# DD1418
## 5a: Part-of-speech tagging

Johan Boye, KTH

# Recall: Word classes

Words can be divided into classes depending on their use in the language.

- Noun, verb, adjective, adverb, preposition, pronoun, conjunction, interjection, determiner, etc.
- These classes are often called *parts-of-speech* (or *POS tags*).

The *part-of-speech tagging* problem is to assign a POS tag to each word in a text.

Some words can belong to more than one class, e.g. *like*:

VERB "*I like her.*"

NOUN "*He got a like on Facebook.*"

ADJ "*The portrait is very like.*"

ADV "*This is, like, crazy!*"

PREP "*It looks like an accident*"

CONJ "*He acted like he was all alone*"

Swedish also has this kind of words, e.g. *var*:

*Den gamle mannen visste inte var han var.*

Swedish also has this kind of words, e.g. *var*:

| Den | gamle | mannen | visste | inte | var | han | var |
|-----|-------|--------|--------|------|-----|-----|-----|
| DET | ADJ | NOUN | VERB | ADV | ADV | PRON | VERB |

## Part-of-speech tagging

Retrieve all possible tags from a dictionary, then decide which
ones are the most likely, e.g. :

| *I* | *like* | *plays* | *about* | *Bolliwogs* |
|:---:|:---:|:---:|:---:|:---:|
| PRON | VB/NOUN/ADJ/ ADV/PREP/CONJ | VB/NOUN | ADV/PREP/ADJ | ? |

$$\Downarrow$$

| | | | | |
|:---:|:---:|:---:|:---:|:---:|
| PRON | VB | NOUN | PREP | NOUN |

Given a sequence of words $w_1 \ldots w_n$, we want to find the most probable sequence of tags $t_1 \ldots t_n$.

$$\arg \max_{t_1 \ldots t_n} P(t_1 \ldots t_n | w_1 \ldots w_n)$$

What information can we use to do POS tagging?

The word itself *man* is often a noun, more rarely an interjection, even more rarely a verb.

The preceding tags Some combinations are common (noun verb), others are rare (determiner determiner).

Bayes' theorem expresses how beliefs should be updated in the light of new evidence.



$$P(S|T) = \frac{P(T|S)P(S)}{P(T)}$$

Likelihood

Prior

Posterior

Marginal

# Rewrite and simplify

Rewrite using Bayes' theorem:

$$\arg \max_{t_1 \ldots t_n} P(t_1 \ldots t_n | w_1 \ldots w_n) =$$

$$= \arg \max_{t_1 \ldots t_n} \frac{P(w_1 \ldots w_n | t_1 \ldots t_n) P(t_1 \ldots t_n)}{P(w_1 \ldots w_n)} =$$

$$= \arg \max_{t_1 \ldots t_n} P(w_1 \ldots w_n | t_1 \ldots t_n) P(t_1 \ldots t_n)$$

## Rewrite and simplify

$$\arg\max_{t_1 \dots t_n} P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n)$$

Markov assumptions: Prior (context model)

- $P(t_1 \dots t_n) = \prod_{i=1}^{n} P(t_i | t_{i-1})$ (in the bigram case)
- $P(t_1 \dots t_n) = \prod_{i=1}^{n} P(t_i | t_{i-k+1} \dots t_{i-1})$ (in the $k$-gram case)

Likelihood (lexical model)

- $P(w_1 \dots w_n | t_1 \dots t_n) = \prod_{i=1}^{n} P(w_i | t_i)$

# Maximum likelihood estimation of parameters

Lexical model:

$$P(w_i|t_i) = \frac{c(w_i, t_i)}{c(t_i)} = \frac{\text{\# of times } w_i \text{ has been tagged as } t_i}{\text{\# of tokens tagged as } t_i}$$

Context model (bigram case):

$$P(t_i|t_{i-1}) = \frac{c(t_{i-1}t_i)}{c(t_{i-1})} = \frac{\text{\# of times a } t_{i-1} \text{ is followed by a } t_i}{\text{\# of tokens tagged as } t_{i-1}}$$

# Decoding problem

Find the most probable tagging (bigram case):

$$\arg\max_{t_1\ldots t_n} \prod_{i=1}^{n} P(t_i|t_{i-1})P(w_i|t_i)$$

*k*-gram case:

$$\arg\max_{t_1\ldots t_n} \prod_{i=1}^{n} P(t_i|t_{i-k+1}\ldots t_{i-1})P(w_i|t_i)$$

