

Natural Language Processing

Lecture 3 - Log Linear Models

April 2025

Eyal Ben-David

*Partially Based on Coursera's NLP
course by Michael Collins.

(DON'T FORGET RECORDING!!)

Course Logistics & updates

- First assignment will be released on this week
 - You should have a partner & a machine
- Office hours will be dedicated to helping you in your HW assignments

(Recap) Language Modeling

Goal:

Assign a probability to a given sentence or sequence of words.

$P(\text{"Insanity is doing the same thing over and over again and expecting different results"}) = ?$

Another way of thinking about this:

Given a sequence of words, what is likely to be the next word?

$P(\text{"expecting"} \mid \text{"Insanity is doing the same thing over and over again and..."}) = ?$

(Recap) Bigram Language Model

Assumes a **first order markov** assumption.

→ The Bigram MLE of text W is:

$$P(W) = \prod_{i=0}^n p(w_i | w_{i-1}) = \prod_{i=0}^n \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_i)}$$

input
vector

(One-hot) –

והיה עליו שם נח.
↓
מקורו.

ב. חמ"ס. תל"ח
הצ"ח.

$$\left. \begin{array}{l} \text{---} \\ \text{---} \end{array} \right\} \text{---}$$

פירא
המלך

$$\left\{ \begin{array}{l} \text{---} \end{array} \right\}$$
$$\log (p(w \mid u))$$

ਪ੍ਰਿਥਵੀ ਨੂੰ ਚੰਦਰਮਾ ਦਾ ਪ੍ਰਭਾਵ

(Recap) Trigram HMM

Definition

For any sentence x_1, x_2, \dots, x_n and any sequence of tags y_1, y_2, \dots, y_{n+1} , the joint probability of the sentence and the tags is:

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

(Recap) Main Drawbacks of N-gram Modeling

- How do we decide what N would be?
- Storage
- Sparsity
 - How to deal with OOV (unigrams, n-grams)?
 - We can perform:
 - Smoothness
 - Linear interpolation of N-gram models

Today's Agenda

- **Motivation - modeling more information**
- Log linear models
- Parameter estimation (log linear)
- Log linear model for tagging (MEMM)
- Inference

Four-gram Language Models

Estimate: $P(\text{"expecting"} \mid \text{"Insanity is doing the same thing over and over again and..."})$

- Given "history" - w_1, w_2, \dots, w_{n-1}
- Here, "history" is based only on n-grams

Thus, we estimate as follows:

$$P(\text{expecting} \mid w_1, w_2, \dots, w_{n-1}) = \lambda_1 \cdot p_{ML}(\text{expecting} \mid \text{over again and}) + \lambda_2 \cdot p_{ML}(\text{expecting} \mid \text{again and}) + \lambda_3 \cdot p_{ML}(\text{expecting} \mid \text{and}) + \lambda_4 \cdot p_{ML}(\text{expecting})$$

hyper-parameter (pointing to λ_1)

back-off (pointing to the sum of terms)

Four-gram Language Models

$$\begin{aligned} P(\textit{expecting} \mid w_1, w_2, \dots, w_{n-1}) = & \lambda_1 \cdot p_{ML}(\textit{expecting} \mid \textit{over again and}) + \\ & \lambda_2 \cdot p_{ML}(\textit{expecting} \mid \textit{again and}) + \\ & \lambda_3 \cdot p_{ML}(\textit{expecting} \mid \textit{and}) + \\ & \lambda_4 \cdot p_{ML}(\textit{expecting}) \end{aligned}$$

Why just words? Other features may be useful also.

List of Useful Features (besides n-grams)

$$P(\textit{expecting} \mid w_1, w_2, \dots, w_{n-1}) =$$

“מילת חיבור”

- “expecting” does not occur on w_1, w_2, \dots, w_{n-1}
- Previous word is a conjunction (מילת חיבור)
- “doing” and “and” appears in w_1, w_2, \dots, w_{n-1}
- who wrote the text
- what is the text about
- the grammatical structure of the sentence

Try suggesting more...

