

0. HW4A

CS 585

UMass Amherst

10/18/16

1. Log-Linear Models  
    Ⓐ Classif    Ⓑ Sequence

2. Perceptron

3. Structural Perceptron

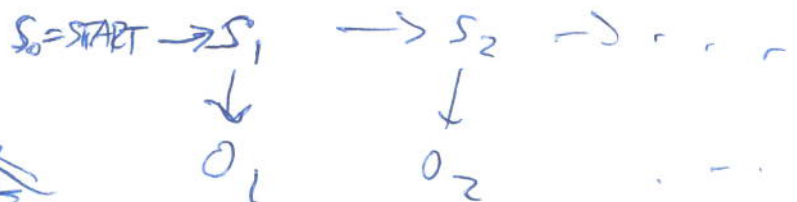
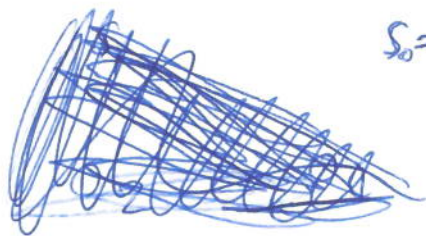
# HW4A

① 1-3: Can answer just for whether holds in general  
    4-5: Numbers                      for most HMM parameters

② implicit  $S_0 = \text{START}$

③ If you know the sequence length is  $n$ :  
    implies  $S_{n+1} = \text{END}$

$$P(o_2 | o_1) \iff P(o_2 | o_1, S_0 = \text{START})$$



# Homework 4 Part A: HMM, Viterbi

CS 585, UMass Amherst, Fall 2016

## Overview

Due **Friday Oct 21**.

Get starter code from the course website's schedule page. You should submit a zipped directory (please don't use other compression formats like .rar) named hw4a\_YOUR-USERNAME that contains:

- your vit.py file
- writeup

Our course's collaboration policy is specified on the website.

## 1 HMM

[15 total points]

Answer the following questions using the transition matrix  $T$  and emission probabilities  $E$  below. Below,  $\Delta$  and  $\square$  are two output variables,  $A$  and  $B$  are two hidden states;  $s_n$  refers to the  $n^{\text{th}}$  hidden state in the sequence and  $o_n$  refers to the  $n^{\text{th}}$  observation.

$$T = \begin{array}{c|ccc} & A & B & \text{END} \\ \hline \text{START} & 0.5 & 0.5 & 0.0 \\ A & 0.2 & 0.3 & 0.5 \\ B & 0.4 & 0.4 & 0.2 \end{array}$$

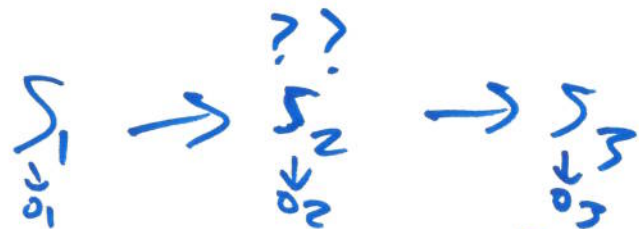
$$E = \begin{array}{c|cc} & \Delta & \square \\ \hline A & 0.5 & 0.5 \\ B & 0.3 & 0.7 \end{array}$$

Handwritten notes:  $s_0 = \text{START}$ ,  $\square \rightarrow s_1$ ,  $s_1 \rightarrow s_2$ ,  $s_2 \rightarrow o_2$ . A diagram shows a box labeled  $s_1$  with a box labeled  $\square$  below it, and an arrow pointing from  $s_1$  to  $s_2$ . Another diagram shows a box labeled  $s_1$  with a box labeled  $\square$  below it, and an arrow pointing from  $s_1$  to  $o_2$ .

- [2 points] Does  $P(o_2 = \Delta | s_1 = B) = P(o_2 = \Delta | o_1 = \square)$ ?
- [2 points] Does  $P(s_2 = B | s_1 = A) = P(s_2 = B | s_1 = A, o_1 = \Delta)$ ?
- [3 points] Does  $P(o_2 = \Delta | s_1 = A) = P(o_2 = \square | s_1 = A, s_3 = A)$ ?
- [3 points] Compute the probability of observing  $\square$  as the first emission of a sequence generated by an HMM with transition matrix  $T$  and emission probabilities  $E$ .
- [5 points] Compute the probability of the first state being  $A$  given that the last token in an observed sequence of length 2 was the token  $\Delta$ .

