

- **Part-of-Speech Tagging**

- Input some sentence
- Output have a tagged sequence each word is given an associated tag

INPUT:
 Profits soared at Boeing Co., easily topping forecasts on Wall Street, as their CEO Alan Mulally announced first quarter results.

OUTPUT:
 Profits/N soared/V at/P Boeing/N Co./N ./, easily/ADV topping/V forecasts/N on/P Wall/N Street/N ./, as/P their/POSS CEO/N Alan/N Mulally/N announced/V first/ADJ quarter/N results/N ./.

N = Noun
 V = Verb
 P = Preposition
 Adv = Adverb
 Adj = Adjective
 ...

- Named Entity Recognition
 - Input some sentence
 - Output named entities
 - A company, location, Person
- Named Entity Extraction as Tagging
 - Input some sentence
 - Output tagged sequence
 - NA = No entity
 - SC = Start company
 - CC = Continue Company
 - SL = Start location
 - CL = Continue Location
- Our Goal
 - Treat as machine learning problem
 - Training set
 - A set of sentences
 - Training sentences annotated by hand

Training set:

1 Pierre/NNP Vinken/NNP ./, 61/CD years/NNS old/JJ ./, will/MD join/VB the/DT board/NN as/IN a/DT nonexecutive/JJ director/NN Nov./NNP 29/CD ./.

2 Mr./NNP Vinken/NNP is/VBZ chairman/NN of/IN Elsevier/NNP N.V./NNP ./, the/DT Dutch/NNP publishing/VBG group/NN ./.

3 Rudolph/NNP Agnew/NNP ./, 55/CD years/NNS old/JJ and/CC chairman/NN of/IN Consolidated/NNP Gold/NNP Fields/NNP PLC/NNP ./, was/VBD named/VBN a/DT nonexecutive/JJ director/NN of/IN this/DT British/JJ industrial/JJ conglomerate/NN ./.

...

38,219 It/PRP is/VBZ also/RB pulling/VBG 20/CD people/NNS out/IN of/IN Puerto/NNP Rico/NNP ./, who/WP were/VBD helping/VBG Hurricane/NNP Hugo/NNP victims/NNS ./, and/CC sending/VBG them/PRP to/TO San/NNP Francisco/NNP instead/RB ./.

- From the training set, induce a function / algorithm that maps new sentences to their tag sequences
- Treat this problem as a supervised learning problem
- Two types of Constraints
 - “Local”: e.g., can is more likely to be a model verb MD rather than a noun NN
 - has bias to be one part of speech over another part of speech
 - “Contextual”: e.g., a noun is much more likely than a verb to follow a determiner
 - Sometimes these preferences are in conflict

- The trash can is in the garage
- **Generative models, and the noisy-channel model, for supervised learning**
- Supervised Learning Problems
 - We have training examples $x^{(i)}$, $y^{(i)}$ for $i = 1 \dots m$. Each $x^{(i)}$ is an input, each $y^{(i)}$ is a label
 - Task is to learn a function f mapping inputs x to labels $f(x)$
 - $x^{(i)}$ = the dog laughs, $y^{(i)}$ DT NN VB
 - Conditional models:
 - Learn a distribution $p(y|x)$ from training examples
 - will have various parameters
 - For any test input x , define $f(x) = \arg \max_y p(y|x)$
- Generative Models
 - We have training examples $x^{(i)}$, $y^{(i)}$ for $i = 1 \dots m$. Task is to learn a function f mapping inputs x to labels $f(x)$
 - Generative models:
 - Learn a distribution $p(x,y)$ from training examples
 - Joint distribution
 - Often we have $p(x,y) = p(y)p(x|y)$
 - Use bayes rule
 - $p(y)$ prior likelihood
 - $p(x|y)$ conditional generative model
 - the probability of generating x given y

► Note: we then have

$$p(y|x) = \frac{p(y)p(x|y)}{p(x)}$$

where $p(x) = \sum_y p(y)p(x|y)$

• $p(y|x)$ DIRECTLY
"DISCRIMINATIVE"

• $p(x,y)$ GENERATIVE

- given a joint distribution can derive a conditional distribution
- Decoding with Generative Models
 - We have training examples $x^{(i)}$, $y^{(i)}$ for $i = 1 \dots m$. Task is to learn a function f mapping inputs x to labels $f(x)$
 - Generative models:
 - Learn a distribution $p(x,y)$ from training examples
 - Often we have $p(x,y) = p(y)p(x|y)$

► Output from the model:

$$\begin{aligned} f(x) &= \arg \max_y p(y|x) \\ &= \arg \max_y \frac{p(y)p(x|y)}{p(x)} \\ &= \arg \max_y p(y)p(x|y) \end{aligned}$$

- $p(x)$ does not vary with y
- **Hidden Markov Model (HMM) taggers**
- Hidden Markov Models
 - have sentence $x = x_1 x_2 x_3 \dots x_n$

- We have a tag sequence $y = y_1 y_2 y_3 \dots y_n$
- We'll use an HMM to define

- ▶ We'll use an HMM to define

$$p(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n)$$

for any sentence $x_1 \dots x_n$ and tag sequence $y_1 \dots y_n$ of the same length.

