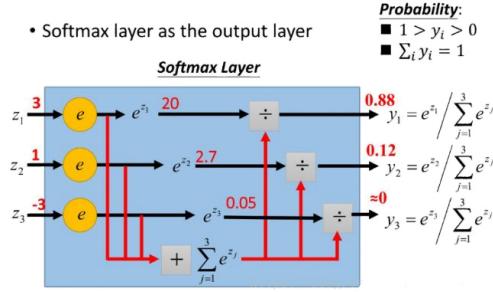


Figure 3.4: Softmax

the input numbers. The standard (unit) softmax function  $\sigma : \mathbb{R}^K \rightarrow \mathbb{R}^K$  is defined by the formula:

$$y_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

In Word2Vec, because the softmax function is calculated over all words in the vocabulary, it is quite expensive computationally.

### 3.2.4 Training a model by optimizing parameters

To train a model, we adjust parameters to minimize a loss. We compute all vector gradients.  $\theta$  represents all model parameters in one long vector. We optimize these parameters by walking down the gradient. In our case with  $d$ -dimensional vectors and  $V$ -many words (every word has two vectors):

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \dots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \dots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

We went through gradient for each center vector  $v$  in a window. We also need gradients for outside vectors  $u$ . Generally in each window we will compute updates for all parameters that are being used in that window.

### 3.2.5 Optimization: gradient descent

We have a cost function  $J(\theta)$  we want to minimize. Gradient descent is an algorithm to minimize  $J(\theta)$ : for current value of  $\theta$  calculate gradient of  $J(\theta)$ , then take small step in direction of negative gradient, repeat.

Update equation (in matrix notation):

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$

For single parameter:

$$\theta^{new} = \theta^{old} - \alpha \frac{\partial}{\partial \theta_j^{old}} J(\theta)$$

$\alpha$  - step size or *learning rate*

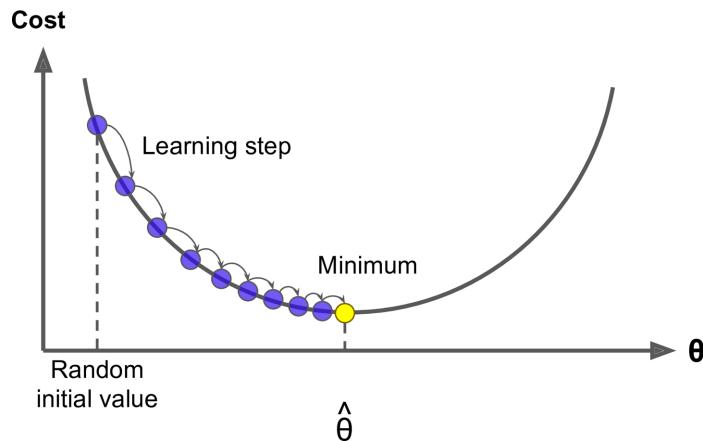


Figure 3.5: Gradient descent

Evaluation of gradient:

$$\begin{aligned}
 J(\theta) &= -\frac{1}{T} \log L(\theta) \\
 &= -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log p(w_{t+j} | w_t; \theta) \\
 &= -\frac{1}{T} \sum_{o \in context(c)} \sum_{c \in corpus} \log p(o | c; u, v) \\
 \frac{\partial}{\partial \theta} J(\theta) &= -\frac{1}{T} \sum_{o \in context(c)} \sum_{c \in corpus} \frac{\partial}{\partial \theta} \log p(o | c; u, v)
 \end{aligned}$$

$$\begin{aligned}
\frac{\partial}{\partial v_c} \log p(o|c; u, v) &= \frac{\partial}{\partial v_c} \log \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} \\
&= \frac{\partial}{\partial v_c} (u_o^T v_c) - \frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^T v_c) \\
&= u_o - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \frac{\partial}{\partial v_c} \sum_{w=1}^V \exp(u_w^T v_c) \\
&= u_o - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \sum_{w=1}^V \frac{\partial}{\partial v_c} \exp(u_w^T v_c) \\
&= u_o - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \sum_{w=1}^V \exp(u_w^T v_c) u_w \\
&= u_o - \sum_{x=1}^V \frac{\exp(u_x^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} u_x \\
&= u_o - \sum_{x=1}^V p(x|c) u_x
\end{aligned}$$

$$\frac{\partial}{\partial \theta} J(\theta) = -\frac{1}{T} \sum_{o \in \text{context}(c)} \sum_{c \in \text{corpus}} \frac{\partial}{\partial \theta} \log p(o|c; u, v)$$

$$\begin{aligned}
\frac{\partial}{\partial v_c} J(\theta) &= -\frac{1}{T} \sum_{o \in \text{context}(c)} \sum_{c \in \text{corpus}} \left[ u_o - \sum_{x=1}^V p(x|c) u_x \right] \\
\frac{\partial}{\partial u_o} J(\theta) &= -\frac{1}{T} \sum_{o \in \text{context}(c)} \sum_{c \in \text{corpus}} \left[ v_c - \sum_{x=1}^V p(o|x) v_x \right]
\end{aligned}$$

### 3.2.6 Stochastic gradient descent

$J(\theta)$  is a function of all windows in the corpus (potentially billions!), so  $\nabla_\theta J(\theta)$  is very expensive to compute. Instead you can repeatedly sample windows from corpus, and update after each one. It is **stochastic gradient descent** (SGD).

We iteratively take gradients at each window for SGD. But in each window, we only have at most  $2m + 1$  words, so  $\nabla_\theta J(\theta)$  is very sparse. The solution is sparse matrix operations to only update certain rows of full embedding matrices  $U$  and  $V$  or keeping a hash for word vectors.

### 3.2.7 Word2vec: details

We have 2 vectors  $u$  and  $v$  for easier optimization. So there are two model variants:

1. **Skip-grams** (SG). Predict context words (position independent) given center word.
2. **Continuous Bag of Words** (CBOW). Predict center word from (bag of) context words.

Method for additional efficiency in training: **negative sampling**.

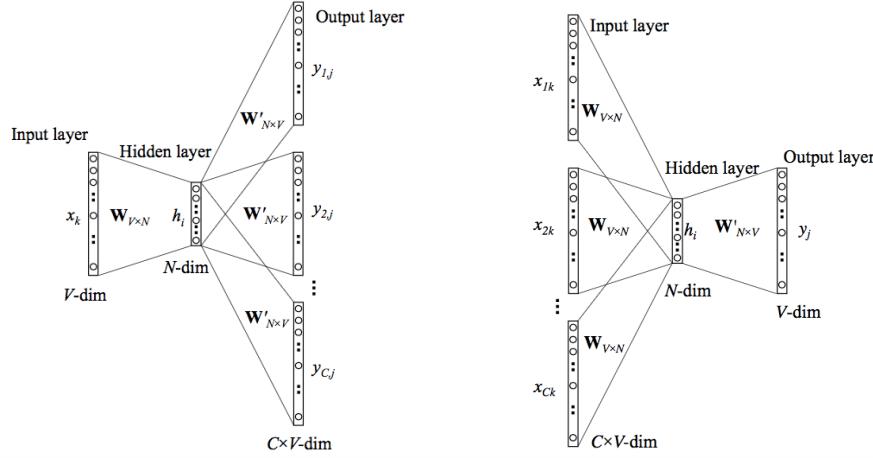


Figure 3.6: Skip-gram Model vs. CBOW Model

### 3.2.8 Character n-gram based model

It is difficult for good representations of *rare words* to be learned with traditional word2vec. There could be words in the NLP task that were not present in the word2vec training corpus. This limitation is more pronounced in case of morphologically rich languages: in French or Spanish, most verbs have more than forty different inflected forms, while the Finnish languages has fifteen different cases for nouns. It is possible to improve vector representations for morphologically rich languages by using *character level* information.

The basic skip-gram model described above ignores the internal structure of the word. However, character n-gram based model incorporates information about the structure in terms of character n-gram embeddings. It is supposed that each word  $w$  is represented as a bag of character n-gram.

Word *where* and  $n = 3$  will be represented by character n-grams: *wh/whe/her/ere/re*; and the special sequence *wherei*. NB: Word *heri* is different from the tri-gram *her* from the word *where*.

**Computing word vector representation:** We represent a word by the sum of the vector representations of its n-gram. We thus obtain the scoring function:

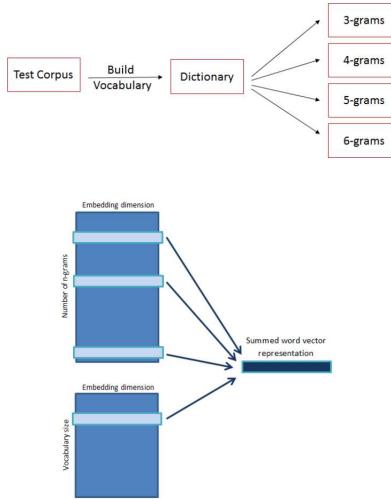
$$s(w, c) = \sum_{g \in \zeta_w} z_g^T v_c$$

where  $w$  - given word,  $\zeta_w$  - the set on n-grams appearing it word  $w$ ,  $z_g$  - vector representation to each n-grams,  $v_c$  - the word vector of the center word  $c$ .

We extract all the n-grams ( $3 \leq n \leq 6$ ):

and compute word vector representation

This model is more robust to the size of the training data, it seems to quickly saturate and adding more and more data does not always lead to improved result. However, the performance of the CBOW model gets better as more and more data is available. On the



other hand, the performance subword information skip-gram with very small dataset better than the performance CBOW.

### 3.2.9 Count based vs. direct prediction

Count based	Direct
+ Fast training	+ Generate improved performance on other tasks
+ Efficient usage of statistics	+ Can capture complex patterns beyond word similarity
- Primarily used to capture word similarity	- Scales with corpus data
- Disproportional importance given to large counts	- Inefficient usage of statistics

## 3.3 GloVe

Ratios of co-occurrence probabilities can encode meaning components:

	$x = \text{solid}$	$x = \text{gas}$	$x = \text{water}$	$x = \text{random}$
$P(x \text{ice})$	large	small	large	small
$P(x \text{steam})$	small	large	large	small
$\frac{P(x \text{ice})}{P(x \text{steam})}$	large	small	$\sim 1$	$\sim 1$

We can capture ratios of co-occurrence probabilities as linear meaning components in a word vector space.

Log-bilinear model:  $w_i \cdot \tilde{w}_j = \log P(i|j)$  with vector differences:  $w_x \cdot (\tilde{w}_a - \tilde{w}_b) = \log \frac{P(x|a)}{P(x|b)}$

**Note:**  $P(i|j) \neq P(j|i)$ ,  $w$  and  $\tilde{w}$  should be defined separately!

Glove cost function is:

$$J = \sum_{i=1}^V f(X_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

- Fast training

- Scalable to huge corpora
- Good performance even with small corpus and small vectors

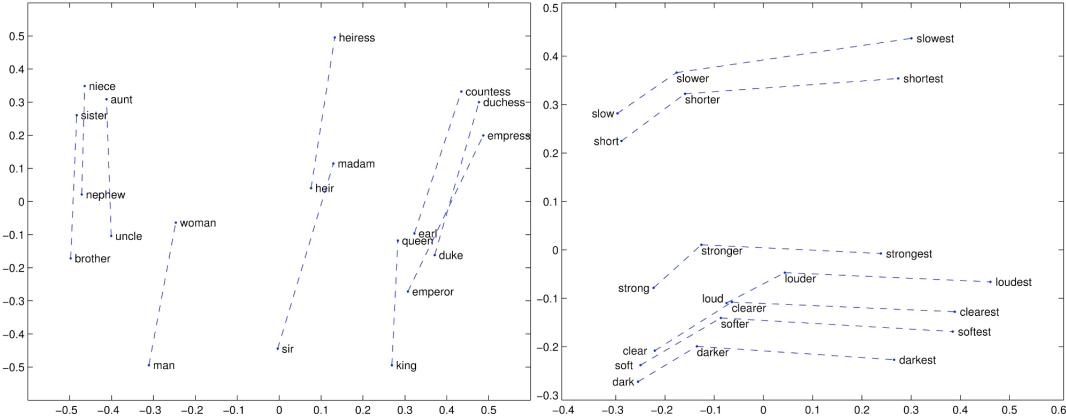


Figure 3.7: Word vector analogies (Glove)

## 3.4 How to evaluate word vectors?

Two approaches related to general evaluation in NLP: intrinsic and extrinsic.