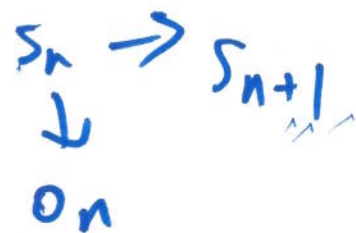
Handwritten note:  $P(s_2 | s_1, \text{I know } n=3) \neq P(s_2 | s_1)$

Handwritten note:  $s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow \text{END}$   
 $\neq P(s_2 | s_1, s_4 = \text{END}) \neq P(s_2 | s_1)$



$$P(S_2 | S_1, S_3) = \frac{P(S_1, S_2, S_3)}{P(S_1, S_3)}$$

$$\begin{aligned}
 &\downarrow \\
 &= P(S_2 | S_1, S_3, o_3) \\
 &= \frac{P(S_1, S_2, S_3)}{\sum_{S_2'} P(S_1, S_2', S_3)}
 \end{aligned}$$



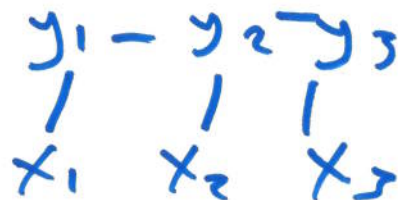
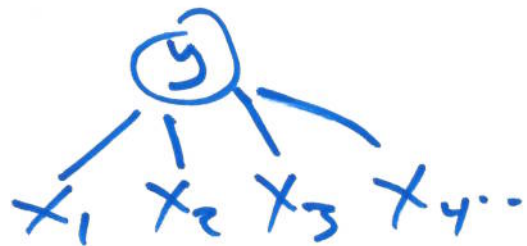
$$\begin{aligned}
 &P(S_n, o_n, S_{n+1}) \\
 &= P(S_n) P(o_n | S_n) P(S_{n+1} | S_n)
 \end{aligned}$$

# Log-Linear (Discriminative) Models

$$P(y|x)$$

y: doc class .... sequence of POS tags

x: text



HMM

$$P(y, x) =$$

$$\prod_t P(y_t | y_{t-1}) P(x_t | y_t)$$

$$\log P(y|x) = C + \theta^T f(x, y)$$

Model Params

Feature Vector for (x, y)

Classif: "Multi-class Log. Reg."

x = "I ♡ cats"

$$f(x, POS) = \begin{pmatrix} \text{POS\_I} & \text{POS\_♡} & \text{POS\_dog} & \dots & \text{NEG\_I} \\ 1 & 1 & 0 & \dots & 0 \end{pmatrix}$$

$$f(x, NEG) = \begin{pmatrix} 0 & 0 & 0 & \dots & 1 \end{pmatrix}$$

$$\theta = \begin{pmatrix} +2.1, & 3.9, & 0.3 & \dots & - & - & - \end{pmatrix}$$

Decision  $\operatorname{argmax}_{y'} P(y'/x) = \theta^\top f(x, y')$

$$\theta^\top f(x, \text{pos}) \quad \text{vs.} \quad \theta^\top f(x, \text{NEG})$$

$$= \sum_{j=1}^J \theta_j f_j(x, \text{pos})$$

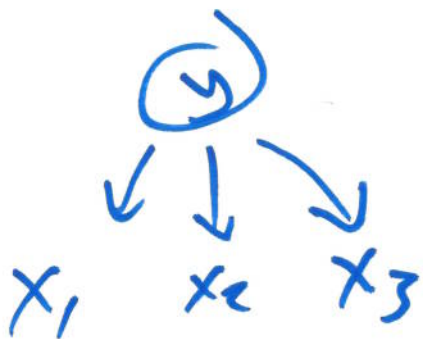
$$\theta^\top f(x, \text{pos})$$

NB as log-Lin?