# DD1418
## 5b: Hidden Markov Models

Johan Boye, KTH

Given a sequence of words $w_1 \ldots w_n$, we want to find the most probable sequence of POS tags $t_1 \ldots t_n$.

$$\arg \max_{t_1 \ldots t_n} P(t_1 \ldots t_n | w_1 \ldots w_n)$$

The POS tags $t_1 \ldots t_n$ are the **hidden states**, and the words $w_1 \ldots w_n$ are the **observations**.

# Hidden Markov Models

A Hidden Markov Model (HMM) consists of:

- a set of (hidden) states $\{q_1, \ldots, q_m\}$
- the transition probability matrix $A$
- a set of observations $\{o_1, \ldots, o_k\}$.
- the observation probability matrix $B$

# Example: A first-order Markov chain for POS tags



$P(\text{START NOUN VERB DET NOUN}) = 0.3 \times 0.3 \times 0.7 \times 1$

# Example



$P(flies|N) = 10^{-6}$  $P(flower|N) = 10^{-5}$  $P(like|N) = 10^{-7}$
$P(flies|V) = 10^{-5}$  $P(flower|V) = 10^{-8}$  $P(like|Prep) = 10^{-5}$
$P(a|N) = 10^{-8}$  $P(a|DET) = 0.36$  $P(like|V) = 10^{-4}$

Find the POS tags for **flies like a flower**.

## Viterbi algorithm

We want to find the most probable sequence of hidden states.

A naive algorithm that enumerates all the possible sequences and computes their probability would have exponential complexity.

However, the **Viterbi** algorithm finds the solution quicker (quadratic complexity), by storing intermediate results.

This is an example of **dynamic programming** (just as the algorithm in assignment 1).

$\text{v}(\text{t},\text{k}) =$ the probability of being in state $q_k$ after having seen the first $t$ observations, and passing through the most probable state sequence $q_{i_1} \ldots q_{i_{t-1}}$.

**flies**

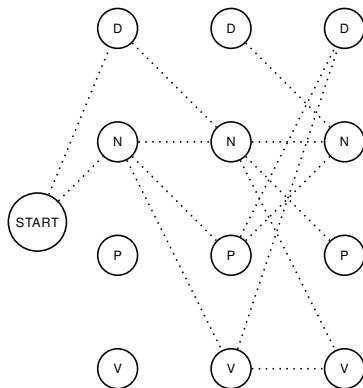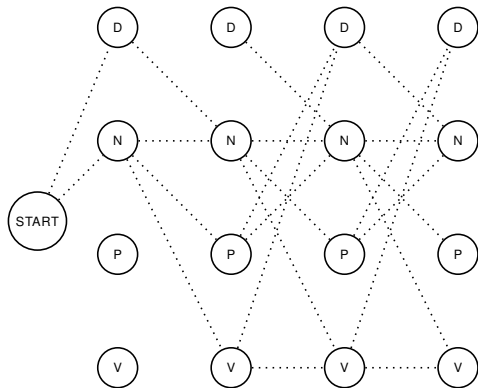**flies**     **like**

# Viterbi trellis



flies    like    a    flower

Bigram case:

$$v(0, k) = p(q_k | \text{START}) p(o_0 | q_k)$$

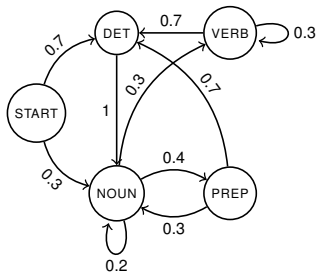$$v(t, k) = \max_i [\, v(t - 1, i) \, p(q_k | q_i) \, p(o_t | q_k) \,]$$

$$\texttt{backptr}(t, k) = \arg \max_i [\, v(t - 1, i) \, p(q_k | q_i) \, p(o_t | q_k) \,]$$

# DD1418:
## 5c: Viterbi decoding

Johan Boye, KTH

$P(\textit{flies}|N) = 10^{-6}$    $P(\textit{flower}|N) = 10^{-5}$    $P(\textit{like}|N) = 10^{-7}$