Intrinsic:

- Evaluation on a specific/intermediate subtask
- Fast to compute
- Helps to understand that system
- Not clear if really helpful unless correlation to real task is established

Extrinsic:

- Evaluation on a real task
- Can take a long lime to compute accuracy
- Unclear if the subsystem is the problem or its interaction of other subsystems
- If replacing exactly one subsystem with another improves accuracy → winning!

### 3.4.1 Intrinsic word vector evaluation

Word Vector Analogies:  $word_a : word_b \text{ as } word_c : ?$ . Evaluate word vectors by how well their cosine distance after addition captures intuitive semantic and syntactic analogy question. Discarding the input words from the search.

$$d = \arg \max_i \frac{(x_b - x_a + x_c)^T x_i}{\|x_b - x_a + x_c\|}$$

**Problem:** what if the information is there but not linear?

Accuracy of SG, CBOW and subword information skip-gram models on word analogy tasks for Czech, German, English and Italian:

		sg	cbow	sisg
Cs	Semantic	25.7	27.6	27.5
	Syntactic	52.8	55.0	77.8
De	Semantic	66.5	66.8	62.3
	Syntactic	44.5	45.0	56.4
En	Semantic	78.5	78.2	77.8
	Syntactic	70.1	69.9	74.9
It	Semantic	52.3	54.7	52.3
	Syntactic	51.5	51.8	62.7

Figure 3.8: Word analogy task, where sg - skip-gram, cbow - continuous bag of words, sisg - subword information skip-gram (treat unseen words by summing the n-gram vectors)

It is observed that morphological information from SISG model significantly improves the syntactic task. In contrast it doesn't help for semantic question, and even degrades the performance for German and Italian.

Another intrinsic word evaluation is word vector distances and their correlation with human judgments.

Example (dataset WordSim353):

Word 1	Word 2	Human (mean)
tiger	tiger	10
tiger	cat	7.35
plane	car	5.77
stock	phone	1.62
stock	jaguar	0.92

Performance different models:

### 3.4.2 Extrinsic word vector evaluation

Extrinsic word vector evaluation: all subsequent task in this class. One example where good word vectors should help directly is named entity recognition (NER): finding a person, organisation or location.

## 3.5 Word senses and word sense ambiguity

Most words have lots of meaning. Especially common words or words that have existed for a long time. Does one vector capture all these meanings or do we have a mess?

		sg	cbow	sisg-	sisg
AR	WS353	51	52	54	<b>55</b>
GUR350		61	62	64	<b>70</b>
DE	GUR65	78	78	<b>81</b>	<b>81</b>
ZG222		35	38	41	<b>44</b>
EN	RW	43	43	46	<b>47</b>
	WS353	72	<b>73</b>	71	71
ES	WS353	57	58	58	<b>59</b>
FR	RG65	70	69	<b>75</b>	<b>75</b>
RO	WS353	48	52	51	<b>54</b>
RU	HJ	59	60	60	<b>66</b>

Figure 3.9: Human similarity judgement, where sg - skip-gram, cbow - continuous bag of words, sisg- - subword information skip-gram (treat unseen words as a null vector), sisg - subword information skip-gram (treat unseen words by summing the n-gram vectors)

We can cluster word windows around words, retrain with each word assigned to multiple different clusters  $bank_1, bank_2, \dots$

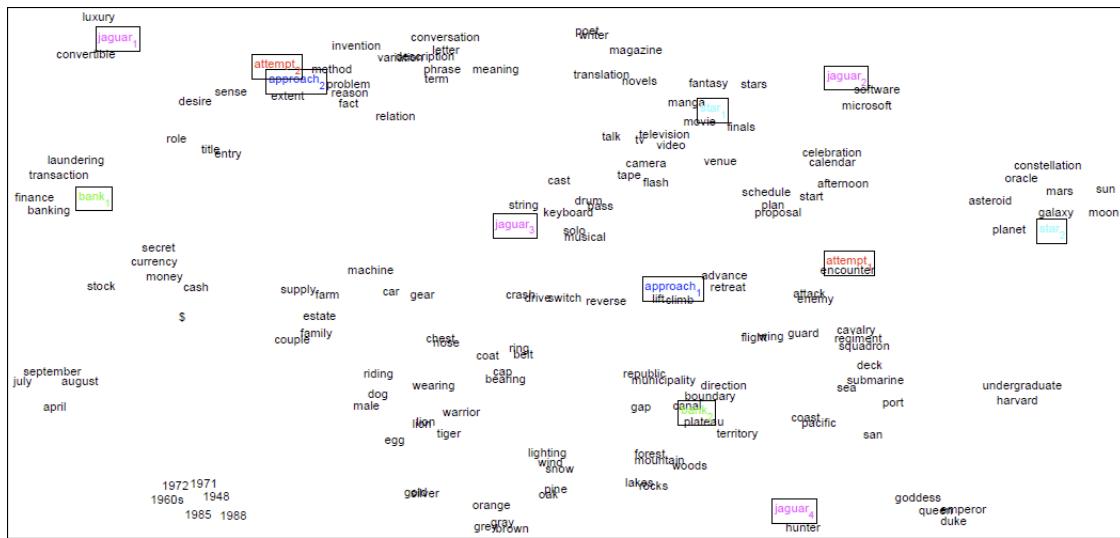


Figure 3.10: Example of clusters

We can use linear algebra. Different senses of a word reside in a linear superposition (weighted sum) in standard word embeddings like word2vec. For example:  $v_{pike} = \alpha_1 v_{pike_1} + \alpha_2 v_{pike_2} + \alpha_3 v_{pike_3}$ , where  $\alpha_i = \frac{f_i}{\sum_i f_i}$  and  $f$  - frequency.

Surprising result: because of ideas from sparse coding you can actually separate out the senses (providing they are relatively common):

tie				
trousers	season	scoreline	wires	operatic
blouse	teams	goalless	cables	soprano
waistcoat	winning	equaliser	wiring	mezzo
skirt	league	clinching	electrical	contralto
sleeved	finished	scoreless	wire	baritone
pants	championship	replay	cable	coloratura

### 3.6 References

---

1. T Mikolov, I Sutskever, K Chen, GS Corrado, J Dean, Distributed representations of words and phrases and their compositionality, NIPS 2013
2. A Joulin, É Grave, P Bojanowski, T Mikolov, Bag of Tricks for Efficient Text Classification. EACL 2017.
3. Huang et al., Improving word representations via global context and multiple word prototypes, 2012
4. Arora et al., Linear algebraic structure of word senses, with applications to polysemy, TACL 2018
5. Piotr Bonjanowski, Edouard Grave, Armand Joulin, Tomas Mikolov, Enriching Word Vectors with Subword Information

# 4. Neural Networks Basics

## 4.1 Neural networks

62

- 4.1.1 Gradient descent
- 4.1.2 Notation in vectorized form: Jacobian Matrix
- 4.1.3 Capacity of multilayer perceptron
- 4.1.4 Bag of words
- 4.1.5 Handling Combinations
- 4.1.6 Convolutional Neural Networks
- 4.1.7 Stacked Convolution and Dilated Convolution

*In this chapter we discuss the basics of neural networks. Where they come from, how to build and optimize them.*

### 4.1 Neural networks

Logistic regression is a linear classifier, which means its decision boundary is a hyperplane:

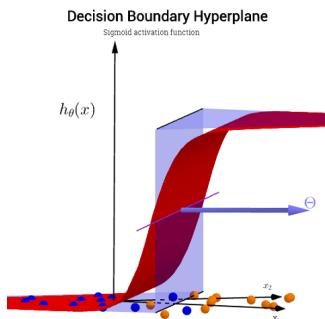


Figure 4.1: Limitation of linear classifiers

However, some of the classification problems are not linear separable:

A typical linear-inseparable problem is the exclusive-or function:

Neural networks (NNs), or artificial neural networks (ANNs), provide an efficient way to solve non-linear-separable problem.

Artificial neural networks (ANNs) or connectionist systems are computing systems vaguely inspired by the biological neural networks that constitute animal brains. – Wikipedia

An ANN is based on a collection of connected units or nodes called artificial neurons, which loosely model the neurons in a biological brain. Each connection, like the synapses in

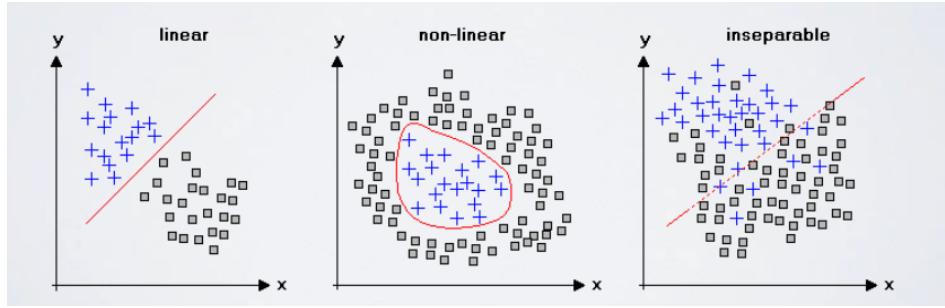


Figure 4.2: Limitation of linear classifiers

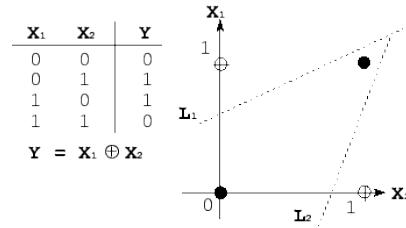


Figure 4.3: Exclusive-or function

a biological brain, can transmit a signal to other neurons. An artificial neuron that receives a signal then processes it and can signal neurons connected to it.

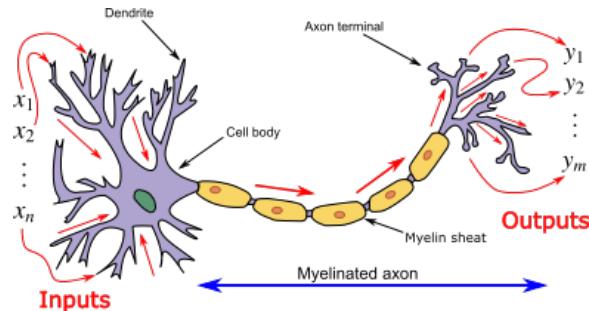


Figure 4.4: Neuron and myelinated axon, with signal flow from inputs at dendrites to outputs at axon terminals

Simple threshold function for artificial neuron is:

$$f(x) = \begin{cases} 0, & \text{if } \sum_i w_i x_i + b \leq 0 \\ 1, & \text{otherwise} \end{cases}$$

Perceptron is also an effective linear classifier but we cannot apply gradient descent algorithm to train it. Logistic regression can be regarded as improved version of perceptron to which gradient descent training can be applied.

Perceptrons (including Logistic regression classifiers) are linear classifiers, which do not work for non-linear-separable problems. Multilayer perceptrons (MLP) are proposed to solve non-linear-separable problems.

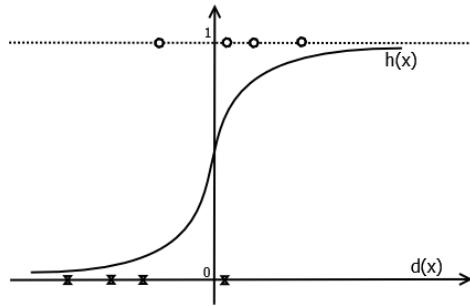


Figure 4.5: Improved perceptron – logistic regression

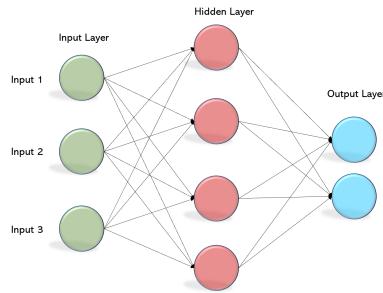


Figure 4.6: From linear classifier to non-linear classifier

- An MLP consists of at least three layers of nodes: an input layer, a hidden layer and an output layer.
- Except for the input nodes, each node is a neuron that uses a nonlinear activation function.
- MLP utilizes a supervised learning technique called backpropagation for training.
- Its multiple layers and non-linear activation distinguish MLP from a linear perceptron.
- It can distinguish data that is not linearly separable.

$$\text{Logistic: } f(x) = \frac{1}{1 + e^{-x}}$$

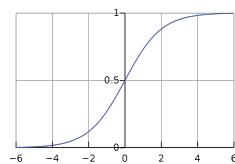


Figure 4.7: Activation functions: logistic

$$\text{Hyperbolic tangent: } f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\text{Rectifier: } f(x) = \max(0, x)$$

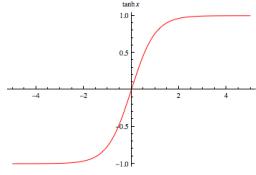


Figure 4.8: Activation functions: hyperbolic tangent

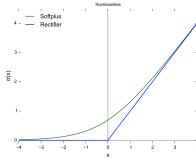


Figure 4.9: Activation functions: rectifier

The MLP consists of three or more layers (an input and an output layer with one or more hidden layers) of nonlinearly-activating nodes. Since MLPs are fully connected, each node in one layer connects with a certain weight  $w_{ij}$  to every node in the following layer. Output layer may use active functions differ from hidden layers, depending on the task nature. For classification problems, typically, a softmax function is used in the output layer.