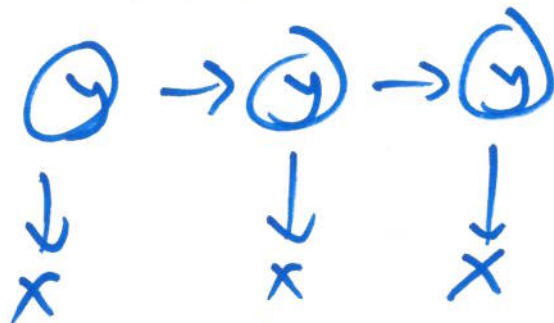
$$f_{w, \text{pos}}(x, y) = \begin{cases} \text{Count of word } w \text{ in } x & \text{when } y = \text{pos} \\ \text{else } 0 \end{cases}$$

$$\theta_{w, \text{pos}} \equiv \log P(w | \text{pos})$$

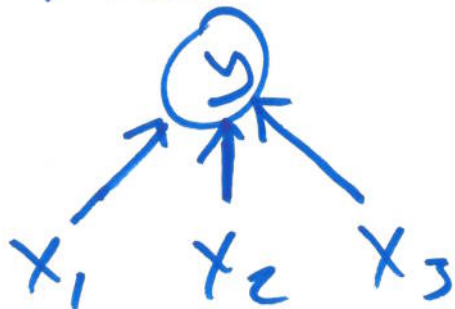
NB



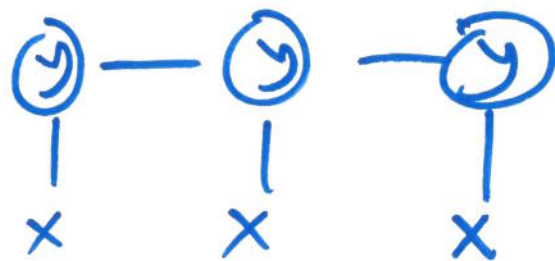
HMM



MLR



(CRF) Conditional Random Fields



features over  
structure  $\vec{y}$

$$\log P(\vec{y} | \vec{x}) = C + \theta^T f(\vec{x}, \vec{y})$$

HMM as log-Linear  $\log P(N|N) \Rightarrow K \times K$

$y = A \quad N \quad N \quad N$   
 $x = \text{happy President Abe Lincoln}$   
 $\log P(\text{Lincoln} | N) \Rightarrow V \times K$

Feat Vector: for  $\log P(y|x) = \dots + \theta^T f(x, y)$

$f(x, y) = \left( \begin{array}{ccc|ccc} \#A, N & A, V & N, N & & & \\ 1 & 0 & 2 & \dots & 1 & 0 & \dots \end{array} \right)$

Trans. Feats  $\Rightarrow A$   
 $K \times K$

obs. feats  $\Rightarrow B$   
 $V \times K$

Params!

↓

Count(tag {k} after {j})    Count(tag {y} appears above {w}')

$\theta = (+2.1, -4.1, +4.4 \dots \dots \dots +110.2 \dots)$

Predict:  $\arg \max_{y'} \theta^T f(x, y') \Leftrightarrow \arg \max_{y'} P(y'|x)$   
Viterbi  $y' = P(y', x)$

Why CRFs?

$$P(y|x) = \frac{1}{Z} \exp(\Theta^T f(x, y))$$

① Discrimina Training of  $\Theta$

② More features

(HMM) Trans from  $j$  to  $k$

(HMM) observe word  $w$  below  $k$

- Word ends in -ing below tag  $k$

- Word begins with a capital letter

- - - - -



Learning: Perceptron Algo  
(Structured)

Rosenblatt 1957

Error-Driven Learning

Converge?

Averaging

Algo

For each pass thru data ( $\sim 30$  iter)

For each  $(x, y^{(gold)})$ :

Pred  $y^* = \operatorname{argmax}_{y'} \theta^T f(x, y')$

Update  $\theta := \theta + r [f(x, y^{(gold)}) - f(x, y^*)]$

$\Rightarrow \theta^T f(x, y^{(gold)})$  gets bigger

What if  $y^* = y^{(gold)}$  ?