- ▶ Then the most likely tag sequence for x is

$$\arg \max_{y_1 \dots y_n} p(x_1 \dots x_n, y_1, y_2, \dots, y_n)$$

- $p()$ defines a joint distribution over word sequences and tag sequences
- the output is the tag sequence that maximizes the probability
 - the number of tag sequences grows exponential with n
 - therefore brute force is not going to work
- Trigram Hidden Markov Models (Trigram HMMs)

For any sentence $x_1 \dots x_n$ where $x_i \in \mathcal{V}$ for $i = 1 \dots n$, and any tag sequence $y_1 \dots y_{n+1}$ where $y_i \in \mathcal{S}$ for $i = 1 \dots n$, and $y_{n+1} = \text{STOP}$, the joint probability of the sentence and tag sequence is

$$p(x_1 \dots x_n, y_1 \dots y_{n+1}) = \prod_{i=1}^{n+1} q(y_i | y_{i-2}, y_{i-1}) \prod_{i=1}^n e(x_i | y_i)$$

where we have assumed that $x_0 = x_{-1} = *$.

Handwritten notes: $\mathcal{V} = \{the, dog, cat, a, bark, \dots\}$, $\mathcal{S} = \{DT, NN, VB, P, ADV, \dots\}$, $e(the | DT) \approx 50 \text{ TAGS}$

Parameters of the model:

- ▶ $q(s|u, v)$ for any $s \in \mathcal{S} \cup \{\text{STOP}\}$, $u, v \in (\mathcal{S} \cup \{*\}) \leftarrow \text{TRIGRAM}$
- ▶ $e(x|s)$ for any $s \in \mathcal{S}$, $x \in \mathcal{V} \leftarrow \text{EMISSION PARAMETERS}$

- Parameters
 - Trigram parameters
 - Conditional probability \leftarrow emission parameters $e(x|s)$
- An example

If we have $n = 3$, $x_1 \dots x_3$ equal to the sentence the dog laughs, and $y_1 \dots y_4$ equal to the tag sequence DT NN V STOP, then

$$\begin{aligned} p(x_1 \dots x_n, y_1 \dots y_{n+1}) &= q(D|*, *) \times q(N|*, D) \times q(V|D, N) \times q(\text{STOP}|N, V) \\ &\quad \times e(\text{the}|D) \times e(\text{dog}|N) \times e(\text{laughs}|V) \end{aligned}$$

- ▶ STOP is a special tag that terminates the sequence
- ▶ We take $y_0 = y_{-1} = *$, where $*$ is a special "padding" symbol

- Why the Name?
 - Hidden Markov Models

HIDDEN MARKOV MODELS

$$p(x_1 \dots x_n, y_1 \dots y_n) = \underbrace{q(\text{STOP} | y_{n-1}, y_n) \prod_{j=1}^n q(y_j | y_{j-2}, y_{j-1})}_{\text{Markov Chain}} \times \underbrace{\prod_{j=1}^n e(x_j | y_j)}_{\substack{\text{ML ESTIMATE} \\ x_j \text{'s are observed}}}$$

$p(x_1, \dots, x_n | y_1, \dots, y_{n+1})$ $p(x, y) = p(y) \times p(x|y)$

- Prior probability over tag sequences
- Second order Markov Chain
 - applied Markov to the problem of $p(y)$
- $p(x_1 \dots x_n | y_1 \dots y_{n+1})$
 - probability of x conditioned on y
 - each word x is chosen only on the value of y
- choose sequence of tags under this model
- x_j 's are observed and y_j 's are unobserved
- Hidden Markov Model (HMM) taggers Parameter estimation**
- Smoothed Estimation

$$q(Vt | Dt, Jj) = \lambda_1 \times \frac{\text{Count}(Dt, Jj, Vt)}{\text{Count}(Dt, Jj)} + \lambda_2 \times \frac{\text{Count}(Jj, Vt)}{\text{Count}(Jj)} + \lambda_3 \times \frac{\text{Count}(Vt)}{\text{Count}()}$$

ML ESTIMATE
BIGRAM ML ESTIMATE
UNIGRAM ML ESTIMATE

$\lambda_1 + \lambda_2 + \lambda_3 = 1$, and for all i , $\lambda_i \geq 0$

$e(x|y) = 0$ for all y , if x is never seen in the training data

$e(\text{base} | Vt) = \frac{\text{Count}(Vt, \text{base})}{\text{Count}(Vt)}$

- Dealing with Low-Frequency Words: An Example
 - A common method is as follows:
 - Step1: Split vocabulary into two sets
 - Frequent words = words occurring ≥ 5 times in training
 - Low Frequency words = all other words
 - Step2: Map low frequency words into small, finite set, depending on prefixes, suffixes etc.