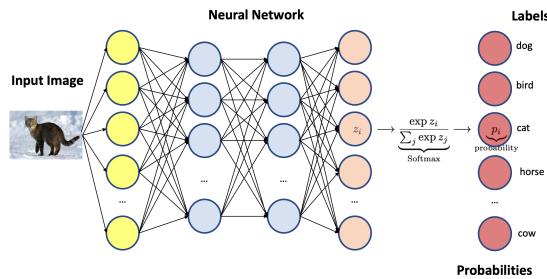


Figure 4.10: Output layer

### 4.1.1 Gradient descent

Weight space for a multilayer neural net: one coordinate for each weight or bias of the network, in all the layers. Conceptually, not any different from what we've seen so far - just higher dimensional and harder to visualize. We want to compute the cost gradient, which is the vector of partial derivatives. This is the average of  $dL/dw$  over all training examples, so we focus on computing  $dL/dw$

We've already been using the univariate Chain Rule. If  $f(x)$  and  $x(t)$  are univariate functions, then  $\frac{d}{dt}f(x(t)) = \frac{df}{dx} \frac{dx}{dt}$ .

Let's compute the loss derivatives for univariate logistic least squares model:

Computing the loss:

$$z = wx + b, y = \sigma(z), L = \frac{1}{2}(y - t)^2$$

$$L = \frac{1}{2}(\sigma(wx + b) - t)^2$$

Computing the derivatives:

$$\begin{aligned}\frac{\partial L}{\partial w} &= \frac{\partial}{\partial w} \left[ \frac{1}{2}(\sigma(wx + b) - t)^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial w} (\sigma(wx + b) - t)^2 \\ &= (\sigma(wx + b) - t) \frac{\partial}{\partial w} (\sigma(wx + b) - t) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) \frac{\partial}{\partial w} (wx + b) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) x\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial b} &= \frac{\partial}{\partial b} \left[ \frac{1}{2}(\sigma(wx + b) - t)^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial b} (\sigma(wx + b) - t)^2 \\ &= (\sigma(wx + b) - t) \frac{\partial}{\partial b} (\sigma(wx + b) - t) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) \frac{\partial}{\partial b} (wx + b) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b)\end{aligned}$$

The disadvantages of this approach are complex form of solution and inefficient computations. A more structured and efficient way to do it:

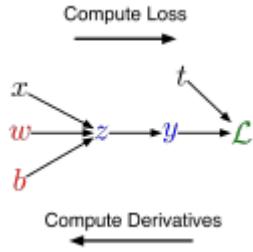
$$z = wx + b, y = \sigma(z), L = \frac{1}{2}(y - t)^2$$

$$\begin{aligned}\frac{dL}{dy} &= y - t \\ \frac{dL}{dz} &= \frac{dL}{dy} \sigma'(z) \\ \frac{\partial L}{\partial w} &= \frac{dL}{dz} x \\ \frac{\partial L}{\partial b} &= \frac{dL}{dz}\end{aligned}$$

We can diagram out the computations using a *computational graph*. The nodes represent all the inputs and computed quantities, and the edges represent which nodes are computed directly as a function of which other nodes.

A slightly more convenient notation:

$\bar{y}$  is the derivative  $dL/dy$ , sometimes it's called the *error signal*. This emphasizes that the error signals are just values our program is computing (rather than mathematical operation). This is not standard notation.



$$\begin{array}{l|l} \begin{aligned} z &= wx + b \\ y &= \sigma(z) \\ L &= \frac{1}{2}(y - t)^2 \end{aligned} & \begin{aligned} \bar{y} &= y - t \\ \bar{z} &= \bar{y}\sigma'(z) \\ \bar{w} &= \bar{z}x \\ \bar{b} &= \bar{z} \end{aligned} \end{array}$$

If the computation graph has fan-out  $> 1$ , this requires the multivariate Chain Rule.

Suppose we have a function  $f(x, y)$  and functions  $x(t)$  and  $y(t)$ . All the variables are scalar-valued. Then:

$$\frac{d}{dt}f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

Example:  $f(x, y) = y + e^{xy}$ ,  $x(t) = \cos t$ ,  $y(t) = t^2$ . Plug in to Chain Rule:

$$\frac{d}{dt}f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} = (ye^{xy})(-\sin t) + 2t(1 + xe^{xy})$$

In the context of backpropagation  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}$  are the values already computed by our program,  $\frac{dx}{dt}, \frac{dy}{dt}$  are the mathematical expressions to be evaluated.

In error signal notation:

$$\bar{t} = \bar{x} \frac{dx}{dt} + \bar{y} \frac{dy}{dt}$$

Full backpropagation algorithm:

Let  $v_1, \dots, v_N$  be a topological ordering of the computational graph, i.e. parents come before children.  $v_N$  denotes the variable we're trying to compute derivatives of (e.g. loss).

- Forward pass: compute  $v_i$  as a function of  $Pa(v_i)$  for  $i = 1, \dots, N$
- Backward pass:  $\bar{v}_i = \sum_{j \in Ch(v_i)} \bar{v}_j \frac{\partial v_j}{\partial v_i}$  for  $i = N - 1, \dots, 1$

### Backprop as message passing

Each node receives a bunch of messages from its children, which it aggregates to get its error signal. It then passes messages to its parents. This provides modularity, since each node only has to know how to compute derivatives with respect to its arguments, and doesn't have to know anything about the rest of the graph.

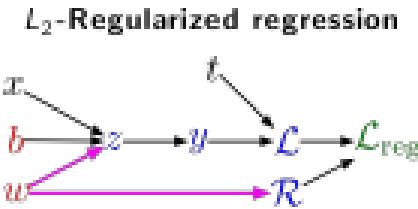
### Computational cost

Computational cost of forward pass: one add-multiply operation per weight.

Computational cost of backward pass: two add-multiply operations per weight. The backward pass is about as expensive as two forward passes.

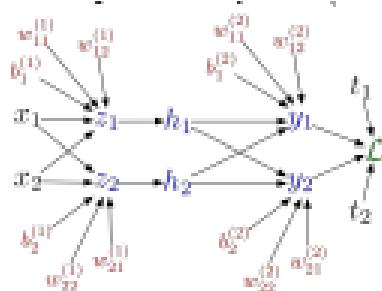
For a multilayer perceptron, this means the cost is linear in the number of layers , quadratic in the number of units per layer.

Backprop is used to train the overwhelming majority of neural nets today. Even optimization algorithms much fancier than gradient descent (e.g. second-order methods) use backprop to compute the gradients. Despite its practical success, backprop is believed to be neurally implausible. There is no evidence for biological signals analogous to error derivatives, and all the biologically plausible alternatives we know about learn much more slowly on computers.



<p>Forward pass:</p> $z = wx + b$ $y = \sigma(z)$ $L = \frac{1}{2}(y - t)^2$ $R = \frac{1}{2}w^2$ $L_{reg} = L + \lambda R$	<p>Backward pass:</p> $\overline{L}_{reg} = 1$ $\overline{R} = \overline{L}_{reg} \frac{dL_{reg}}{dR} = \overline{L}_{reg} \lambda$ $\overline{L} = \overline{L}_{reg} \frac{dL_{reg}}{dL} = \overline{L}_{reg}$ $\overline{y} = \overline{L} \frac{dL}{dy} = \overline{L}(y - t)$ $\overline{z} = \overline{y} \frac{dy}{dz} = \overline{y}\sigma'(z)$ $\overline{w} = \overline{z} \frac{\partial z}{\partial w} + \overline{R} \frac{dR}{dw} = \overline{z}x + \overline{R}w$ $\overline{b} = \overline{z} \frac{\partial z}{\partial b} = \overline{z}$
---	---

### Multilayer perceptron:



Forward pass: $z_i = \sum_j w_{ij}^{(1)} x_j + b_j^{(1)}$ $h_i = \sigma(z_i)$ $y_k = \sum_i w_{ki}^{(2)} h_i + b_k^{(2)}$ $L = \frac{1}{2} \sum_k (y_k - t_k)^2$	Backward pass: $\bar{L} = 1$ $\bar{y}_k = \bar{L}(y_k - t_k)$ $\overline{w_{ki}^{(2)}} = \bar{y}_k h_i$ $\overline{b_k^{(2)}} = \bar{y}_k$ $\bar{h}_i = \sum_k \bar{y}_k \overline{w_{ki}^{(2)}}$ $\bar{z}_i = \bar{h}_i \sigma'(z_i)$ $\overline{w_{ij}^{(1)}} = \bar{z}_i x_j$ $\overline{b_i^{(1)}} = \bar{z}_i$
---	--

**Backpropagation in vectorized form:**

Forward pass: $\mathbf{z} = \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}$ $\mathbf{h} = \sigma(\mathbf{z})$ $\mathbf{y} = \mathbf{W}^{(2)} \mathbf{h} + \mathbf{b}^{(2)}$ $L = \frac{1}{2} \ \mathbf{t} - \mathbf{y}\ ^2$	Backward pass: $\bar{L} = 1$ $\bar{\mathbf{y}} = \bar{L}(\mathbf{y} - \mathbf{t})$ $\overline{\mathbf{W}^{(2)}} = \bar{\mathbf{y}} \mathbf{h}^T$ $\overline{\mathbf{b}^{(2)}} = \bar{\mathbf{y}}$ $\bar{\mathbf{h}} = (\mathbf{W}^{(2)})^T \bar{\mathbf{y}}$ $\bar{\mathbf{z}} = \bar{\mathbf{h}} \circ \sigma'(\mathbf{z})$ $\overline{\mathbf{W}^{(1)}} = \bar{\mathbf{z}} \mathbf{x}^T$ $\overline{\mathbf{b}^{(1)}} = \bar{\mathbf{z}}$
--	--

#### 4.1.2 Notation in vectorized form: Jacobian Matrix

Previously, we derived backprop equations in terms of sums and indices, and then vectorized them. But we'd like to implement our primitive operations in vectorized form.

The Jacobian is the matrix of partial derivatives:

$$\mathbf{J} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

The backprop equation (single child node) can be written as a vector-Jacobian product (VJP):

$$\bar{x}_j = \sum_i \bar{y}_i \frac{\partial y_i}{\partial x_j}$$

$$\bar{\mathbf{x}} = \bar{\mathbf{y}}^T \mathbf{J}$$

That gives a row vector. We can treat it as a column vector by taking  $\bar{\mathbf{x}} = \mathbf{J}^T \bar{\mathbf{y}}$

We never explicitly construct the Jacobian. It's usually simpler and more efficient to compute the VJP directly.

### Examples

Matrix-vector product:  $\mathbf{z} = \mathbf{Wx}, \mathbf{J} = \mathbf{W}, \bar{\mathbf{x}} = \mathbf{W}^T \bar{\mathbf{z}}$

Elementwise operations:

$$\mathbf{y} = \exp(\mathbf{z})$$

$$\mathbf{J} = \begin{pmatrix} \exp(z_1) & & 0 \\ & \ddots & \\ 0 & & \exp(z_D) \end{pmatrix}$$

$$\bar{\mathbf{z}} = \exp(\mathbf{z}) \circ \bar{\mathbf{y}}$$

### 4.1.3 Capacity of multilayer perceptron

**Universal approximation theorem** (Hornik, 1991):

A single hidden layer neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden units.

The result applies for a sigmoid, tanh and many other hidden layer activation functions.

This is the good result, but it doesn't mean there is a learning algorithm than can find the necessary parameter values.