How To Naively Implement This

$$\begin{aligned} P(\text{expecting} \mid w_1, w_2, \dots, w_{n-1}) = & \lambda_1 \cdot p_{ML}(\text{expecting} \mid w_1, w_2, \dots, w_{n-1} \neq \text{expecting}) + \\ & \lambda_2 \cdot p_{ML}(\text{expecting} \mid \overset{\text{Part of Speech.}}{POS}(w_{n-1}) = CONJ) + \\ & \lambda_3 \cdot p_{ML}(\text{expecting} \mid \{\text{doing, and}\} \in w_1, w_2, \dots, w_{n-1}) + \\ & \lambda_4 \cdot p_{ML}(\text{expecting} \mid \text{text writer}) + \\ & \lambda_5 \cdot p_{ML}(\text{expecting} \mid \text{text subject}) + \\ & \lambda_6 \cdot p_{ML}(\text{expecting} \mid \text{grammar}) \end{aligned}$$

The Naive Implementation (as a thought experiment)

$$\begin{aligned} P(\text{expecting} \mid w_1, w_2, \dots, w_{n-1}) = & \lambda_1 \cdot p_{ML}(\text{expecting} \mid w_1, w_2, \dots, w_{n-1} \neq \text{expecting}) + \\ & \lambda_2 \cdot p_{ML}(\text{expecting} \mid POS(w_{n-1}) = CONJ) + \\ & \lambda_3 \cdot p_{ML}(\text{expecting} \mid \{\text{doing, and}\} \in w_1, w_2, \dots, w_{n-1}) + \\ & \lambda_4 \cdot p_{ML}(\text{expecting} \mid \text{text writer}) + \\ & \lambda_5 \cdot p_{ML}(\text{expecting} \mid \text{text subject}) + \\ & \lambda_6 \cdot p_{ML}(\text{expecting} \mid \text{grammar}) \end{aligned}$$

Adds many hyperparameters (lambdas) to our model.

- In practice, this quickly becomes unwieldy (מגושם, מסורבל)

In Today's lecture: Log linear models can incorporate massive features in an elegant way



A Second Example: Part of Speech Tagging

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

...

A Second Example: Part of Speech Tagging

token classification

Hispaniola/**NNP** quickly/**RB** became/**VB** an/**DT** important/**JJ**
base/**??** from which Spain expanded its empire into the rest of the
Western Hemisphere .

- There are many possible tags in the position **??**
{NN, NNS, Vt, Vi, IN, DT, ...}

- The task: model the distribution

$$p(t_i | t_1, \dots, t_{i-1}, w_1 \dots w_n)$$

where t_i is the i 'th tag in the sequence, w_i is the i 'th word

A Second Example: Part of Speech Tagging

Hispaniola/**NNP** quickly/**RB** became/**VB** an/**DT** important/**JJ**
base/**??** from which Spain expanded its empire into the rest of the
Western Hemisphere .

- The task: model the distribution

$$p(t_i | t_1, \dots, t_{i-1}, w_1 \dots w_n)$$

where t_i is the i 'th tag in the sequence, w_i is the i 'th word

- Again: many “features” of $t_1, \dots, t_{i-1}, w_1 \dots w_n$ may be relevant

$q_{ML}(\text{NN}$		$w_i = \text{base})$
$q_{ML}(\text{NN}$		$t_{i-1} \text{ is JJ})$
$q_{ML}(\text{NN}$		$w_i \text{ ends in “e”})$
$q_{ML}(\text{NN}$		$w_i \text{ ends in “se”})$
$q_{ML}(\text{NN}$		$w_{i-1} \text{ is “important”})$
$q_{ML}(\text{NN}$		$w_{i+1} \text{ is “from”})$



Today's Agenda

- Motivation - modeling more information
- **Log linear models**
- Parameter estimation (log linear)
- Log linear model for tagging (MEMM)
- Inference

Log-Linear Models

The general problem:

- We have an input X
- Have a finite label set Y
- Aim to provide a conditional probability $p(y|x)$

For any $x \in X$ and $y \in Y$

****Note, this is in contrast to HMM that provides a joint probability $p(y,x)$.**

Example - Language Modeling

- X is a “history” w_1, w_2, \dots, w_{n-1} .

e.g.,

Third, the notion “grammatical in English” cannot be identified in any way with the notion “high order of statistical approximation to English”. It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) has ever occurred in an English discourse. Hence, in any statistical

- y is an “outcome” w_n

Feature Vector Representation

- Aim is to provide a conditional probability $p(y \mid x)$ for “decision” y given “history” x
- A feature is a function $f_k(x, y) \in \mathbb{R}$
(Often **binary features** or **indicator functions** $f_k(x, y) \in \{0, 1\}$).
- Say we have m features f_k for $k \in 1, 2, \dots, m$
→ Our feature vector $f(x, y) \in \mathbb{R}^m$

Language Modeling

Third, the notion “grammatical in English” cannot be identified in any way with the notion “high order of statistical approximation to English”. It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) has ever occurred in an English discourse. Hence, in any statistical