$P(\textit{flies}|V) = 10^{-5}$    $P(\textit{flower}|V) = 10^{-8}$    $P(\textit{like}|\textit{Prep}) = 10^{-5}$
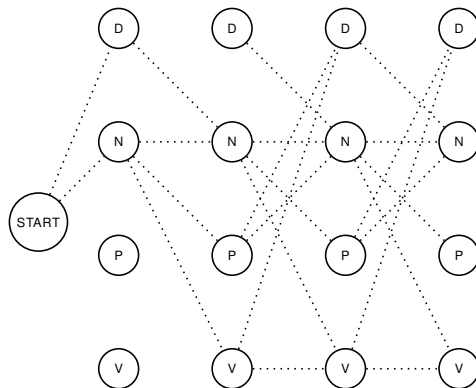
$P(a|N) = 10^{-8}$    $P(a|DET) = 0.36$    $P(\textit{like}|V) = 10^{-4}$

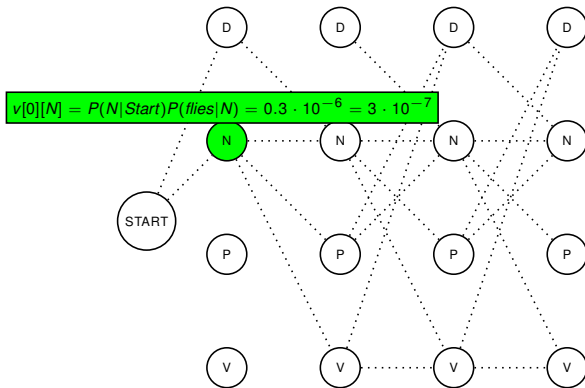Find the POS tags for **flies like a flower**.

# Viterbi trellis

Bigram case:

$$v(0, k) = p(q_k|\text{START})p(o_0|q_k)$$

$$v(t, k) = \max_i [\, v(t-1, i)\, p(q_k|q_i)\, p(o_t|q_k)\,]$$

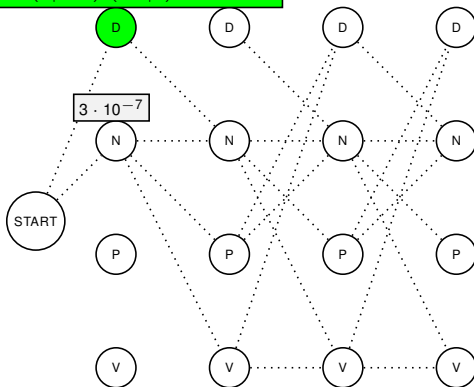$$\text{backptr}(t, k) = \arg\max_i [\, v(t-1, i)\, p(q_k|q_i)\, p(o_t|q_k)\,]$$

**flies**     **like**     **a**     **flower**

$v[0][N] = P(N|Start)P(flies|N) = 0.3 \cdot 10^{-6} = 3 \cdot 10^{-7}$

# Viterbi decoding



**flies    like    a    flower**

$v[0][D] = P(D|Start)P(flies|D) = 0.7 \cdot 0 = 0$

$3 \cdot 10^{-7}$

# Viterbi decoding



**flies    like    a    flower**

$v[1][N] =$
$\max(v[0][N]P(N|N), v[0][D]P(N|D))P(like|N) =$
$\max(3 \cdot 10^{-7} \cdot 0.2, 0) \cdot 10^{-6} =$
$6 \cdot 10^{-14}$

# Viterbi decoding

**flies** **like** **a** **flower**

$$v[1][P] = v[0][N] \cdot P(P|N) \cdot P(\textit{like}|P) = 3 \cdot 10^{-7} \cdot 0.4 \cdot 10^{-5} = 1.2 \cdot 10^{-12}$$
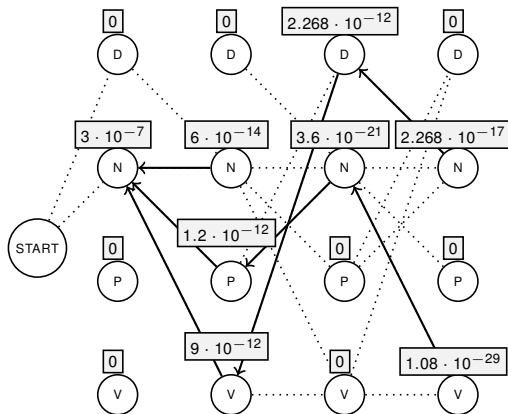
# Viterbi decoding