### 4.1.4 Bag of words

We have a sentence classification task. Some items are simple:

We can apply Bag of Words (BOW):

We can get some probabilities, but the following example will be classified incorrectly

Continuous Bag of Words (CBOW):

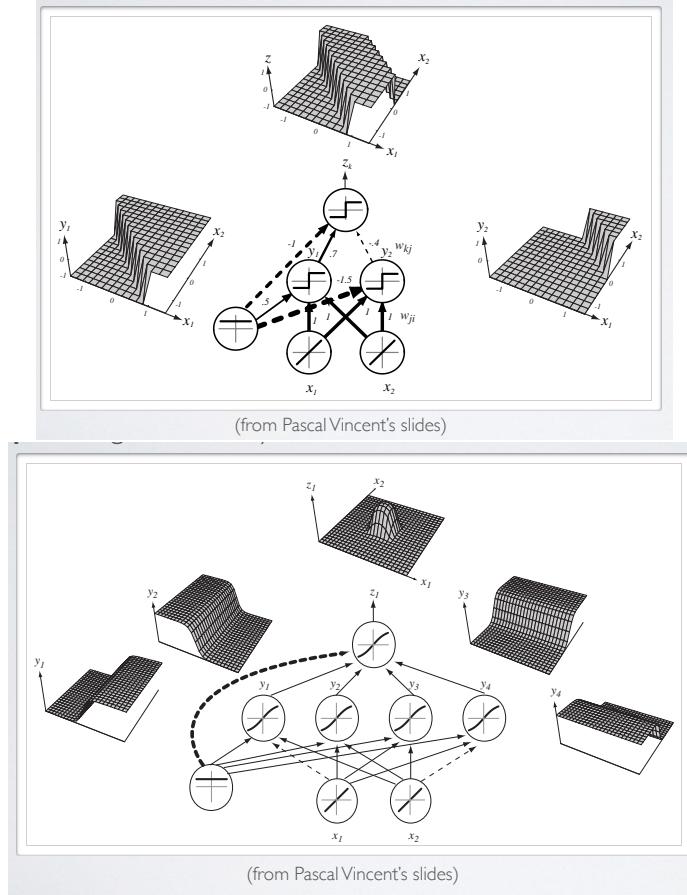


Figure 4.11: Capacity of multilayer perceptron

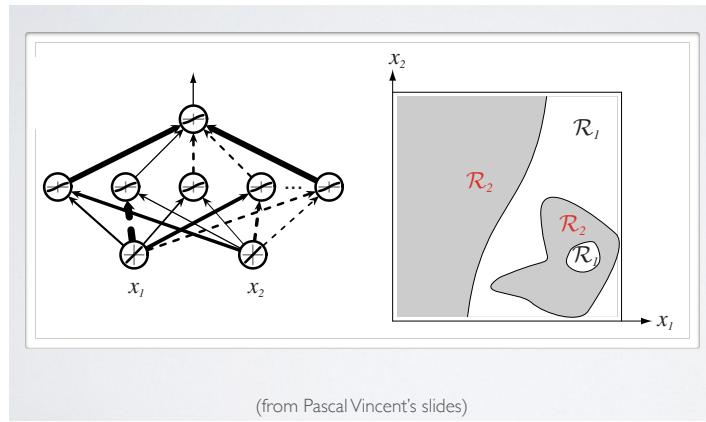
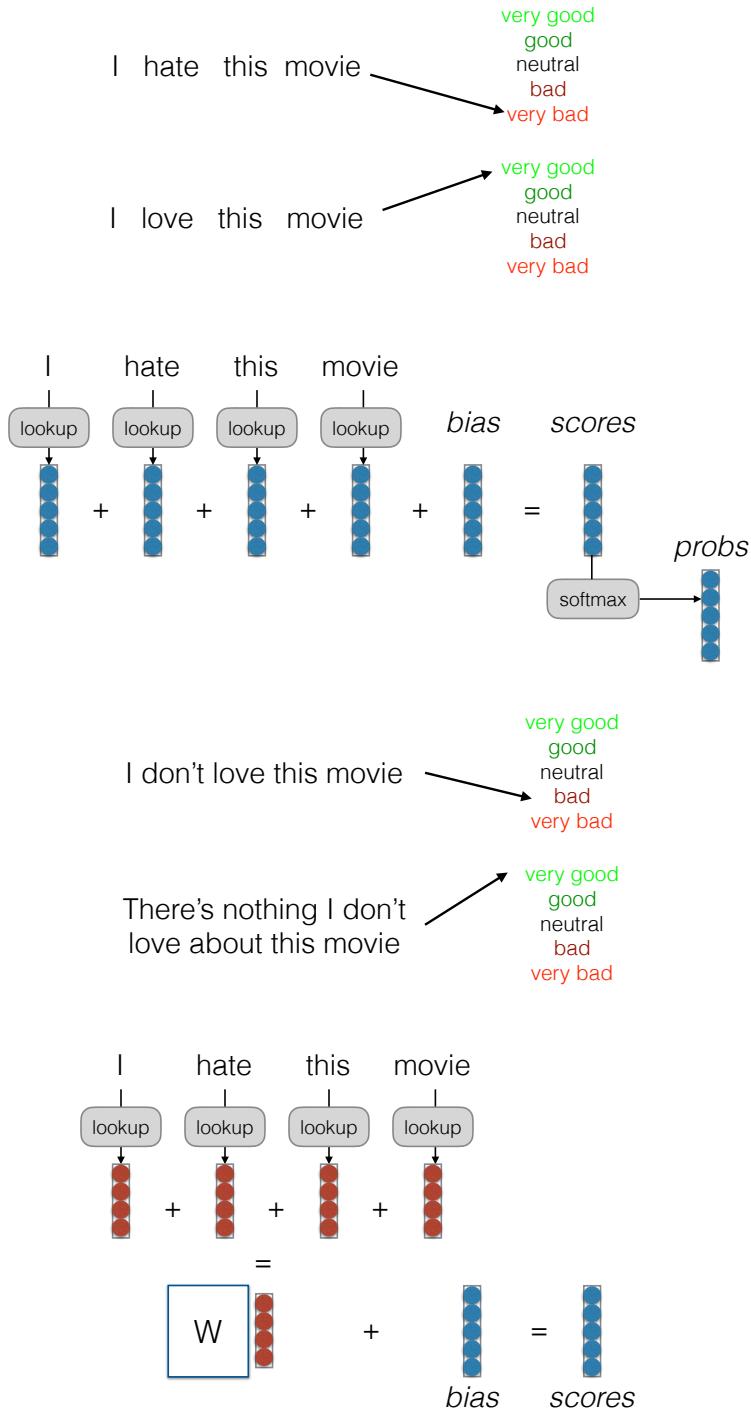


Figure 4.12: Capacity of multilayer perceptron

### Deep CBOW

We can learn feature combinations: a node in a second layer might be "feature 1 AND feature 5 are active". It captures things such as "not" AND "hate", but cannot handle "not hate"

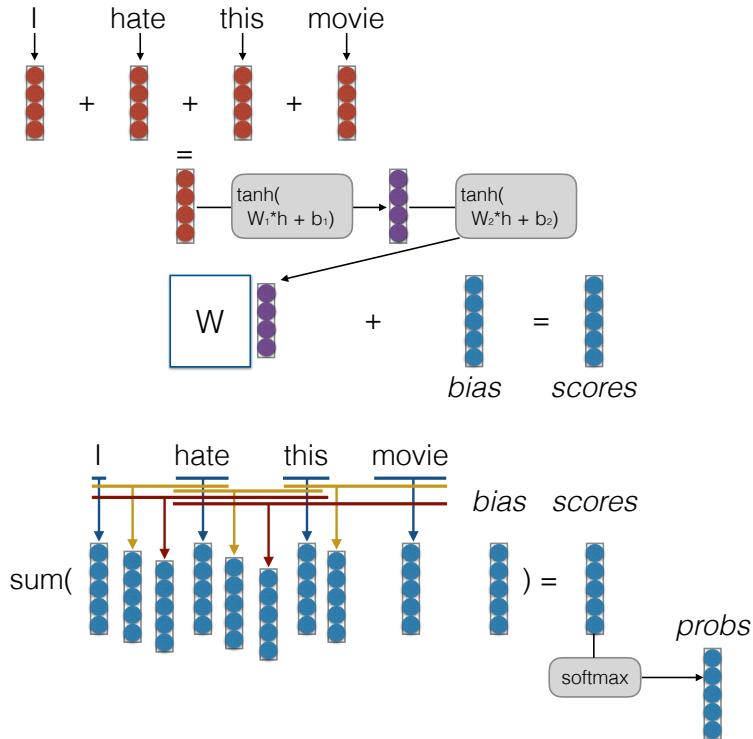


#### 4.1.5 Handling Combinations

We can apply bag of n-grams:

It allows us to capture combination features in a simple way "don't love", "not the best", and it works pretty well.

Problems are the same as before: parameter explosion and no sharing between similar words/n-grams.



#### 4.1.6 Convolutional Neural Networks

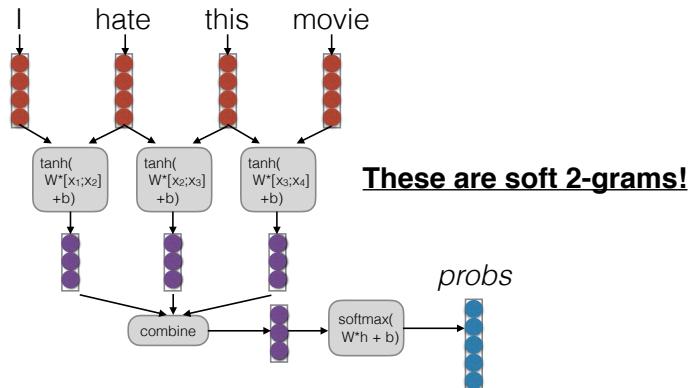


Figure 4.13: 1-dimensional Convolutions/Time-delay networks (Waibel et al. 1989)

CNNs for Text (Collobert and Weston 2011):

It generally based on 1D convolutions, but often uses terminology/functions borrowed from image processing for historical reasons.

Two main paradigms:

- **Context window modeling:** get the surrounding context before tagging
- **Sentence modeling:** do convolution to extract n-grams, pooling to combine over whole sentence.

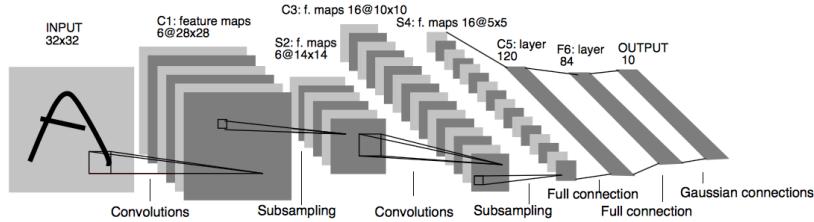


Figure 4.14: 2-dimensional Convolutional networks (LeCun et al. 1997) for image processing. Parameter extraction performs a 2D sweep, not 1D

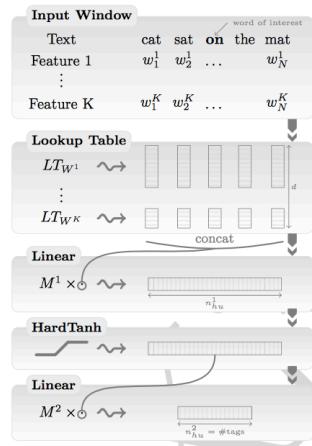


Figure 4.15: CNNs for Tagging

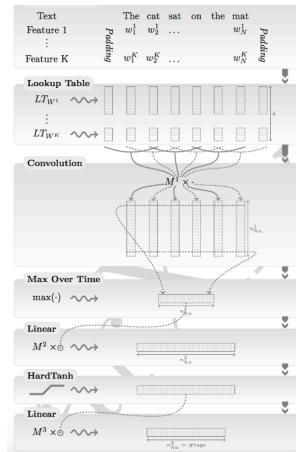


Figure 4.16: CNNs for sentence modeling

Standard 2D convolution function takes input + parameters. Input is a 3D tensor: rows (words), columns, features ("channels"). Parameters/filters are 4D tensor: rows, columns, input features, output features.