- Example features:

$$\begin{aligned}
 f_1(x, y) &= \begin{cases} 1 & \text{if } y = \text{model} \\ 0 & \text{otherwise} \end{cases} \\
 f_2(x, y) &= \begin{cases} 1 & \text{if } y = \text{model and } w_{i-1} = \text{statistical} \\ 0 & \text{otherwise} \end{cases} \\
 f_3(x, y) &= \begin{cases} 1 & \text{if } y = \text{model, } w_{i-2} = \text{any, } w_{i-1} = \text{statistical} \\ 0 & \text{otherwise} \end{cases}
 \end{aligned}$$

האם יש קשר בין המודל לטקסט? $f(x, y)$ מייצגת את ההסתברות של y בהינתן x .

$$F(x, y) = \begin{bmatrix} f_1(x, y) \\ f_2(x, y) \\ f_3(x, y) \end{bmatrix} \approx \lambda_k$$

האם יש קשר בין המודל לטקסט? $f(x, y)$ מייצגת את ההסתברות של y בהינתן x .

Language Modeling

Third, the notion “grammatical in English” cannot be identified in any way with the notion “high order of statistical approximation to English”. It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) has ever occurred in an English discourse. Hence, in any statistical

- Example features:

$$f_1(x, y) = \begin{cases} 1 & \text{if } y = \text{model} \\ 0 & \text{otherwise} \end{cases}$$

$$f_2(x, y) = \begin{cases} 1 & \text{if } y = \text{model and } w_{i-1} = \text{statistical} \\ 0 & \text{otherwise} \end{cases}$$

$$f_3(x, y) = \begin{cases} 1 & \text{if } y = \text{model, } w_{i-2} = \text{any, } w_{i-1} = \text{statistical} \\ 0 & \text{otherwise} \end{cases}$$

Q1: What is x? הסמל
Q2: What is y? מילה לבחירה

Language Modeling

Third, the notion “grammatical in English” cannot be identified in any way with the notion “high order of statistical approximation to English”. It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) has ever occurred in an English discourse. Hence, in any statistical

- Example features:

$$\text{(Unigram)} \quad f_1(x, y) = \begin{cases} 1 & \text{if } y = \text{model} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{(Bigram)} \quad f_2(x, y) = \begin{cases} 1 & \text{if } y = \text{model and } w_{i-1} = \text{statistical} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{(Trigram)} \quad f_3(x, y) = \begin{cases} 1 & \text{if } y = \text{model, } w_{i-2} = \text{any, } w_{i-1} = \text{statistical} \\ 0 & \text{otherwise} \end{cases}$$

Language Modeling

$$f_4(x, y) = \begin{cases} 1 & \text{if } y = \text{model}, w_{i-2} = \text{any} \\ 0 & \text{otherwise} \end{cases}$$

$$f_5(x, y) = \begin{cases} 1 & \text{if } y = \text{model}, w_{i-1} \text{ is an adjective} \\ 0 & \text{otherwise} \end{cases}$$

$$f_6(x, y) = \begin{cases} 1 & \text{if } y = \text{model}, w_{i-1} \text{ ends in "ical"} \\ 0 & \text{otherwise} \end{cases}$$

$$f_7(x, y) = \begin{cases} 1 & \text{if } y = \text{model}, \text{author} = \text{Chomsky} \\ 0 & \text{otherwise} \end{cases}$$

$$f_8(x, y) = \begin{cases} 1 & \text{if } y = \text{model}, \text{"model"} \text{ is not in } w_1, \dots, w_{i-1} \\ 0 & \text{otherwise} \end{cases}$$

$$f_9(x, y) = \begin{cases} 1 & \text{if } y = \text{model}, \text{"grammatical"} \text{ is in } w_1, \dots, w_{i-1} \\ 0 & \text{otherwise} \end{cases}$$

In Practice

- We follow a markovian assumption to the define the history (X) for w_i
- E.g., a trigram based history: $\text{The cat sat on the mat}$

$$X_i = \text{History for } w_i = \langle t_{i-2}, t_{i-1}, t_i, w_{i-2}, w_{i-1}, w_i \rangle$$

- We would probably include a feature for each tri-gram that appeared enough times in the training data.

$$f_{N(u,v,w)}(x, y) = \begin{cases} 1 & \text{if } y = w, w_{i-2} = u, w_{i-1} = v \\ 0 & \text{otherwise} \end{cases}$$

**** $N(u,v,w)$ maps each trigram to a unique integer.**

The Final Result

- We can come up with any question regarding the history and define a feature.

Note, the history includes the tags and the words.

