

MAXIMUM ENTROPY MARKOV MODELS FOR INFORMATION EXTRACTION AND SEGMENTATION

Andrew McCallum, Dayne Freitag, Fernando Pereira, ICMIL 2000

Presented by
Ruogu Fang

Cornell University
Department of Electrical and Computer Engineering

February 23, 2010 CS 6784 Advanced Topics in Machine Learning

OUTLINE

- Conditional Model
 - Hidden Markov Model
 - Maximum Entropy
- Maximum Entropy Markov Model Framework
 - Model :
 - Exponential Model for Transition
 - Learning:
 - Generalized Iterative Scaling
 - Prediction:
 - State estimation from Observations
- Experimental Results

February 23, 2010 CS 6784 Advanced Topics in Machine Learning

HIDDEN MARKOV MODEL (HMM)

- A hidden Markov model (HMM) is a triple (Π, A, B)
- States: $y \in \{s_1, \dots, s_k\}$
- Output symbols: $x \in \{o_1, \dots, o_m\}$
- Initial State Probabilities $\Pi = \{\pi_i, i \in S\} \quad P(Y_i = y_i)$
 - Specifies where the sequence starts
- State Transition Probabilities $A = \{a_{ij}, i, j \in S\} \quad P(Y_i = x_i | Y_{i-1} = y_{i-1})$
 - Probability that one state succeeds another
- Symbol Emission Probabilities $B = \{b_{ij}, i, j\} \quad P(X_i = x_i | Y_i = y_i)$
 - Probability that observation is generated in this state
- Every output + state sequence has a probability

$$P(x_1, \dots, x_n, y_1, \dots, y_n) = [P(y_1)P(x_1 | y_1)] \prod_{i=2}^n [P(y_i | y_{i-1})P(x_i | y_i)]$$

Limitation ?

February 23, 2010 CS 6784 Advanced Topics in Machine Learning

OBJECTIVE OF MEMM

- Maximum Entropy: Incorporate multiple overlapping features
 - Indentation
 - Numbered questions
 - Styles of paragraph breaks
 - Line length
- HMM: Conditional probability of state sequences given observations

DET → N

↑

hear

The bear chased the cat

DET N V DET N

February 23, 2010 CS 6784 Advanced Topics in Machine Learning

EXPERIMENT SETUP

Experimental Data

- 58 files belonging to 7 Usenet multipart FAQs
- <head>X-NNTP-Poster: NewsHound v1.33
- <head>
- <head>Archive-name: acoofs/part2
- <head>Frequency: monthly
- <head>
- <content>#2:0 What configuration of serial cable should I use
- <answer>
- <answer> Here follows a diagram of the necessary connections
- <answer> programs to work properly. They are as far as I know t
- <answer> based upon by commercial comma software developers fo
- <answer>
- <answer> Pins 1, 4, and 8 must be connected together inside
- <answer> is to avoid the well known serial port chip bugs. Thi
- Procedure: For each FAQ, train on one part, test on others; average

Features in Experiment

- begins-with-number
- begins-with-ordinal
- begins-with-punctuation
- begins-with-question-word
- begins-with-question-word
- begins-with-subject
- blank indented-1-to-4
- contains-alphabet
- contains-bracketed-number
- contains-http
- contains-non-space
- contains-number
- contains-pipe
- contains-question-mark
- contains-question-word
- ends-with-question-mark
- rst-alpha-is-capitalized
- Indented
- indented-1-to-4
- indented-5-to-10
- more-than-one-third-space
- only-punctuation
- prev-is-blank
- prev-begins-with-ordinal
- shorter-than-30

February 23, 2010 CS 6784 Advanced Topics in Machine Learning

MAXIMUM ENTROPY MARKOV MODELS

HMM

$P(y, x) = [P(y_1) \cdot P(x_1 | y_1)] [P(y_2 | y_1) \cdot P(x_2 | y_2)] \dots [P(y_n | y_{n-1}) \cdot P(x_n | y_n)]$

MEMM

$P(y | x) = [P(y_1 | x_1, y_1)] [P(y_2 | x_2, y_2)] \dots [P(y_n | x_n, y_{n-1})]$