**Padding.** After convolution, the rows and columns of the output tensor are either:

- = to the rows/columns of the input tensor ("same" convolution)
- = to rows/columns of the input tensor minus the size of the filter plus one ("valid" or "narrow")
- = to rows/columns of input tensor plus filter minus one ("wide")

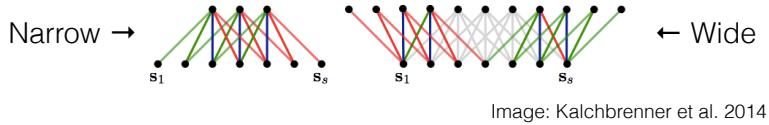


Figure 4.17: Padding

**Striding.** Skip some of the outputs to reduce length of extracted feature vector.

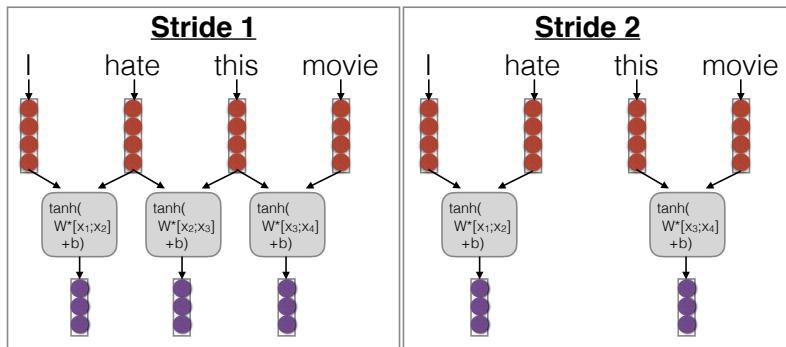


Figure 4.18: Striding

**Pooling.** Pooling is like convolution, but calculates some reduction function feature-wise.

- **Max pooling:** "Did you see feature anywhere in the range?" (most common)
- **Average pooling:** "How prevalent is this feature over the entire range?"
- **k-Max pooling:** "Did you see this feature up to k times?"
- **Dynamic pooling:** "Did you see this feature in the beginning? In the middle? In the end?"

#### 4.1.7 Stacked Convolution and Dilated Convolution

**Stacked convolution.** Feeding in convolution from previous layer results in larger area of focus for each feature.

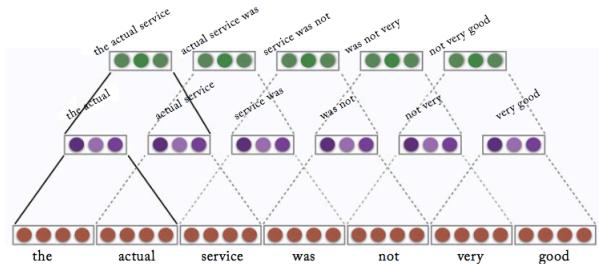


Image Credit: Goldberg Book

Figure 4.19: Stacked convolution

**Dilated convolution.** Kalchbrenner et al. 2016. Gradually increase stride, every time step (no reduction in length).

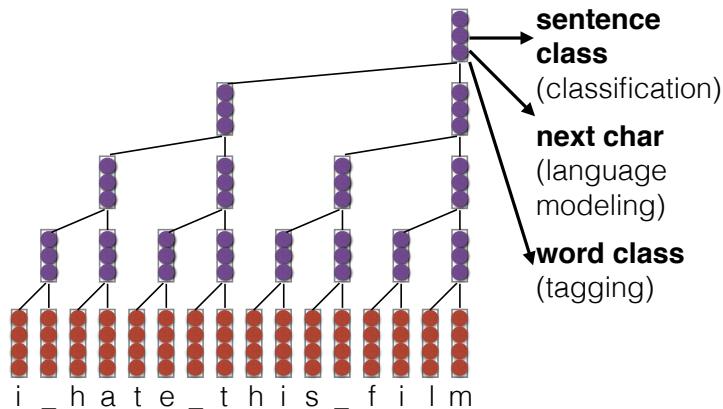


Figure 4.20: Dilated convolution

Why (dilated) convolution for modeling sentences? In contrast to RNNs it requires fewer steps from each word to the final representation ( $O(N)$  for RNN,  $O(\log N)$  for dilated CNN) and easier to parallelize on GPU.

Problems:

- Slightly less natural for arbitrary-length dependencies
- A bit slower on CPU

**Iterated dilated convolution.** Strubell+ 2017. Multiple iterations of the same stack of dilated convolutions. Wider context, more parameter efficient.

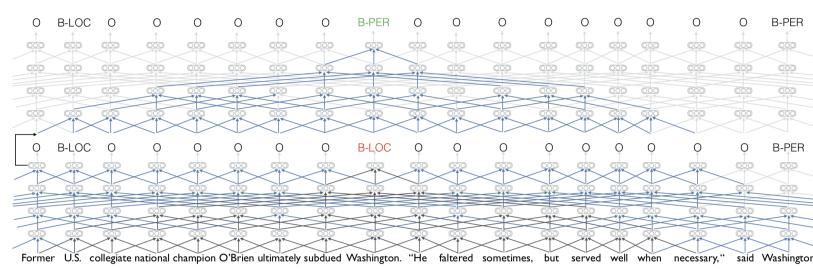


Figure 4.21: Iterated dilated convolution

# 5. Tagging

In this chapter we discuss the basics of neural networks. Where they come from, how to build and optimize them.

## 5.1 Parts Of Speech (POS)

5.1	Parts Of Speech (POS)	78
5.2	POS tagging	79
5.3	POS tagging as a sequence classification task	80
5.4	Named Entity Recognition (NER)	81
5.5	NER applications	82
5.6	The NER task	82
5.7	Three standard approaches to NER	82
5.8	Encoding classes for sequence labeling	85
5.9	Sequence Labeling as Classification Using Outputs as Inputs <i>Backward Procedure</i>	87
5.10	Hidden Markov Models (HMMs)	90
5.11	Inference in HMM	91
5.12	Forward procedure for HMM <i>Backward Procedure</i>	92
5.13	Parameter Estimation <i>Backward Procedure</i>	93
5.14	Decoding Solution Summary <i>Backward Procedure</i>	95
5.15	Viterbi algorithm <i>Backward Procedure</i>	95
5.16	Parameter Estimation <i>Backward Procedure</i>	97
	Open class (lexical) words	98
	Nouns Proper <i>IBM</i> <i>Italy</i>	98
	Common <i>cat / cats</i> <i>snow</i>	98
	Verbs Main <i>see</i> <i>registered</i>	98
	Adjectives <i>old</i> <i>older</i> <i>oldest</i>	98
	Adverbs <i>slowly</i>	99
	Numbers <i>122,312</i> <i>one</i>	99
	... more	99
	Closed class (functional)	100
	Determiners <i>the</i> <i>some</i>	100
	Conjunctions <i>and</i> <i>or</i>	100
	Pronouns <i>he</i> <i>its</i>	100
	Prepositions <i>to</i> <i>with</i>	100
	Particles <i>off</i> <i>up</i>	100
	Interjections <i>Ow</i> <i>Eh</i>	100
	... more	101

5.25	Shortcomings of Hidden Markov Model(1): locality of features <i>Backward Procedure</i>	102
5.26	MEMM	102
5.27	From MEMM to CRF <i>Backward Procedure</i>	106
5.28	CRF: Inference	106
5.29	CRF learning	107

## 5.2 POS tagging

---

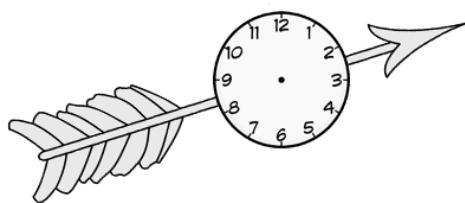


Figure 5.2: \*  
When  $\text{POS}(\text{'files'}) = \text{VB}$  (verb)

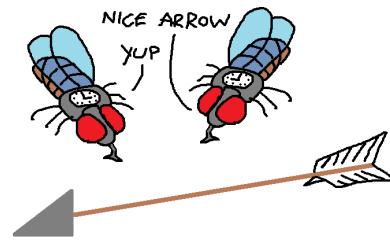


Figure 5.3: \*  
When  $\text{POS}(\text{'files'}) = \text{NN}$  (noun)

Figure 5.4: \*  
An example: Time flies like an arrow

The POS tagging problem is to determine the POS tag for a particular instance of a word. It is complicated by the fact that words often have more than one POS, for example word "back":

- The back door = JJ
- On my back = NN
- Win the votes back = RB
- Promise to back the bill

Input:	Plays	well	with	others
Ambiguity:	NNS/VBZ	UH/JJ/NN/Rb	IN	NNS
Output:	Plays/VBZ	well/RB	with/IN	others/NNS

### Uses:

- MT: reordering of adjectives and nouns (say from Spanish to English)
- Text-to-speech (how do we pronounce "lead")?
- Can write regexps like (Det) Adj\* N+ over the output for phrases, etc.
- Input to a syntactic parser

**How difficult is POS tagging?** About 11% of the word types in the Brown corpus are ambiguous with regard to part of speech. But they tend to be very common word. E.g., **that**:

- I know **that** he is honest = IN
- Yes, **that** play was nice = DT
- You can't go **that** far = RB

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	+%, &
CD	cardinal number	<i>one, two</i>	TO	"to"	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential 'there'	<i>there</i>	VB	verb base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, sing.	<i>IBM</i>	\$	dollar sign	\$
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	#
PDT	predeterminer	<i>all, both</i>	"	left quote	' or "
POS	possessive ending	<i>'s</i>	"	right quote	" or '
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	[, {, <
PRPS	possessive pronoun	<i>your, one's</i>	)	right parenthesis	], }, >
RB	adverb	<i>quickly, never</i>	,	comma	,
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	. ! ?
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	: ; ... - -
RP	particle	<i>up, off</i>			

Figure 5.5: \*  
The Penn TreeBank tagset

40% of the word tokens are ambiguous.

What are the main sources of information for POS tagging?

**Knowledge of neighboring words:** Bill saw that man yesterday  
 NNP NN DT NN NN  
 VB VB(D) IN VB NN

### Feature-based tagger

- Can do surprisingly well just looking at a word by itself:  
 Word the: the  $\Rightarrow$  DT  
 Lowercased word Importantly: Importantly  $\Rightarrow$  RB  
 Prefixes unfathomable: un-  $\Rightarrow$  JJ  
 Suffixes Importantly: -ly  $\Rightarrow$  RB  
 Capitalization Meridian: CAP  $\Rightarrow$  NNP  
 Word shapes 35-years: d-x  $\Rightarrow$  JJ
- Then build a classifier to predict tag. Maxent  $P(t|w)$ : 93.7% overall / 82.6% unknown

## 5.3 POS tagging as a sequence classification task

Let's say we are given a sentence (an "observation" or "sequence of observations") like:

- Secretariat is expected to race tomorrow
- She promised to back the bill

What is the best sequence of tags which corresponds to this sequence of observations? If we use probabilistic view we will:

- Consider all possible sequence of tags
- Out of this sequences, choose the tag sequence which is most probable for given the observation sequence of n words  $w_1 \dots w_n$ .

## 5.4 Named Entity Recognition (NER)

---

A very important sub-task is to *find* and *classify* names in text, lets find named entities and classify them in text below:

The decision by the independent MP **Andrew Wilkie** to withdraw his support for the minority **Labor** government sounded dramatic but it should not further threaten its stability. When, after the **2010** election, **Wilkie, Rob, Oakeshott, Tony Windsor** and the **Greens** agreed to support **Labor**, they gave just two guarantees: confidence and supply.

Classes are marked this way: **Person, Date, Location, Organization**.

```
<PER>Prof. Jerry Hobbs</PER> taught CS544 during <DATE>February 2010</DATE>.  
<PER>Jerry Hobbs</PER> killed his daughter in <LOC>Ohio</LOC>.  
<ORG>Hobbs corporation</ORG> bought <ORG>FbK</ORG>.
```

Figure 5.6: Example of NER dataset markup

NER task report can look this way:

- Person Names: Prof. Jerry Hobbs, Jerry Hobbs
- Organisations: Hobbs corporation, FbK
- Locations: Ohio
- Date and time expressions: February 2010
- E-mail: mkg@gmail.com
- Web address: www.usc.edu
- Names of drugs: paracetamol
- Names of ships: Queen Mary
- ...

## 5.5 NER applications

---

The uses:

- Named entities can be indexed, linked off, etc.
- Sentiment can be attributed to companies or products.
- A lot of IE relations are associations between named entities.
- For question answering, answers are often named entities.

Concretely:

- Many web pages tag various entities, with links to bio or topic pages, etc: Reuters' OpenCalais, Evri, AlchemyAPI, Yahoo's Term Extraction, ...
- Apple/Google/Microsoft/... smart recognizers for document content.