For a given history $x \in X$ each label $y \in Y$

is mapped to a different feature vector:

$$\begin{aligned}
 f(\langle x_{i-2} = \text{any}, x_{i-1} = \text{statistical} \rangle, x_i = \text{outcome}) &= 1001011001001100110 \\
 f(\langle x_{i-2} = \text{any}, x_{i-1} = \text{statistical} \rangle, x_i = \text{time}) &= 0110010101011110010 \\
 f(\langle x_{i-2} = \text{any}, x_{i-1} = \text{statistical} \rangle, x_i = \text{background}) &= 0001111101001100100 \\
 f(\langle x_{i-2} = \text{any}, x_{i-1} = \text{statistical} \rangle, \underbrace{x_i}_{y_i} = \text{course}) &= 00010110110000000010
 \end{aligned}$$

For POS Tagging

- Each X is a “history” $\langle t_1, t_2, \dots, t_{n-1}, w_1, w_2, \dots, w_n, i \rangle$
- y is a POS tag. E.g., NN, NNS, DT
- We have m features $f_k(x, y)$ for $k=1, \dots, m$

time-stamp
→ #word

Example features:

$$f_1(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & \text{if current word } w_i \text{ is base and } y = \text{Vt} \\ 0 & \text{otherwise} \end{cases}$$

$$f_2(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & \text{if current word } w_i \text{ ends in ing and } y = \text{VBG} \\ 0 & \text{otherwise} \end{cases}$$

...

The Full Feature Set in Ratnaparkhi, 1996

- ▶ Word/tag features for all word/tag pairs, e.g.,

$$f_{100}(x, y) = \begin{cases} 1 & \text{if current word } w_i \text{ is base and } y = \text{Vt} \\ 0 & \text{otherwise} \end{cases}$$

לפי מילה קיימת שורה בנקודה.

לפי היותו של x|V| - עבר למחרת של המילה לראשון בסוף מילה שלל.

- ▶ Spelling features for all prefixes/suffixes of length ≤ 4 , e.g.,

$$f_{101}(x, y) = \begin{cases} 1 & \text{if current word } w_i \text{ ends in ing and } y = \text{VBG} \\ 0 & \text{otherwise} \end{cases}$$

לפי: נסתיק יתוף סומא של ארץ 4
וק שומאלי כפחה 6 פתחם.

$$f_{102}(h, t) = \begin{cases} 1 & \text{if current word } w_i \text{ starts with pre and } y = \text{NN} \\ 0 & \text{otherwise} \end{cases}$$

$$F(x, y) = \begin{bmatrix} f_{100} \\ f_{101} \\ f_{102} \end{bmatrix}$$

The Full Feature Set in Ratnaparkhi, 1996

► Contextual Features, e.g.,

$$f_{103}(x, y) = \begin{cases} 1 & \text{if } \langle t_{i-2}, t_{i-1}, y \rangle = \langle \text{DT}, \text{JJ}, \text{Vt} \rangle \\ 0 & \text{otherwise} \end{cases}$$

$$f_{104}(x, y) = \begin{cases} 1 & \text{if } \langle t_{i-1}, y \rangle = \langle \text{JJ}, \text{Vt} \rangle \\ 0 & \text{otherwise} \end{cases}$$

$$f_{105}(x, y) = \begin{cases} 1 & \text{if } \langle y \rangle = \langle \text{Vt} \rangle \\ 0 & \text{otherwise} \end{cases}$$

$$f_{106}(x, y) = \begin{cases} 1 & \text{if previous word } w_{i-1} = \textit{the} \text{ and } y = \text{Vt} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{107}(x, y) = \begin{cases} 1 & \text{if next word } w_{i+1} = \textit{the} \text{ and } y = \text{Vt} \\ 0 & \text{otherwise} \end{cases}$$

The Final Result (POS)

- We can come up with any question regarding the history and define a feature.

Note, the history includes the tags and the words.

For a given history $x \in X$ each label $y \in Y$ is mapped to a different feature vector:

$$f(x = \langle JJ, DT, \langle Hispaniola, \dots \rangle, 6 \rangle, y = Vt) = 1001011001001100110$$

$$f(x = \langle JJ, DT, \langle Hispaniola, \dots \rangle, 6 \rangle, y = JJ) = 0110010101011110010$$

$$f(x = \langle JJ, DT, \langle Hispaniola, \dots \rangle, 6 \rangle, y = NN) = 0001111101001100100$$

$$f(x = \langle JJ, DT, \langle Hispaniola, \dots \rangle, 6 \rangle, y = IN) = 00010110110000000010$$

Parameter Vector

- Given features f_k for $k \in 1, 2, \dots, m$

Define a **parameter vector** $v \in \mathbb{R}^m$

Each (x,y) pair is mapped into a **score**: $v \cdot f(x, y) = \sum_{k=1}^m v_k \cdot f_k(x, y)$

Scoring
function

In our language modeling example:

$$v \cdot f(x, model) = 5.6$$

$$v \cdot f(x, is) = 1.5$$

$$v \cdot f(x, models) = 4.5$$

$$v \cdot f(x, the) = -3.2$$

$$v \cdot f(x, of) = 1.3$$

...