Transition Prob
 $P(S_t | S_{t-1})$

Emission Prob
 $P(O_t | S_t)$

→

$P(S_t | S_{t-1}, O_t)$

→

$P(S_{t-1} | S_t, O_t)$

|S|

February 23, 2010 CS 6784 Advanced Topics in Machine Learning

EXPONENTIAL MODEL FOR TRANSITION

- Model transition in terms of multiple, non-independent features of observations
- Objective function: $H(p)$
- Goal: Among all the distributions that satisfy the constraints, choose the one, p^* , that maximizes $H(p)$.

$$p^* = \arg \max_{p \in P} H(p)$$
- Question: How to represent constraints?
- Feature: a binary valued function on events

$$f_j: \mathcal{E} \rightarrow \{0,1\}, \quad \mathcal{E} = A \times B$$
- A: the set of possible classes (e.g., tags in POS tagging)
- B: space of contexts (e.g., neighboring words/ tags in POS tagging)

$$f_j(a,b) = \begin{cases} 1 & \text{if } a = \text{DET} \ \& \ \text{curWord}(b) = \text{"that"} \\ 0 & \text{o.w.} \end{cases}$$

February 23, 2010 CS 6781: Advanced Topics in Machine Learning 7

In courtesy of Fei Xia, Maximum Entropy Model, 02/2006

EXPONENTIAL MODEL FOR TRANSITION (CONT)

- Task: Find $p^* \quad p^* = \arg \max_{p \in P} H(p)$
 where $P = \{p \mid E_p f_j = E_p f_j, j = \{1, \dots, k\}\}$
- Objective function: $H(p)$
- Constraints: $\{E_p f_j = E_p f_j = d_j, j = \{1, \dots, k\}\}$

$\sum_x p(x) = 1 \rightarrow$ Add a feature $f_0(a,b) = 1 \quad \forall a,b$

- Is P empty? $E_p f_0 = E_p f_0 = 1$

$$p(x) = \frac{e^{\sum_{j=1}^k \lambda_j f_j(x)}}{Z}$$

Questions:

- Does p^* exist?
- Is p^* unique?
- What is the form of p^* ?
- How to find p^* ?

Z is the normalizing factor that makes the distribution sum to one across all next states s

February 23, 2010 CS 6781: Advanced Topics in Machine Learning 8

STATE ESTIMATION

HMM

MEMM

- forward probability**
 - probability of producing the observation sequence up to time t and being in state s at time t
$$\alpha_{t+1}(s) = \sum_{s' \in S} \alpha_t(s') \cdot P(s|s') \cdot P(o_{t+1}|s)$$
- Forward probability $\alpha_t(s)$**
 - Probability of being in state s at time t given the observation sequence up to time t
$$\alpha_{t+1}(s) = \sum_{s' \in S} \alpha_t(s') \cdot P_p(s|s', o_{t+1})$$
- Backward probability**
 -
$$\beta_t(s) = \sum_{o_{t+1}} P(s|s', o_{t+1}) \cdot \beta_{t+1}(s)$$

February 23, 2010 CS 6781: Advanced Topics in Machine Learning 9

MODELS TESTED

- MaxEntStateless**
 - A single maximum entropy classifier applied to each line independently.
- TokenHMM**
 - A fully connected HMM with four states, one for each of the line categories, each of which generates individual tokens (groups of alphanumeric characters and individual punctuation characters).
- FeatureHMM**
 - Identical to *tokenHMM*, only the lines in a document are first converted to sequences of features.
- MaxEntMM**
 - The maximum entropy Markov model described in this paper.