**flies    like    a    flower**



$v[2][D] =$
$\max(v[1][V]P(D|V), v[1][P]P(D|P))P(a|D) =$
$\max(9 \cdot 10^{-12} \cdot 0.7, 1.2 \cdot 10^{-12} \cdot 0.7) \cdot 0.36 =$
$2.268 \cdot 10^{-12}$

$3 \cdot 10^{-7}$
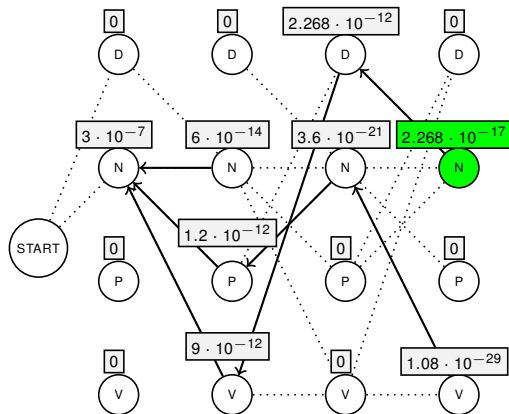$6 \cdot 10^{-14}$
$1.2 \cdot 10^{-12}$
$9 \cdot 10^{-12}$

**flies**     **like**     **a**     **flower**
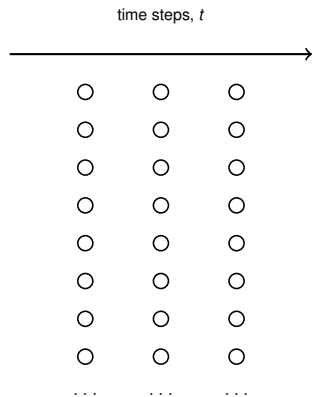
# Viterbi decoding

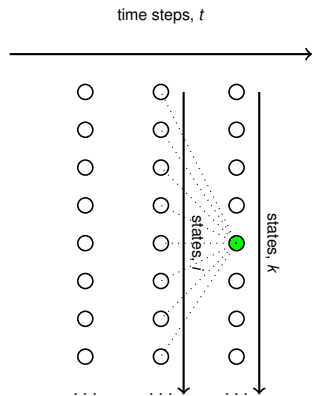flies      like      a      flower

# The Viterbi loop, bigram case



time steps, $t$

time steps, *t*



states, *k*

# The Viterbi loop, bigram case



time steps, $t$

states, $k$

$\mathtt{v[t][k]} =$

# The Viterbi loop, bigram case



time steps, $t$

states, $j$

states, $k$

$v[t][k] =$

# Viterbi decoding, data structures

Data structures in the bigram case:

- $a[i][k]$ = probability of going from $q_i$ to $q_k$

- $b[i][k]$ = probability of state $q_k$ emitting $o_i$

- $v[t][k]$ = probability of being in state $q_k$ after having seen the first $t$ observations, and passing through the most probable state sequence $q_{i_1} \ldots q_{i_{t-1}}$.

# DD1418:
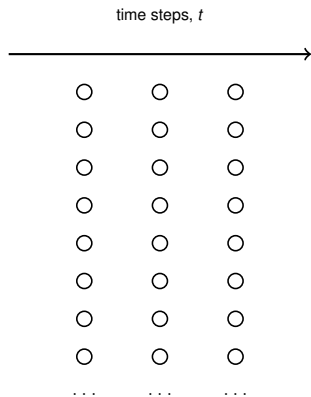## 5d: More on Viterbi decoding

Johan Boye, KTH

# Viterbi algorithm

Trigram case:

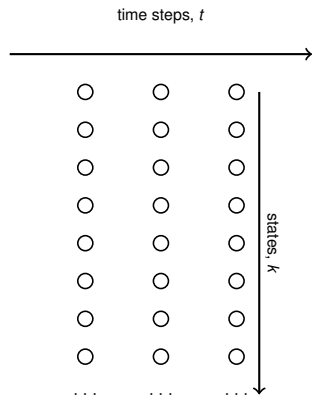$$v(0, \text{START}, k) = p(q_k|\text{START START})p(o_0|q_k)$$

$$v(t, j, k) = \max_i [\, v(t-1, i, j)\, p(q_k|q_i q_j)\, p(o_t|q_k)\,]$$

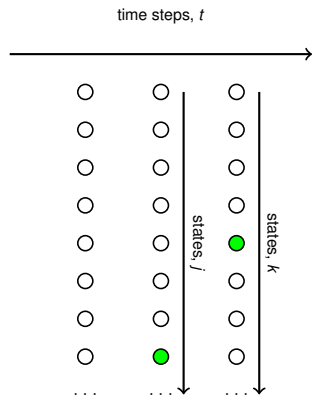$$\texttt{backptr}(t, j, k) = \arg\max_i [\, v(t-1, i, j)\, p(q_k|q_i q_j)\, p(o_t|q_k)\,]$$
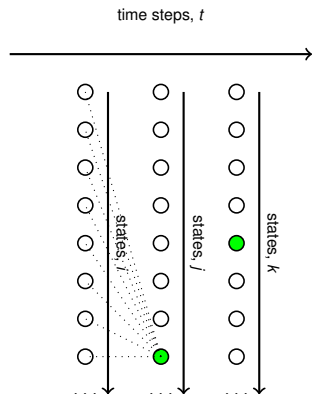
time steps, $t$

# The Viterbi loop, trigram case

# The Viterbi loop, trigram case



time steps, $t$

states, $j$

states, $k$

$\text{v}[\text{t}][\text{j}][\text{k}] =$

# The Viterbi loop, trigram case



time steps, $t$

states, $i$

states, $j$

states, $k$

$\mathrm{v[t][j][k]} =$
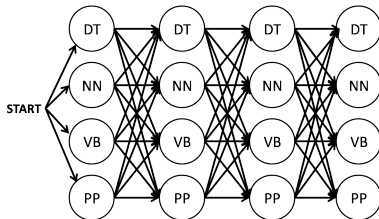
Data structures in the trigram case:

- $a[i][j][k]$ = probability of going to $q_k$ given that the preceding state was $q_j$, and the state before that was $q_i$

- $b[i][k]$ = probability of state $q_k$ emitting $o_i$

- $v[t][j][k]$ = probability of being in state $q_k$ after having seen the first $t$ observations, and passing through the most probable state sequence $q_{i_1} \ldots q_{i_{t-2}} q_j$.

# Smoothing

Normally the context model is **smoothed**, just as in the case of language models.

That is: Even if we have never seen the combination DET DET in our data, we will still give $P(\text{DET}|\text{DET})$ a non-zero probability.

This means that we have to consider all possible state transitions in every step of the Viterbi algorithm.
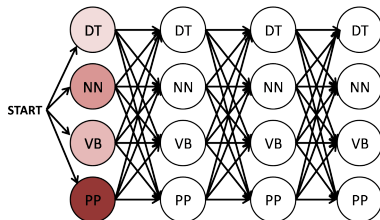
## Beam search

If we consider all possible state transitions, the Viterbi algorithm can still be quite time-consuming.
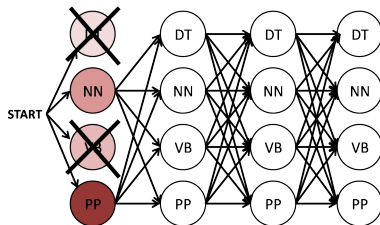
Common solution: **Beam search**

In every time step $t$, only consider the $\beta$ states with the highest *v*-value. Paths ending at other states at time step $t$ are terminated.

The parameter $\beta$ is called the **beam width**.
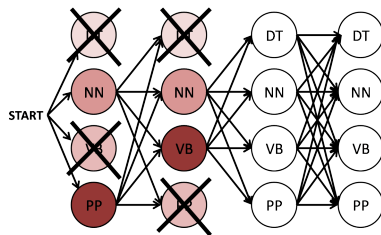
# Beam search (example)



1a. In time step 1, the state PP is best, NN is second best.

1b. Remove the other states in step 1, and their outbound transitions.

# Beam search (example)



2. In time step 2, $\texttt{VB}$ and $\texttt{NN}$ are best. Remove $\texttt{DT}$ and $\texttt{PP}$.

etc.