## 5.6 The NER task

---

Task: Predict entities in a text:

Foreign	ORG	
Ministry	ORG	
spokesman	O	Standard
Shen	PER	
Guofang	PER	
told	O	evaluation is per entity, <i>not</i> per token
Reuters	ORG	
:	:	

## 5.7 Three standard approaches to NER

---

1. Hand-written regular expressions
  - Perhaps stacked
2. Using classifiers
  - Generative: Naive Bayes

- Discriminative: Maxent model
3. Sequence models
- HMMs
  - CMMs/MEMMs
  - CRFs

\*Generative - observes whole data and labels

\*Discriminative - observes only some distributions on input date

### **Rule based NER**

With regular expressions you can extract telephone numbers, e-mails, Capitalized names.

Regular expressions provide a flexible way to match strings of text, such as particular characters, words or patterns of characters. Suppose, you are looking for a word that:

1. starts with a capital letter "P"
2. is the first word of a line
3. the second letter is a lower case letter
4. is exactly three letters long
5. the third letter is a vowel

the regular expression would be ^P[a-z][aeiou], where: ^- indicates the beginning for the string, [a-z] - any letter in the range a to z, [aeiou] - any vowel.

### **Special perl RegEx:**

- \w (word char) any alpha-numeric
- \d (digit char) any digit
- \s (space char) any whitespace
- . (wildcard) anything
- \b word boundary
- ^ beginning of string
- \$ end of string
- ? for 0 or 1 occurrences

- + for 1 or more occurrences
- specific range of number of occurrences: **min, max**
  - A1,5 One to five A's
  - A5, Five or more A's
  - A5 Exactly five A's

**Example: blocks of digits separated by hyphens - RegEx =  $(\d\{-\})+\d{+}$ :**

- matches valid phone numbers like 900-865-1125 and 7225-1234
- incorrectly extracts social security numbers 123-45-6789
- fails to identify numbers like 800.865.1125 and (800)865-CARE

**Improved RegEx =  $(\d\{3\}[-.\s])\{1,2\}[\dA-Z]\{4\}$**

**As for locations extraction:**

- Capitalized word + {city, center, river} indicates location. Ex.: New York city, Hudson river
- Capitalized word + {street, boulevard, avenue} indicates location. Ex.: Fifth avenue.

**Use context patterns:**

- [PERSON] earned [MONEY] - Frank earned \$20
- [PERSON] joined [ORGANIZATION] - Sam joined IBM
- [PERSON], [JOBTITLE] - Mary, the teacher
- [PERSON|ORGANIZATION|ANIMAL] fly to [LOCATION|PERSON|EVENT] - Jerry flies to the party / Delta flies to Europe / bird flies to trees / bee flies to the wood

**But simple things often doesn't work. Example - capitalization is a very tricky:**

- first word of a sentence is capitalized
- sometimes titles in web pages are all capitalized
- nested named entities contain non-capital words: University of Southern California is Organisation
- all nouns in German are capitalized

***The ML sequence model approach to NER***

Training

1. Collect a set of representative training documents
2. Label each token for its entity class or other (O)
3. Design feature extractors appropriate to the text and classes
4. Train a sequence classifier to predict the labels from the data

### Testing

1. Receive a set of testing documents
2. Run sequence model inference to label each token
3. Appropriately output the recognized entities

## 5.8 Encoding classes for sequence labeling

---

Many NLP tasks are sequence labeling tasks.

	IO encoding	IOB encoding
Fred	PER	B-PER
showed	O	O
Sue	PER	B-PER
Mengqiu	PER	B-PER
Huang	PER	I-PER
's	O	O
new	O	O
painting	O	O

**Input:** a sequence of tokens/words: *Pierre Vinken, 61 years old, will join IBM's board as a nonexecutive director Nov. 29.*

**Output:** a sequence of labeled tokens/words:

**POS-tagging:** Pierre\_NNP Vinken\_NNP ,\_CD years\_NNS old\_JJ, will\_MD join\_VB IBM\_NNP 's\_POS board\_NN as\_IN a\_DT nonexecutive\_JJ director\_NN Nov.\_NNP 29\_CD ...

**Named Entity Recognition:** Pierre\_B-PERS Vinken\_I-PERS ,\_O 61\_O years\_O old\_O, will\_O join\_O IBM\_B-ORG 's\_O board\_O as\_O a\_O nonexecutive\_O director\_O Nov.\_B-DATE 29\_I-DATE ..O

We define three new tags:

- B-NP: beginning of a noun phrase chunk
- I-NP: inside of a noun phrase chunk
- O: outside of a noun phrase chunk

*Task:assign POS tags to words*

Pierre \_NNP Vinkem\_NNP,\_ 61\_CD years\_NNS old\_JJ ,\_ will\_MD join\_VB IBM\_NNP  
's\_POS board\_NN as\_IN a\_DT nonexecutive\_JJ director\_NN Nov.\_NNP 29\_CD .\_.

*Task:identify all non-recursive NP chunks*

[NP Pierre Vinken] , [NP 61 years] old , will join [NP IBM] 's [NP board] as [NP a nonexecutive director] [NP Nov. 2] .

*The BIO encoding*

Pierre\_B-NP Vinken\_I-NP , \_O 61\_B-NP years\_I-NP old\_O, \_O join\_O IBM\_B-NP  
's\_O board\_B-NP as\_O a\_B-NP nonexecutive\_I-NP director\_I-NP Nov.\_B-NP 29\_I-NP  
.O

↓

[NP Pierre Vinken] , [NP 61 years] old , will join [NP IBM] 's [NP board] as [NP a nonexecutive director] [NP Nov. 2] .

*Identify all non-recursive NP, verb ('VP') and preposition ('PP') chunks*

NP Pierre Vinken , NP 61 years old , will join NP IBM 's NP board as NP a nonexecutive director NP Nov. 2 .

↓

[NP Pierre Vinken] , [NP 61 years] old , will join [NP IBM] 's [NP board] as [NP a nonexecutive director] [NP Nov. 2] .

**The BIO encoding for shallow parsing** We define several new tags:

- **B-NP, B-VP, B-PP:** beginning of an NP, "VP", "PP" chunk
- **I-NP, I-VP, I-PP:** inside of an NP, "VP", "PP" chunk
- **O:** outside of any chunk

Pierre\_B-NP Vinken\_I-NP , \_O 61\_B-NP years\_I-NP old\_O, \_O join\_O IBM\_B-NP  
's\_O board\_B-NP as\_O a\_B-NP nonexecutive\_I-NP director\_I-NP Nov.\_B-NP 29\_I-NP  
.O

↓

[NP Pierre Vinken] , [NP 61 years] old , will join [NP IBM] 's [NP board] as [NP a nonexecutive director] [NP Nov. 2] .

*Task: identify all mentions of named entities*

NP Pierre Vinken , NP 61 years old , will join NP IBM 's NP board as NP a nonexecutive director NP Nov. 2 .

↓

[PERS Pierre Vinken] , NP 61 years old , will join [ORG IBM] 's NP board as NP a nonexecutive director [DATE Nov. 2] .

**The BIO encoding for NER** We define many new tags:

- **B-PERS, B-DATE, ...:** beginning of a mention of person/date...
- **I-PERS, I-DATE, ...:** inside of a mention of a person/date...
- **O:** outside of any mention of a named entity

Pierre\_B-PERS Vinken\_I-PERS,\_O 61\_O years\_O old\_O ,\_O will\_O join\_O IBM\_B-ORG 's\_O board\_O as\_O a\_O nonexecutive\_O director\_O Nov. 29\_B-DATE .

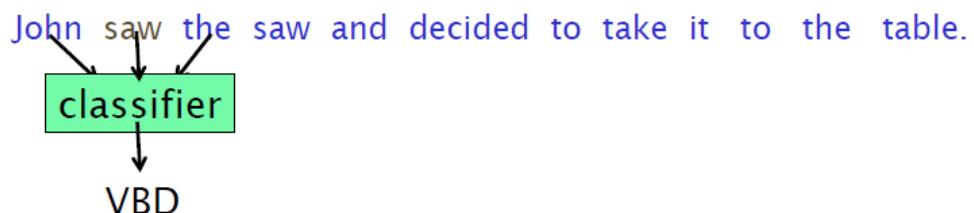
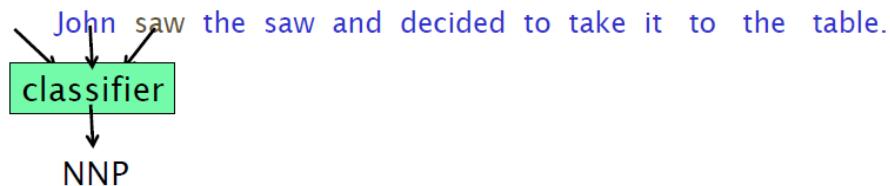
↓

[PERS Pierre Vinken] , NP 61 years old , will join [ORG IBM] 's NP board as NP a nonexecutive director [DATE Nov. 2] .

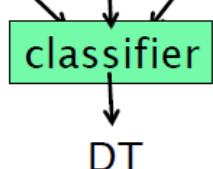
## 5.9 Sequence Labeling as Classification Using Outputs as Inputs

---

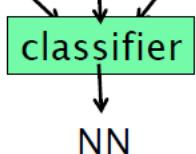
Better input features are usually the categories of the surrounding tokens, but these are not available yet. Can use category of either the preceding or succeeding tokens by going forward or back and using previous output.



John saw the saw and decided to take it to the table.



John saw the saw and decided to take it to the table.



**Word-window classification** Idea is to **classify a word in its context window** of neighboring words. For example, **Named Entity Classification** of a word in context: Person, Location, Organization, None. A simple way to classify a word in context might be to **average** the word vectors in a window and to classify the average vector. But here is a problem: that would **lose position information**. Also we can train a softmax classifier to classify a center word by taking **concatenation of word vectors surrounding it** in a window. Example: Classify "Paris" in the context of this sentence with window length 2:

$$\dots \text{ museums } \text{ in } \text{ Paris } \text{ are } \text{ amazing } \dots$$

$$\bullet\bullet\bullet \quad \bullet\bullet\bullet \quad \bullet\bullet\bullet \quad \bullet\bullet\bullet \quad \bullet\bullet\bullet$$

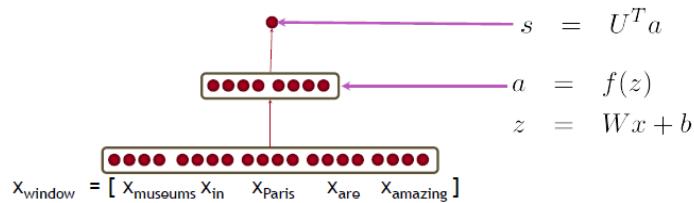
$$x_{\text{window}} = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]^T$$

Resulting vector  $x_{\text{window}} = x \in \mathbb{R}^{5d}$ , a column vector!

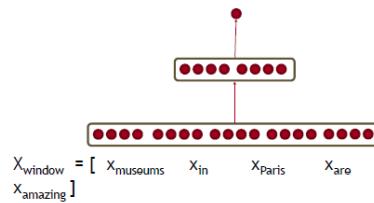
**Simplest window classifier: Softmax** With  $x = x_{\text{window}}$  we can use the same softmax classifier as before:  $\hat{y} = p(y|x) = \frac{\exp(W_y x)}{\sum_{c=1}^C \exp(W_c x)}$ . With cross entropy error as before:  $J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log(\frac{\exp(f_{y_i})}{\sum_{c=1}^C \exp(f_c)})$ . But how do you update the word vector? Short answer: Just take derivatives like last week and optimize.

**Slightly more complex: Multilayer Perceptron** Introduce an additional layer in our softmax classifier with a non-linearity. MLPs are fundamental building blocks of more complex neural systems. Assume we want to classify whether the center word is a Location. Similar to word2vec, we will go over all positions in a corpus. But this time, it will be supervised s.t. positions that are true NER Locations should assign high probability to that class, and others should assign low probability.

**Neural Network Feed-forward Computation** We compute a window  $\text{score}(x) = U^T a \in R$  with a 3-layer neural net:  $s = \text{score}(\text{"museum in Paris are amazing"})$ :  $s = U^T f(Wx + b)$ ,  $x \in$



$R^{20 \times 1}, W \in R^{8 \times 20}, U \in R^{8 \times 2}$ . The middle layer learns **non-linear interactions** between the input word vectors.