Log-Linear Models

- We have an input domain \mathcal{X} and a finite label space \mathcal{Y}
- A **feature** vector $f(x, y) \in \mathbb{R}^m$
- A **parameter** vector $v \in \mathbb{R}^m$
- We aim to provide a **conditional probability** $p(y | x)$ for any $x \in \mathcal{X}$ and $y \in \mathcal{Y}$
- We'll use the **log-linear model** to define

$$p(y | x; v) = \frac{e^{v \cdot f(x, y)}}{\sum_{y' \in \mathcal{Y}} e^{v \cdot f(x, y')}}$$

Soft max
distribution

Name Origin - Log Linear Models

$$\log (p(y \mid x; v)) = \underbrace{v \cdot f(x, y)}_{\text{Linear term}} - \underbrace{\log \left(\sum_{y' \in \mathcal{Y}} e^{v \cdot f(x, y')} \right)}_{\text{Normalization term}}$$

כלל min/max נכנס \leftarrow (היחס) ויחידות \rightarrow \bar{p}

Today's Agenda

- Motivation - modeling more information
- Log linear models
- **Parameter estimation (log linear)**
- Log linear model for tagging (MEMM)
- Inference

Maximum Likelihood Estimation

- Given a training sample $(x^{(i)}, y^{(i)})$ for $i = 1, \dots, n$,
each $(x^{(i)}, y^{(i)}) \in \mathcal{X} \times \mathcal{Y}$

$$v_{ML} = \arg \max_{v \in \mathbb{R}^m} L(v)$$

gain → not loss

$$L(v) = \sum_{i=1}^n \log \left(p(y^{(i)} \mid x^{(i)}; v) \right) = \sum_{i=1}^n \left[v \cdot f(x^{(i)}, y^{(i)}) - \log \left(\sum_{y' \in \mathcal{Y}} e^{v \cdot f(x^{(i)}, y')} \right) \right]$$

Calculating the Maximum Likelihood Estimates

- Need to maximize:

$$L(v) = \sum_{i=1}^n v \cdot f(x^{(i)}, y^{(i)}) - \sum_{i=1}^n \log \sum_{y' \in \mathcal{Y}} e^{v \cdot f(x^{(i)}, y')}$$

- Calculating gradients:

$$\begin{aligned} \frac{dL(v)}{dv_k} &= \sum_{i=1}^n f_k(x^{(i)}, y^{(i)}) - \sum_{i=1}^n \frac{\sum_{y' \in \mathcal{Y}} f_k(x^{(i)}, y') e^{v \cdot f(x^{(i)}, y')}}{\sum_{z' \in \mathcal{Y}} e^{v \cdot f(x^{(i)}, z')}} \\ &= \sum_{i=1}^n f_k(x^{(i)}, y^{(i)}) - \sum_{i=1}^n \sum_{y' \in \mathcal{Y}} f_k(x^{(i)}, y') \frac{e^{v \cdot f(x^{(i)}, y')}}{\sum_{z' \in \mathcal{Y}} e^{v \cdot f(x^{(i)}, z')}} \end{aligned}$$

החלק הזה

$$= \underbrace{\sum_{i=1}^n f_k(x^{(i)}, y^{(i)})}_{\text{Empirical counts}} - \underbrace{\sum_{i=1}^n \sum_{y' \in \mathcal{Y}} f_k(x^{(i)}, y') p(y' | x^{(i)}; v)}_{\text{Expected counts}}$$

bias $\rightarrow \hat{\theta} - \theta$

למה ב-f
הוא הפונקציה

החלק הזה
הוא הפונקציה

Empirical counts

Expected counts

החלק הזה
הוא הפונקציה

Gradient Ascent

- We need to maximize $L(v)$, where:

$$\frac{dL(v)}{dv} = \sum_{i=1}^n f(x^{(i)}, y^{(i)}) - \sum_{i=1}^n \sum_{y' \in \mathcal{Y}} f(x^{(i)}, y') p(y' | x^{(i)}; v)$$

Initialization: $v = 0$

Iterate until convergence:

- ▶ Calculate $\Delta = \frac{dL(v)}{dv}$
- ▶ Calculate $\beta_* = \operatorname{argmax}_{\beta} L(v + \beta \Delta)$ (Line Search)
- ▶ Set $v \leftarrow v + \beta_* \Delta$

L2 Regularization

- We add a penalty for large weights:

- ▶ Modified loss function

$$L(v) = \sum_{i=1}^n v \cdot f(x^{(i)}, y^{(i)}) - \sum_{i=1}^n \log \sum_{y' \in \mathcal{Y}} e^{v \cdot f(x^{(i)}, y')} - \frac{\lambda}{2} \sum_{k=1}^m v_k^2$$

- ▶ Calculating gradients:

$$\frac{dL(v)}{dv_k} = \underbrace{\sum_{i=1}^n f_k(x^{(i)}, y^{(i)})}_{\text{Empirical counts}} - \underbrace{\sum_{i=1}^n \sum_{y' \in \mathcal{Y}} f_k(x^{(i)}, y') p(y' | x^{(i)}; v)}_{\text{Expected counts}} - \lambda v_k$$

Handwritten notes in Hebrew:
 - מנסים להימנע
 - מנתונים רבים
 - Overfitting

Today's Agenda

- Motivation - modeling more information
- Log linear models
- Parameter estimation (log linear)
- **Log linear model for tagging (MEMM)**
- Inference

Log Linear Models for Tagging
Maximum-entropy Markov Models
(MEMM)

Recap: Part of Speech (POS)

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

...

Recap: Named Entity Recognition (NER)

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

OUTPUT: Profits soared at [Company Boeing Co.], easily topping forecasts on [Location Wall Street], as their CEO [Person Alan Mulally] announced first quarter results.

Recap: Named Entity Recognition (NER)

INPUT:

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

OUTPUT:

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

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

Goal

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** pro/

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

Log Linear Models for Tagging

- ▶ We have an input sentence $w_{[1:n]} = w_1, w_2, \dots, w_n$
(w_i is the i 'th word in the sentence)

- ▶ We have a tag sequence $t_{[1:n]} = t_1, t_2, \dots, t_n$
(t_i is the i 'th tag in the sentence)

- ▶ We'll use an log-linear model to define

$$p(t_1, t_2, \dots, t_n | w_1, w_2, \dots, w_n)$$

for any sentence $w_{[1:n]}$ and tag sequence $t_{[1:n]}$ of the same length.

(Note: contrast with HMM that defines $p(t_1 \dots t_n, w_1 \dots w_n)$)

- ▶ Then the most likely tag sequence for $w_{[1:n]}$ is

$$t_{[1:n]}^* = \operatorname{argmax}_{t_{[1:n]}} p(t_{[1:n]} | w_{[1:n]})$$

Trigram Log Linear Tagger

A Trigram Log-Linear Tagger:

$$p(t_{[1:n]} | w_{[1:n]}) = \prod_{j=1}^n p(t_j | w_1 \dots w_n, t_1 \dots t_{j-1}) \quad \text{Chain rule}$$

$$= \prod_{j=1}^n p(t_j | w_1, \dots, w_n, t_{j-2}, t_{j-1})$$

Independence assumptions

- ▶ We take $t_0 = t_{-1} = *$
- ▶ Independence assumption: each tag only depends on previous two tags

$$p(t_j | w_1, \dots, w_n, t_1, \dots, t_{j-1}) = p(t_j | w_1, \dots, w_n, t_{j-2}, t_{j-1})$$

Example - History Tuple

Hispaniola/**NNP** quickly/**RB** became/**VB** an/**DT** important/**JJ**
base/**??** from which Spain expanded its empire into the rest of the
Western Hemisphere .

- ▶ $t_{-2}, t_{-1} = \text{DT, JJ}$
- ▶ $w_{[1:n]} = \langle \text{Hispaniola, quickly, became, } \dots, \text{Hemisphere, .} \rangle$
- ▶ $i = 6$
- ▶ A **history** is a 4-tuple $\langle t_{-2}, t_{-1}, w_{[1:n]}, i \rangle$
- ▶ t_{-2}, t_{-1} are the previous two tags.
- ▶ $w_{[1:n]}$ are the n words in the input sentence.
- ▶ i is the index of the word being tagged

Feature Set (Ratnaparkhy, 1996)

- ▶ Word/tag features for all word/tag pairs, e.g.,

$$f_{100}(h, t) = \begin{cases} 1 & \text{if current word } w_i \text{ is base and } t = \text{Vt} \\ 0 & \text{otherwise} \end{cases}$$

- ▶ Spelling features for all prefixes/suffixes of length ≤ 4 , e.g.,

$$f_{101}(h, t) = \begin{cases} 1 & \text{if current word } w_i \text{ ends in ing and } t = \text{VBG} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{102}(h, t) = \begin{cases} 1 & \text{if current word } w_i \text{ starts with pre and } t = \text{NN} \\ 0 & \text{otherwise} \end{cases}$$

....