Dealing with Low-Frequency Words: An Example

[Bikel et. al 1999] (named-entity recognition)

Word class	Example	Intuition
-twoDigitNum	90	Two digit year
-fourDigitNum	1990	Four digit year
-containsDigitAndAlpha	A8956-67	Product code
-containsDigitAndDash	09-96	Date
-containsDigitAndSlash	11/9/89	Date
-containsDigitAndComma	23,000.00	Monetary amount
-containsDigitAndPeriod	1.00	Monetary amount, percentage
-othernum	456789	Other number
-allCaps	BBN	Organization
-capPeriod	M.	Person name initial
-firstWord	first word of sentence	no useful capitalization information
-initCap	Sally	Capitalized word
-lowercase	can	Uncapitalized word
-other	,	Punctuation marks, all other words

- lower frequency words mapped to pseudo words
 - first word
 - initCap
 - etc...

Dealing with Low-Frequency Words: An Example

Profits/NA soared/NA at/NA Boeing/SC Co./CC./NA easily/NA
 topping/NA forecasts/NA on/NA Wall/SL Street/CL./NA as/NA their/NA
 CEO/NA Alan/SP Mulally/CP announced/NA first/NA quarter/NA
 results/NA ./NA

↓

firstword/NA soared/NA at/NA initCap/SC Co./CC./NA easily/NA
 lowercase/NA forecasts/NA on/NA initCap/SL Street/CL./NA as/NA
 their/NA CEO/NA Alan/SP initCap/CP announced/NA first/NA
 quarter/NA results/NA ./NA

NA = No entity
 SC = Start Company
 CC = Continue Company
 SL = Start Location
 CL = Continue Location

- build a Hidden Markov Model
 - $e(\text{firstword|na})$
 - $e(\text{initCap|SC})$
- closing the vocabulary
- mapping the low frequency words to a smaller method maintaining the spelling
- downside is it is heuristic
- **Hidden Markov Model The Viterbi algorithm**
- The Viterbi Algorithm

Problem: for an input $x_1 \dots x_n$, find

$$\arg \max_{y_1 \dots y_{n+1}} p(x_1 \dots x_n, y_1 \dots y_{n+1})$$
 where the arg max is taken over all sequences $y_1 \dots y_{n+1}$ such that $y_i \in S$ for $i = 1 \dots n$, and $y_{n+1} = \text{STOP}$.
 We assume that p again takes the form

$$p(x_1 \dots x_n, y_1 \dots y_{n+1}) = \prod_{i=1}^{n+1} q(y_i | y_{i-2}, y_{i-1}) \prod_{i=1}^n e(x_i | y_i)$$

Recall that we have assumed in this definition that $y_0 = y_{-1} = *$, and $y_{n+1} = \text{STOP}$.

- Brute Force Search is Hopelessly Inefficient

Problem: for an input $x_1 \dots x_n$, find

$$\arg \max_{y_1 \dots y_{n+1}} p(x_1 \dots x_n, y_1 \dots y_{n+1})$$
 where the arg max is taken over all sequences $y_1 \dots y_{n+1}$ such that $y_i \in S$ for $i = 1 \dots n$, and $y_{n+1} = \text{STOP}$.

the dog laughs \rightarrow

D	D	D	STOP $\rightarrow 0.3$
D	D	N	STOP $\rightarrow 0.01$
D	D	V	STOP $\rightarrow 0.0001$
D	N	D	STOP \rightarrow

GENERAL CASE:
 $|S|^n$ POSSIBLE SEQUENCES

- The Viterbi Algorithm
 - Define n to be the length of the input sentence

The Viterbi Algorithm

- Define n to be the length of the sentence
- Define S_k for $k = -1 \dots n$ to be the set of possible tags at position k :

$$S_{-1} = S_0 = \{*\}$$

$$S_k = S \text{ for } k \in \{1 \dots n\}$$

- Define

$$r(y_{-1}, y_0, y_1, \dots, y_k) = \prod_{i=1}^k q(y_i | y_{i-2}, y_{i-1}) \prod_{i=1}^k e(x_i | y_i)$$

$x_1 \ x_2 \ x_3 \dots x_n$
 $* \ * \ D \ N \ V$
 $y_{-1} \ y_0 \ y_1 \ y_2 \ y_3$
 $k=3$

- r takes a sequence of tags as input
 - calculates its probability under the HMM under a truncated expression
- Define a dynamic programming table