February 23, 2010 CS 6781: Advanced Topics in Machine Learning 10

EXPERIMENTAL RESULTS

FAQ Dataset

| Metric | ME-Stateless | TokenHMM | FeatureHMM | MEMM |
|-----------|--------------|----------|------------|------|
| COAP | 0.52 | 0.88 | 0.95 | 0.98 |
| SegPrec | 0.05 | 0.28 | 0.42 | 0.88 |
| SegRecall | 0.38 | 0.15 | 0.52 | 0.68 |

$$COAP \quad P_p(act, pred) = \sum_{i,j} D_p(i,j) \delta_{act}(i,j) \oplus \delta_{pred}(i,j)$$

February 23, 2010 CS 6781: Advanced Topics in Machine Learning 11

SUMMARY

- Combine HMMs and maximum-entropy models into a general model
 - Allow state transitions to depend on non-independent features of the sequence under analysis
 - Conditional model that represents the probability of reaching a state given an observation and the previous state
- Training algorithm
 - Direct generalizations of HMMs algorithms: forward-backward, Baum-Welch
 - Generalized iterative scaling for $P(s|s', o)$
 - Generalized Iterative Scaling (GIS)
- Prediction: Dynamic programming

February 23, 2010 CS 6781: Advanced Topics in Machine Learning 12

February 23, 2010 CS 6784: Advanced Topics in Machine Learning

ADDITIONAL SLIDES

13

February 23, 2010 CS 6784: Advanced Topics in Machine Learning

MAXIMUM ENTROPY

- Goal: Estimate p
- Choose p with maximum entropy (or “uncertainty”) subject to the constraints (or “evidence”)

$$H(p) = - \sum_{a \in A, b \in B} p(x) \log p(x) \quad x = (a, b), \text{ where } a \in A \wedge b \in B$$
- Training
 - a : state to be predicted
 - b : context
 - Example: a =NN, b =tags for current and previous words
 - Learn probabilities of each (a, b) : $P(a, b)$
- Maximum Entropy = Minimum Commitment
 - Model all is known: satisfy a set of constraints
 - Assume nothing about what is unknown: choose the most “uniform” distribution

14

February 23, 2010 CS 6784: Advanced Topics in Machine Learning

MAXIMUM ENTROPY: COIN FLIP

- Toss a coin: $P(H) = p_1, P(T) = p_2$
- Constraint: $p_1 + p_2 = 1$
- Question: What is your estimation of $p=(p_1, p_2)$?
- Answer: Choose the p that maximizes $H(p)$

$$H(p) = - \sum_x p(x) \log p(x)$$

In courtesy of Fei Xia, Maximum Entropy Model, 02/2006

15

February 23, 2010 CS 6784: Advanced Topics in Machine Learning

MAXIMUM ENTROPY: COIN FLIP

- Maximum Entropy = Minimum Commitment
 - Model all is known: satisfy a set of constraints
 - Assume nothing about what is unknown: choose the most “uniform” distribution

In courtesy of Fei Xia, Maximum Entropy Model, 02/2006

16

February 23, 2010 CS 6784: Advanced Topics in Machine Learning

MAXIMUM ENTROPY MARKOV MODEL

17

February 23, 2010 CS 6784: Advanced Topics in Machine Learning

OUTLINE

- Conditional Model
 - Hidden Markov Model
 - Maximum Entropy
- Maximum Entropy Markov Model Framework
 - Model :
 - Exponential Model for Transition
 - Learning:
 - Generalized Iterative Scaling
 - Prediction:
 - State estimation from Observations
- Variations
- Experimental Results

18

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

EXPONENTIAL MODEL FOR TRANSITION (CONT)

- Finite training sample of events: S
- Observed probability of x in S : $\tilde{p}(x)$
- The model p 's probability of x : $p(x)$
- The j^{th} feature:
- Observed expectation of f_j $E_p f_j = \sum_{x \in S} \tilde{p}(x) f_j(x)$
 - Empirical count of f_j
- Model expectation of f_j $E_p f_j = \sum_{x \in \mathcal{X}} p(x) f_j(x)$
- Model's feature expectation = observed feature expectation $E_p f_j = E_{\tilde{p}} f_j$
- How to calculate $E_{\tilde{p}} f_j$? $E_{\tilde{p}} f_j = \sum_{x \in S} \tilde{p}(x) f_j(x) = \frac{\sum_{i=1}^N f_j(x_i)}{N}$

$$f_j(a,b) = \begin{cases} 1 & \text{if } a = \text{DET} \ \& \ \text{curWord}(b) = \text{'thar' } \\ 0 & \text{o.w.} \end{cases}$$