Today's Agenda

- Motivation - modeling more information
- Log linear models
- Parameter estimation (log linear)
- Log linear model for tagging (MEMM)
- **Inference**

The Viterbi Algorithm - MEMM

Problem: for an input $w_1 \dots w_n$, find

$$\arg \max_{t_1 \dots t_n} p(t_1 \dots t_n \mid w_1 \dots w_n)$$

We assume that p takes the form

$$p(t_1 \dots t_n \mid w_1 \dots w_n) = \prod_{i=1}^n q(t_i \mid t_{i-2}, t_{i-1}, w_{[1:n]}, i)$$

(In our case $q(t_i \mid t_{i-2}, t_{i-1}, w_{[1:n]}, i)$ is the estimate from a log-linear model.)

The Viterbi Algorithm - MEMM

- ▶ Define n to be the length of the sentence
- ▶ Define

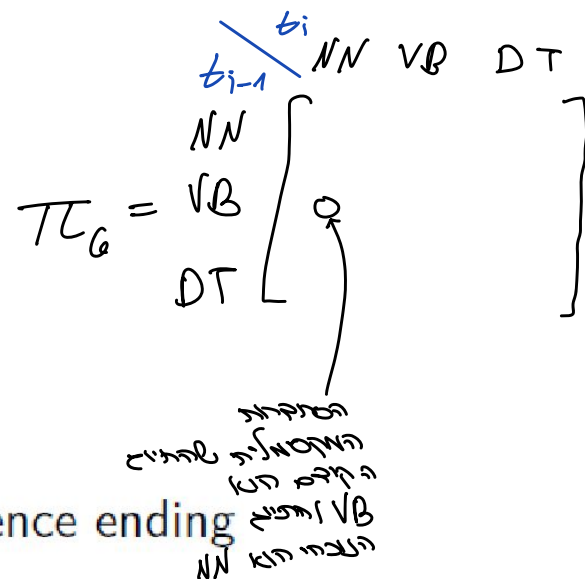
$$r(t_1 \dots t_k) = \prod_{i=1}^k q(t_i | t_{i-2}, t_{i-1}, w_{[1:n]}, i)$$

- ▶ 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,

$$\pi(k, u, v) = \max_{\langle t_1, \dots, t_{k-2} \rangle} r(t_1 \dots t_{k-2}, u, v)$$



A Recursive Definition

Base case:

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

(with arrows pointing to the two asterisks)

$$\pi_0 = \begin{pmatrix} \dots & * \\ - & 0 & - \\ * & & 1 \end{pmatrix}$$

$$\pi_1 = \begin{pmatrix} \dots & * \\ \dots & \dots & \dots \\ 0 & 0 & \dots \\ \vdots & \vdots & \vdots \\ * & \dots & 0 \end{pmatrix}$$

(with blue arrows pointing to the top-left and top-right elements of the matrix, and a green underline under the bottom row)

Recursive definition:

For any $k \in \{1 \dots n\}$, for any $u \in \mathcal{S}_{k-1}$ and $v \in \mathcal{S}_k$:

$$\pi(k, u, v) = \max_{t \in \mathcal{S}_{k-2}} (\pi(k-1, t, u) \times q(v|t, u, w_{[1:n]}, k))$$

where \mathcal{S}_k is the set of possible tags at position k

The full Algorithm (with backpointers)

Input: a sentence $w_1 \dots w_n$, log-linear model that provides $q(v|t, u, w_{[1:n]}, i)$ for any tag-trigram t, u, v , for any $i \in \{1 \dots n\}$

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

Algorithm:

- ▶ For $k = 1 \dots n$,
 - ▶ For $u \in \mathcal{S}_{k-1}, v \in \mathcal{S}_k$,

$$\pi(k, u, v) = \max_{t \in \mathcal{S}_{k-2}} (\pi(k-1, t, u) \times q(v|t, u, w_{[1:n]}, k))$$

$$bp(k, u, v) = \arg \max_{t \in \mathcal{S}_{k-2}} (\pi(k-1, t, u) \times q(v|t, u, w_{[1:n]}, k))$$

- ▶ Set $(t_{n-1}, t_n) = \arg \max_{(u,v)} \pi(n, u, v)$
- ▶ For $k = (n-2) \dots 1$, $t_k = bp(k+2, t_{k+1}, t_{k+2})$
- ▶ **Return** the tag sequence $t_1 \dots t_n$

Beam Search - Heuristic

At each time step (k):

- calculate scores for all possible states (π)
- sort them according to their score values
- pass only the top B (beam width) possible states to the next time step.