**Example:** only if "museums" is first vector should it matter that "in" is in the second position.

**Weakness of word classification** A weakness of the word classification method for a sequence labeling is that it is not capable to make use of the dependencies between POS tags in prediction,

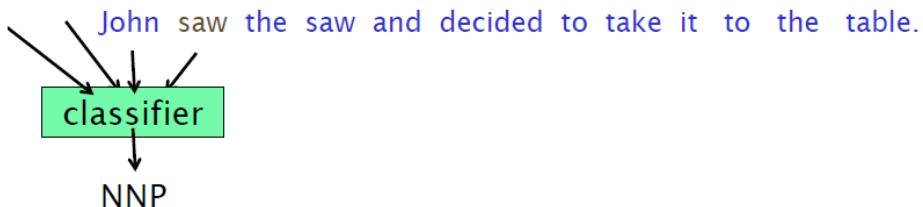
An Extreme example:

It is true for all that that that that that that that refers to is not the same that that that refers to.

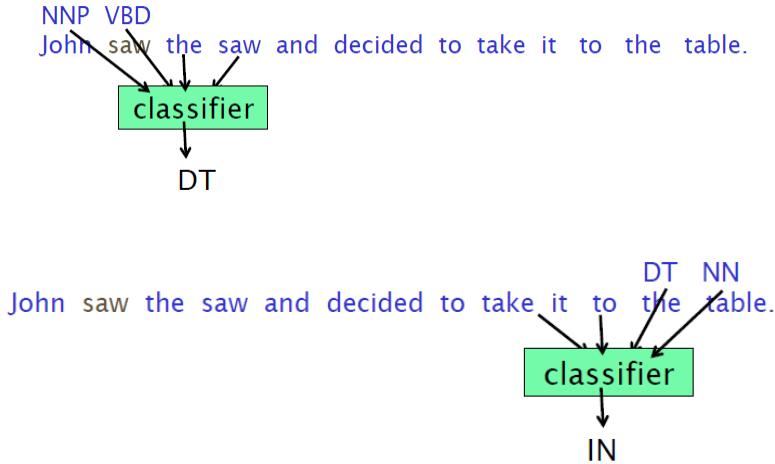
It	is	true	for	all	that	that	that	that	that	that	that	refers	to	is	not	the	same	that	that	that	that	refers	to	.
					pron.	conj.	det.	noun	rel.pron.	det.	noun						noun	rel.pron.	det.	noun				

The word windows for some of the occurrences of "that" in this sentence are the same, if the window size is not greater than 2. The word window classification method will not be able to distinguish those "that". The POS tags of the previous words may be helpful!

Better inputs are usually the categories of the surrounding tokens, but these are not available yet. Also we can use category of either the preceding or succeeding tokens by going forward or back and using previous output.



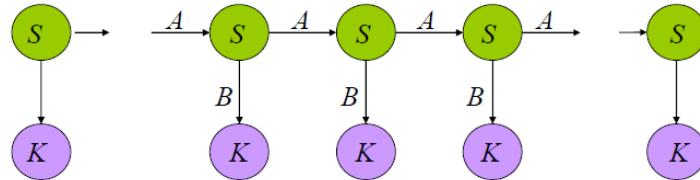
Disambiguating "to" in this case would be even easier backward:



## 5.10 Hidden Markov Models (HMMs)

---

What is HMM? As wiki says: Hidden Markov Model (HMM) is a statistical Markov model in which the system being modeled is assumed to be a Markov process - call it  $X$  - with unobservable ("hidden") states. HMM assumes that there is another process  $Y$  whose behavior "depends" on  $X$ . The goal is to learn about  $X$  by observing  $Y$ . Also it can be described graphically:



Circles here indicate states, arrow indicates probabilistic dependencies between states. Green circles are hidden states, that are dependent only on the previous state. Purple nodes are observer states, they depend only on corresponding hidden state. In formal way:

- $\{S, K, \Pi, A, B\}$
- $S$  are the values for hidden states
- $K$  are the values for the observations
- $\Pi = \pi_i$  are the initial state probabilities
- $A = a_{ij}$  are the state transition probabilities
- $B = b_{ik}$  are the observation state probabilities

Some other diagram examples:

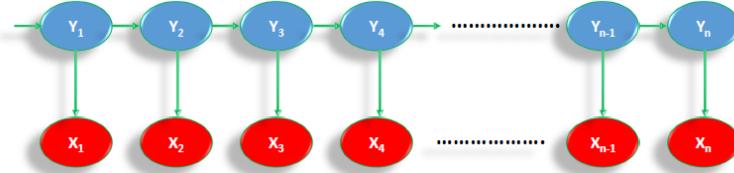


Figure 5.7: \*  
Trellis diagram for sequence of states

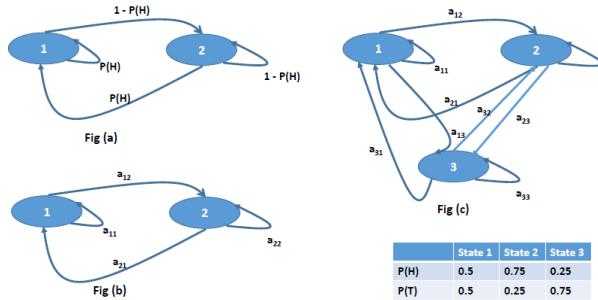


Figure 5.8: \*  
HMMs diagram for coin tossing

## 5.11 Inference in HMM

So, we have seen how we can graphically describe HMMs schemes. But how can we calculate the probability of a given observation sequence? How can we calculate most likely hidden state sequence with given observation sequence? How can we choose a model most closely fitting the data with an observation sequence given?

**Decoding** With a given observation sequence and a model we can compute the probability of the observation sequence:

**Let**  $\mu = (A, B, \Pi); O = (o_1 \dots o_T)$ .

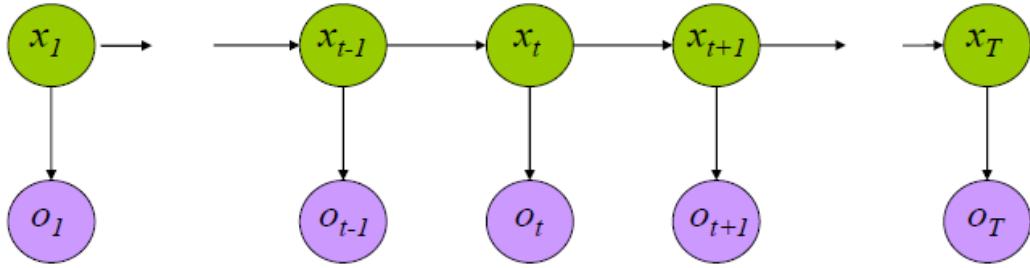
$$P(O|X, \mu) = b_{x_1 o_1} b_{x_2 o_2} \dots b_{x_T o_T} \quad (5.1)$$

$$P(X|\mu) = p_{x_1} a_{x_1 x_2} a_{x_2 x_3} \dots a_{x_{T-1} x_T} \quad (5.2)$$

$$P(O, X|\mu) = P(O|X, \mu)P(X|\mu) \quad (5.3)$$

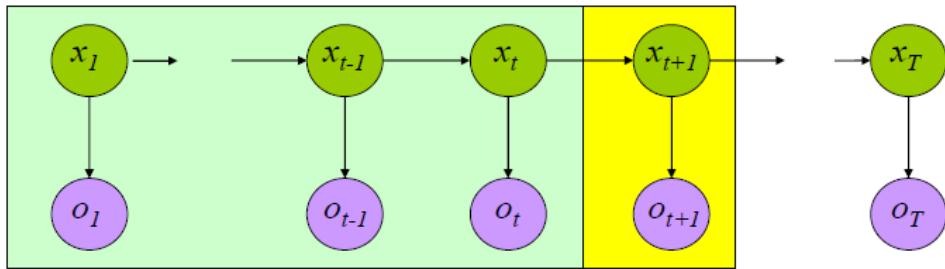
$$P(O|\mu) = \sum_X P(O|X, \mu)P(X|\mu) \quad (5.4)$$

$$P(O|\mu) = \sum_{\{x_1 \dots x_T\}} p_{x_1} b_{x_1 o_1} \prod_{t=1}^{T-1} a_{x_t x_{t+1}} b_{x_{t+1} o_{t+1}} \quad (5.5)$$



## 5.12 Forward procedure for HMM

---



HMM has a special structure that gives us an effective solution using dynamic programming. Intuitively probability of the first  $t$  observations is the same for all possible  $t+1$  length state sequences. Define:  $a_i(t) = P(o_1 \dots o_t, x_t = i | \mu)$ .

Then:

$$\begin{aligned}
 a_j(t+1) &= \\
 &= P(o_1 \dots o_{t+1}, x_{t+1} = j) \\
 &= P(o_1 \dots o_{t+1} | x_{t+1} = j) P(x_{t+1} = j) \\
 &= P(o_1 \dots o_t | x_{t+1} = j) P(o_{t+1} | x_{t+1} = j) P(x_{t+1} = j) \\
 &= P(o_1 \dots o_t, x_{t+1} = j) P(o_{t+1} | x_{t+1} = j) \\
 &= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i, x_{t+1} = j) P(o_{t+1} | x_{t+1} = j) \\
 &= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_{t+1} = j | x_t = i) P(x_t = i) P(o_{t+1} | x_{t+1} = j) \\
 &= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i) P(x_{t+1} = j | x_t = i) P(o_{t+1} | x_{t+1} = j) \\
 &= \sum_{i=1 \dots N} a_i(t) a_{ij} b_{j o_{t+1}}
 \end{aligned} \tag{5.6}$$

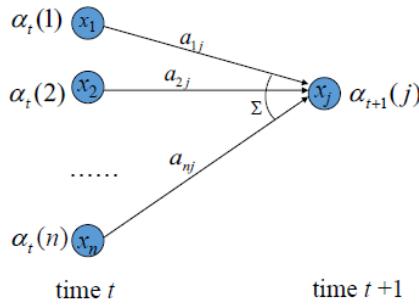


Figure 5.9: \*

Graphically it can be described this way

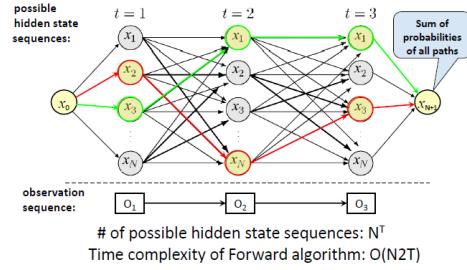


Figure 5.10: \*

And more general

### Forward algorithm Summary

Equation

$$\alpha_i = \pi_i b_i(O_1), 1 \leq i \leq N$$

Multiplications

$$= N$$

for  $t = 1, 2, \dots, T - 1, 1 \leq j \leq N$

$$\alpha_{t+1}(j) = [\sum_{i=1}^N \alpha_t(i) * a_{ij}] * b_j(O_{t+1})$$

$$= (N+1)N(T-1)$$

Finally we have:

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$$

$$= 0$$

Total:  $N + (N+1)N(T-1)$

### 5.13 Backward Procedure

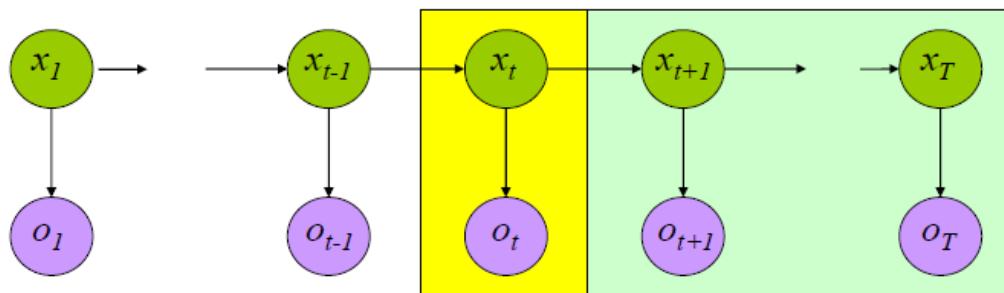


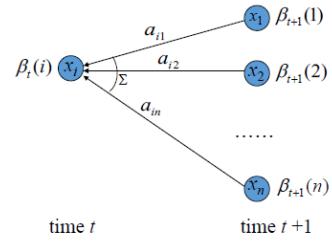
Figure 5.11: \*

Now let's take a look back.