19

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

OUTLINE

- Conditional Model
 - Hidden Markov Model
 - Maximum Entropy
- Maximum Entropy Markov Model Framework
 - Model :
 - Exponential Model for Transition
 - Learning:
 - Generalized Iterative Scaling
 - Prediction:
 - State estimation from Observations
- Variations
- Experimental Results

20

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

PARAMETER ESTIMATION

GENERALIZED ITERATIVE SCALING (DARROCH & RATCLIFF, 1972)

- Setup
 - Obey form of model and constraints

$$p(x) = \frac{e^{\sum_{j=1}^k \lambda_j f_j(x)}}{Z} \quad F_a = \frac{1}{m_a} \sum_{k=1}^m f_a(o_i, s_i)$$

- Additional constraint: $\forall x \in \mathcal{E} \quad \sum_{j=1}^k f_j(x) = C$
- Let $C = \max_{x \in \mathcal{E}} \sum_{j=1}^k f_j(x)$
- Add a new feature f_{k+1} : $\forall x \in \mathcal{E} \quad f_{k+1}(x) = C - \sum_{j=1}^k f_j(x)$

21

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

GENERALIZED ITERATIVE SCALING ALGORITHM

- Compute $F_a, j=1, \dots, k+1$
- Initialize $\lambda_j^{(0)} = 1$ (any values, e.g., 1)
- REPEAT until converge
 - FOR each j
 - Compute $E_{p^{(n)}} f_j = \sum_{x \in \mathcal{E}} p^{(n)}(x) f_j(x)$ **E-Step**
 - where $p^{(n)}(x) = \frac{e^{\sum_{j=1}^k \lambda_j^{(n)} f_j(x)}}{Z}$
 - Update $\lambda_j^{(n+1)} = \lambda_j^{(n)} + \frac{1}{C} (\log \frac{F_a}{E_{p^{(n)}} f_j})$ **M-Step**
 - END FOR
- END REPEAT

22

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

WHAT IS THE FORM OF P*

(RATNAPARKHI, 1997)

$$P = \{p \mid E_p f_j = E_{\tilde{p}} f_j, j = \{1, \dots, k\}\}$$

$$Q = \{p \mid p(x) = \pi \prod_{j=1}^k \alpha_j^{f_j(x)}, \alpha_j > 0\}$$

- Theorem: if $p^* \in P \cap Q$ then $p^* = \arg \max_{p \in P} H(p)$
- Furthermore, p^* is unique.

23

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

APPROXIMATION FOR CALCULATING FEATURE EXPECTATION

$$E_p f_j = \sum_{x \in \mathcal{E}} p(x) f_j(x) = \sum_{a \in A, b \in B} p(a,b) f_j(a,b)$$

$$= \sum_{a \in A, b \in B} p(b) p(a|b) f_j(a,b)$$

$$\approx \sum_{a \in A, b \in B} \tilde{p}(b) p(a|b) f_j(a,b)$$

$$= \sum_{b \in B} \tilde{p}(b) \sum_{a \in A} p(a|b) f_j(a,b)$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{a \in A} p(a|b_i) f_j(a,b_i)$$

24

PROPERTIES OF GIS

- $L(p^{(n+1)}) \geq L(p^{(n)})$
- The sequence is guaranteed to converge to p^* .
- The converge can be very slow.

- The running time of each iteration is $O(NPA)$:
 - N: the training set size
 - P: the number of classes
 - A: the average number of features that are active for a given event (a, b).

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

25

OUTLINE

- Conditional Model
 - Hidden Markov Model
 - Maximum Entropy
- Maximum Entropy Markov Model Framework
 - Model :
 - Exponential Model for Transition
 - Learning:
 - Generalized Iterative Scaling
 - **Prediction:**
 - State estimation from Observations
- Variations
- Experimental Results

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

26

EXPERIMENTAL RESULTS

February 23, 2010 CS 678I: Advanced Topics in Machine Learning

27