Beam Search - Example (B=2)

$$\pi(k, u, v) = \max_{t \in \mathcal{S}_{k-2}} (\pi(k-1, t, u) \times q(v|t, u, w_{[1:n]}, k))$$

(t,u) could only be in {(VB,MD), (MD,VB)}

$\pi(k-1)$	<BOS>	NNP	MD	VB	JJ
<BOS>	0	0	0	0	0
NNP	0	0.1	0.1	0	0
MD	0	0	0	0.2	0.05
VB	0	0	0.2	0	0
JJ	0	0.15	0.5	0	0.15

$\pi(k)$	<BOS>	NNP	MD	VB	JJ
<BOS>	0	0	0	0	0
NNP	0	0			
MD	0		??		
VB	0				
JJ	0				

Beam Search - Example (B=2)

$$\pi(k, u, v) = \max_{t \in \mathcal{S}_{k-2}} (\pi(k-1, t, u) \times q(v|t, u, w_{[1:n]}, k))$$

(t,u) could only be in {(VB,MD), (MD,VB)}

For u not in {MD,VB}, π is automatically set to 0.

$\pi(k-1)$	<BOS>	NNP	MD	VB	JJ
<BOS>	0	0	0	0	0
NNP	0	0.1	0.1	0	0
MD	0	0	0	0.2	0.05
VB	0	0	0.2	0	0
JJ	0	0.15	0.5	0	0.15

$\pi(k)$	<BOS>	NNP	MD	VB	JJ
<BOS>	0	0	0	0	0
NNP	0	0	0	0	0
MD	0		??		
VB	0				
JJ	0	0	0	0	0

Beam Search - Example (B=2)

$$\pi(k, u, v) = \max_{t \in \mathcal{S}_{k-2}} (\pi(k-1, t, u) \times q(v|t, u, w_{[1:n]}, k))$$

(t,u) could only be in {(VB,MD), (MD,VB)}

For u not in {MD,VB}, π is automatically set to 0.

Instead of iterating through all t's, only iterate {MD, VB}.

$\pi(k-1)$	<BOS>	NNP	MD	VB	JJ
<BOS>	0	0	0	0	0
NNP	0	0.1	0.1	0	0
MD	0	0	0	0.2	0.05
VB	0	0	0.2	0	0
JJ	0	0.15	0.5	0	0.15

$\pi(k)$	<BOS>	NNP	MD	VB	JJ
<BOS>	0	0	0	0	0
NNP	0	0	0	0	0
MD	0		??		
VB	0				
JJ	0	0	0	0	0

Evaluating POS

- Accuracy - $\text{\#correct_pred} / \text{\#num_of_words}$
- Confusion matrix
- F1

Confusion Matrix

Ground Truth: The_**DT** boy_**NNP** walked_**VBD**

Our Tagger Predicted: The_**DT** boy_**DT** walked_**NNP**

The Confusion matrix:

<u>Truth \ Predicted</u>	<i>DT</i>	<i>NNP</i>	<i>VBD</i>
<i>DT</i>	1	0	0
<i>NNP</i>	1	0	0
<i>VBD</i>	0	1	0

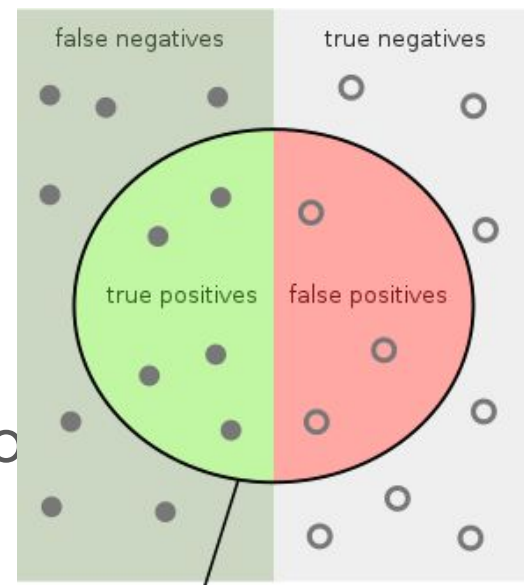
Accuracy:

$$\frac{\# \text{ correct}}{\# \text{ total}} = \frac{\text{sum}(\text{diagonal})}{\text{sum}(\text{matrix})} = \frac{\sum_{i=1}^{|S|} a_{i,i}}{\sum_{i=1}^{|S|} \sum_{j=1}^{|S|} a_{i,j}} = \frac{1}{3}$$

F1

- For each label, we measure its F1 score
 - Based on precision & recall
 - Averaged across all words
 - We average across all words (per label)
- Then, we weight across labels to get a single score

$$2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Next week

- Introduction to deep learning (the basic stuff)
- We need this for better text representation
- See you next week!