**Here:**  $b_i(t)$  is probability of the rest of the states given the first state.

$$b_i(T+1) = 1 \quad (5.7)$$

$$b_i(t) = P(o_t \dots o_T | x_t = i) b_i(t) = \sum_{j=1 \dots N} a_{ij} b_{io_t} b_j(t+1) \quad (5.8)$$



### Backward algorithm Summary

Equation   Multiplications

$$\beta_T(i) = 1, 1 \leq i \leq N$$

$$= 0$$

for  $t = T - 1, T - 2, \dots, 1 \leq j \leq N$

$$\beta_t(i) = \sum_{i=1}^N a_{ij} * b_j(O_{t+1}) * \beta_{t+1}(j)$$

$$= (2N)N(T-1)$$

Finally we have:

$$P(O|\lambda) = \sum_{i=1}^N \pi_i b_i(O_1) \beta_1(i)$$

$$= 2N$$

$$\text{Total: } 2N + (2N)N(T-1)$$

## 5.14 Decoding Solution summary

**Forward procedure:**  $P(O|\mu) = \sum_{i=1}^N a_i(T)$

**Backward procedure:**  $P(O|\mu) = \sum_{i=1}^N \pi_i \beta_i(1)$

**Combination:**  $P(O|\mu) = \sum_{i=1}^N a_i(t) \beta_i(t)$

## 5.15 Viterbi algorithm

Define  $\delta$  as the state sequence which maximizes the probability of seeing the observations to time  $t-1$ , landing in state  $j$ , and seeing the observation at time  $t$ :

$$\delta_j(t) = \max_{x_1 \dots x_{t-1}} P(x_1 \dots x_{t-1}, o_1 \dots o_{t-1}, x_t = j, o_t) \quad (5.9)$$

$$\delta_j(t+1) = \max_i \delta_i(t) a_{ij} b_{j o_{t+1}} \quad (5.10)$$

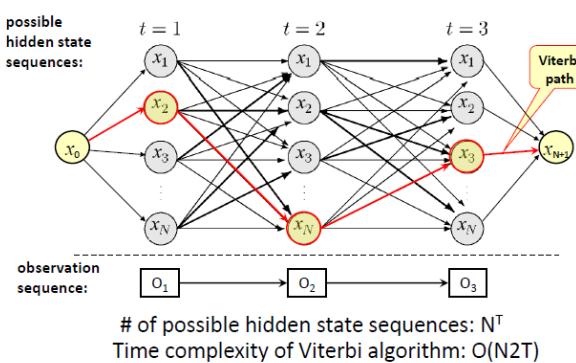
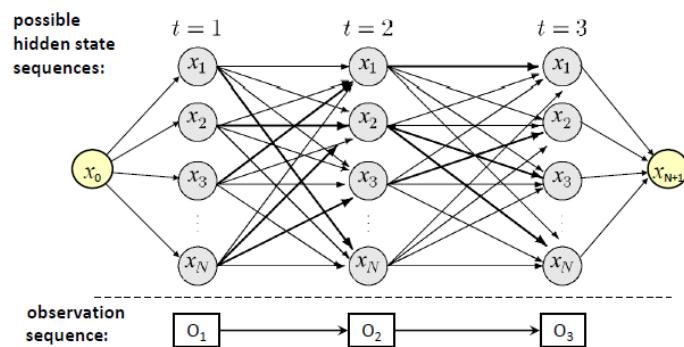
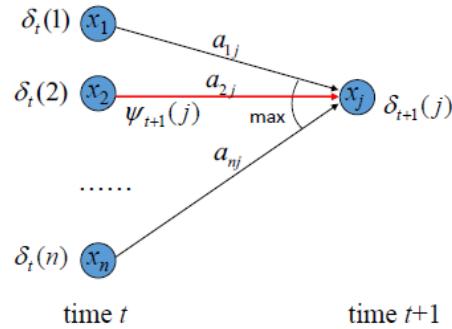
$$\psi_j(t+1) = \arg \max_i \delta_i(t) a_{ij} b_{j o_{t+1}} \quad (5.11)$$

**Computing the most likely state sequence by working backwards:**

$$\hat{X}_T = \arg \max_i \delta_i^T \quad (5.12)$$

$$\hat{X}_t = \psi_{\hat{X}_{t+1}}(t+1) \quad (5.13)$$

$$P(\hat{X}) = \arg \max_i \delta_i^T \quad (5.14)$$



**Viterbi algorithm** Finding The **best single** sequence means computing  $\arg \max_Q P(Q|O, \lambda)$ , equivalent to  $\arg \max QP(Q, O|\lambda)$ . The **Viterbi algorithm** (dynamic programming) defines  $\delta_j(t)$ , i.e. the highest probability of a single path of length  $t$  which accounts for the observations and ends in state  $S_j$ :

$$\delta_j(t) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1, q_2 \dots q_t = j, O_1, O_2 \dots O_t | \lambda) \quad (5.15)$$

By induction:

$$\delta_j(1) = \pi_j b_j o_1 \quad , 1 \leq j \leq N \quad (5.16)$$

$$\delta_j(t+1) = (\max_i \delta_i(t) a_{ij}) b_{j O_{t+1}}, \quad 1 \leq t \leq T-1 \quad (5.17)$$

With **backtracking** (keeping the maximizing argument for each  $t$  and  $j$ ) we find the optimal solution.

## 5.16 Parameter Estimation

---

Given an observation sequence, find the model that is most likely to produce that sequence. No analytic method exists for this task. Given a model and observation sequence, update the model parameters to better fit the observations.

### Learning HMM: two scenarios

**Supervised learning:** estimation when the "right answer" is known. Examples of given:

- a genomic region  $x = x_1 \dots x_{1,000,000}$  where we have good (experimental) annotations of the CpG islands.
- the casino player allows us to observe him one evening, as he changes dice and produces 10,000 rolls

E.g. recall that for complete observation tabular BN:

$$\theta_{ijk}^{ML} = \frac{n_{ijk}}{\sum_{i,j',k} n_{ij'k}} \quad (5.18)$$

$$a_{ij}^{ML} = \frac{\#(i \rightarrow j)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=2}^T y_{n,t-1}^i y_{n,t}^j}{\sum_n \sum_{t=2}^T y_{n,t-1}^i} \quad (5.19)$$

$$b_{ik}^{ML} = \frac{\#(i \rightarrow k)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=1}^T y_{n,t}^i y_{n,t}^k}{\sum_n \sum_{t=1}^T y_{n,t}^i} \quad (5.20)$$

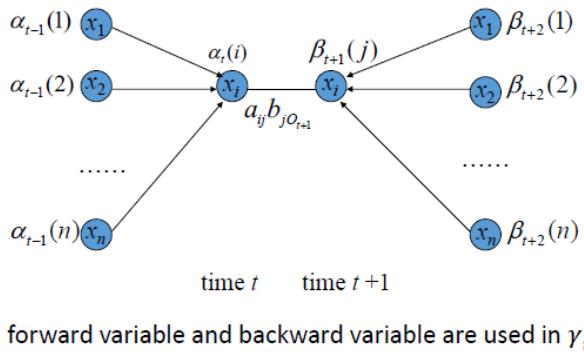
**Unsupervised learning:** when the true state path is unknown, we can fill in the missing values using inference recursions: the Baum Welch algorithm (i.e., EM). It is guaranteed to increase the log likelihood of the model after each iteration. And it converges to local optimum, depending on initial conditions.

**Probability of traversing an arc:**

$$p_t(i, j) = \frac{\alpha_i(t) \alpha_{ij} b_{j o_{t+1}} \beta_j(t+1)}{\sum_{m=1 \dots N} \alpha_m(t) \beta_m(t)} \quad (5.21)$$

**Probability of being in state i:**

$$\gamma_i(t) = \sum_{j=1 \dots N} p_t(i, j) \quad (5.22)$$



**Now we can compute the new estimates of the parameters:**

$$\hat{\pi}_i = \gamma_i(1) \quad (5.23)$$

$$\hat{a}_{ij} = \frac{\sum_{t=1}^T p_t(i, j)}{\sum_{t=1}^T \gamma_t(i)} \quad (5.24)$$

$$\hat{b}_{ik} = \frac{\sum_{\{t: o_t=k\}} \gamma_t(i)}{\sum_{t=1}^T \gamma_t(i)} \quad (5.25)$$

## 5.17 HMM Applications

---

- Generating parameters for n-gram models
- Tagging speech
- Speech recognition

## 5.18 Directed graphical models

---

Graphical models are a *notation for probability models*. In a **directed** graphical model, each **node** represents a distribution over a random variable:  $P(X) = X$

**Arrows** represent dependencies (they define what other random variables the current node is conditioned on)

$$P(Y)P(X|Y) = Y \rightarrow X$$

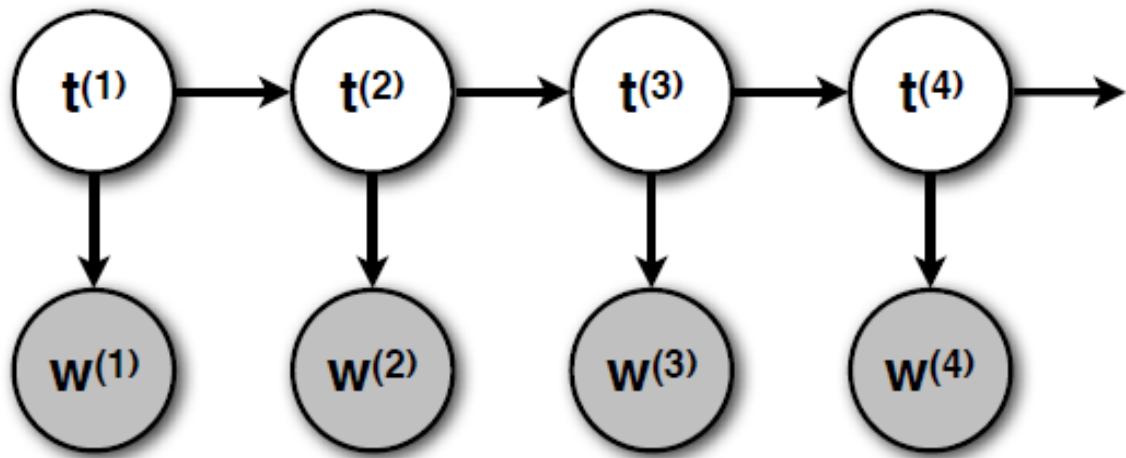
$$P(Y)P(Z)P(X|Y,Z) = \begin{array}{c} Y \longrightarrow X \\ \uparrow \\ Z \end{array}$$

**Shaded nodes** represent observed variables. **White nodes** represent hidden variables  
 $P(Y)P(X|Y)$  with Y hidden and X observed =  $Y \rightarrow X$

## 5.19 HMMs as graphical models

---

HMMs are generative models of the observed input string  $w$ . They 'generate'  $w$  with  $P(w, t) = \prod_i P(t_i | t_{i-1})P(w_i | t_i)$ . When we use an HMM to tag, we observe  $w$ , and need to find  $t$ .



## 5.20 Models for sequence labeling

---

**Sequence labeling:** Given an input sequence  $w = w^{(1)} \dots w^{(n)}$ , predict the best (most likely) label sequence  $t = t^{(1)} \dots t^{(n)}$ :  $\arg \max_t P(t|w)$ .

**Generative models:** use Bayes Rule:

$$\begin{aligned}
 \arg \max_t P(t|w) &= \arg \max_t \frac{P(t, w)}{P(w)} \\
 &= \arg \max_t P(t, w) \\
 &= \arg \max_t P(t)P(w|t)
 \end{aligned} \tag{5.26}$$

**Discriminative (conditional) models:** model  $P(t|w)$  directly.

## 5.21 Advantages of discriminative models

---

We are usually not really interested in  $P(w|t)$ . If  $w$  is given we don't need to predict it. Why not model what we are actually interested in:  $P(t|w)$ .

**Modelling  $P(w|t)$  well is quite difficult:**

- Prefixes (capital letters) or suffixes are good predictors for certain classes of  $t$  (proper nouns, adverbs, ...)
- So we don't want to model words as atomic symbols, but in terms of features
- These features may also help us deal with unknown words
- But features may not be independent

**Modeling  $P(t|w)$  with features should be easier** : Now we can incorporate arbitrary features of the word, because we don't need to predict  $w$  anymore.

## 5.22 Discriminative probability models

---

A discriminative or **conditional** model of the labels  $t$  given the observed input string  $w$  models  $P(t|w) = \prod_i P(t_i|w_i, t_{i-1})$  directly.

There are two main types of discriminative probability models: Maximum Entropy Markov Models (MEMMs) and Conditional Random Fields (CRFs). They are both based on logistic regression, have the same graphical model, require the Viterbi algorithm for tagging. They differ in that MEMMs consist of independently learned distributions, while CRFs are trained to maximize the probability of the entire sequence.

## 5.23 Probabilistic classification

---

**Classification:** Predict a class (label)  $c$  for input  $x$ . There are only a (small) finite number of possible class labels.