$\pi(k, u, v) =$ maximum probability of a tag sequence ending in tags u, v at position k

that is,

$$\pi(k, u, v) = \max_{(y_{-1}, y_0, y_1, \dots, y_k) : y_{k-1}=u, y_k=v} r(y_{-1}, y_0, y_1 \dots y_k)$$

- An Example

$\pi(k, u, v) =$ maximum probability of a tag sequence ending in tags u, v at position k

$S = \{D, N, V, P\}$

$* \ * \ D \ N \ V \ P \ D$

The man saw the dog with the telescope

$\pi(7, P, D)$

$q(\cdot)$
 $e(\cdot)$

- $\pi()$ will be the maximum probability for any sequence that ends at position 6 and 7
- Each sequence will have a probability by multiplying the q terms and e terms
- A Recursive Definition

Base case:

$$\pi(0, *, *) = 1$$

Recursive definition:

For any $k \in \{1 \dots n\}$, for any $u \in S_{k-1}$ and $v \in S_k$:

$$\pi(k, u, v) = \max_{w \in S_{k-2}} (\pi(k-1, w, u) \times q(v|w, u) \times e(x_k|v))$$

$\{1, \dots, n\}$
 $\in S_{k-1}$
 $\in S_k$

- Justification of the Recursive Definition

For any $k \in \{1 \dots n\}$, for any $u \in S_{k-1}$ and $v \in S_k$:

$$\pi(k, u, v) = \max_{w \in S_{k-2}} (\pi(k-1, w, u) \times q(v|w, u) \times e(x_k|v))$$

$S_5 = \{D, N, V, P\}$

$\pi(6, N, P) \times q(D|N, P) \times e(\text{the}|D)$

$\begin{matrix} D \\ N \\ V \\ P \end{matrix} \quad P \quad D$

1 2 3 4 5 6 7 8

The man saw the dog with the telescope

$$\pi(7, P, D) = \max_{w \in \{D, N, V, P\}} (\pi(6, w, P) \times q(D|w, P) \times e(\text{the}|D))$$

- The Viterbi Algorithm

The Viterbi Algorithm

INPUT = x_1, x_2, \dots, x_n

OUTPUT = $\max_{y_1, \dots, y_{n+1}} p(x_1, \dots, x_n, y_1, \dots, y_{n+1})$

Input: a sentence $x_1 \dots x_n$, parameters $q(s|u, v)$ and $e(x|s)$.

Initialization: Set $\pi(0, *, *) = 1$

Definition: $S_{-1} = S_0 = \{*\}$, $S_k = S$ for $k \in \{1 \dots n\}$

Algorithm:

- For $k = 1 \dots n$,
 - For $u \in S_{k-1}$, $v \in S_k$,

$$\pi(k, u, v) = \max_{w \in S_{k-2}} (\pi(k-1, w, u) \times q(v|w, u) \times e(x_k|v))$$
- Return $\max_{u \in S_{n-1}, v \in S_n} (\pi(n, u, v) \times q(\text{STOP}|u, v))$

- The Viterbi Algorithm with Backpointers

The Viterbi Algorithm with Backpointers

INPUT = x_1, \dots, x_n

OUTPUT = $\arg \max_{y_1, \dots, y_{n+1}} p(x_1, \dots, x_n, y_1, \dots, y_{n+1})$

Input: a sentence $x_1 \dots x_n$, parameters $q(s|u, v)$ and $e(x|s)$.

Initialization: Set $\pi(0, *, *) = 1$

Definition: $S_{-1} = S_0 = \{*\}$, $S_k = S$ for $k \in \{1 \dots n\}$

Algorithm:

- For $k = 1 \dots n$,
 - For $u \in S_{k-1}$, $v \in S_k$,

$$\pi(k, u, v) = \max_{w \in S_{k-2}} (\pi(k-1, w, u) \times q(v|w, u) \times e(x_k|v))$$

$$bp(k, u, v) = \arg \max_{w \in S_{k-2}} (\pi(k-1, w, u) \times q(v|w, u) \times e(x_k|v))$$
- Set $(y_{n-1}, y_n) = \arg \max_{(u, v)} (\pi(n, u, v) \times q(\text{STOP}|u, v))$
- For $k = (n-2) \dots 1$, $y_k = bp(k+2, y_{k+1}, y_{k+2})$
- Return the tag sequence $y_1 \dots y_n$

- use dynamic programming to recover the arg max
- runtime complexity is $O(n|S|^3)$
 - linear in the length of the sequence and brute force was exponential in time to the length of the sequence
- Pros and Cons
 - Hidden Markov Model taggers are very simple to train (just need to compile counts from the training corpus)

- Perform relatively well
- Main difficulty is modeling
 - $e(\text{word}|\text{tag})$ can be very difficult if “words” are complex