

Lecture 1, COMS E6998-3, Spring 2011

Michael Collins

January 19, 2011



A Machine-Learning Example: Hand-Written Digit Recognition

- ▶ The problem: given a hand-written digit, decide whether it is 0, 1, 2, ... or 9
- ▶ A learning approach:
 1. Collect several hundred/thousand example digits, and label them by hand to form a *training set*
 2. Automatically learn a digit recognition *model* from the training set
 3. Apply the model to new, previously unseen hand-written digits
- ▶ Systems built in this way are in widespread use in the U.S. postal service (ZIP-code recognition), and in automatic check-reading



Today's Lecture

- ▶ Introduction:
 - ▶ Example problems from machine learning for NLP
 - ▶ Topics we'll cover in the course
 - ▶ Background required for the course
 - ▶ Projects/homework assignments
- ▶ Topic 1: Hidden Markov models
- ▶ Topic 2: Log-linear models



Related Problems

- ▶ Identifying faces within an image (see the [Viola and Jones](#) face detector)
- ▶ Text classification/spam filtering
- ▶ Medical applications: e.g., classification of cancer type
- ▶ Information retrieval: e.g., ranking web-pages in order of relevance to a given query



Supervised Learning Problems

- ▶ Goal: Learn a function $f : \mathcal{X} \rightarrow \mathcal{Y}$

- ▶ We have n training examples

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

where each $x_i \in \mathcal{X}$, and each $y_i \in \mathcal{Y}$

- ▶ Often (not always) $\mathcal{X} = \mathbb{R}^d$ for some integer d

- ▶ Some possibilities for \mathcal{Y} :

- ▶ $\mathcal{Y} = \{-1, +1\}$ (binary classification)
- ▶ $\mathcal{Y} = \{1, 2, \dots, k\}$ for some $k > 2$ (multi-class classification)
- ▶ $\mathcal{Y} = \mathbb{R}$ (regression)

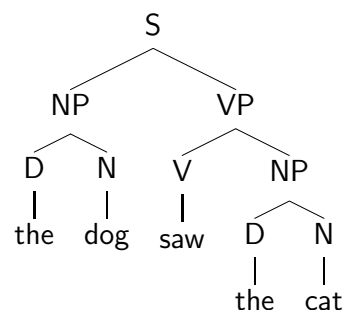


Context-Free Parsing

- ▶ The task: learn a function that maps a sentence, e.g.,

the dog saw the cat

to a parse tree,



Sequence Labeling Problems

- ▶ Task: learn a function that maps an input sequence

$$x_1, x_2, \dots, x_m$$

to an output sequence

$$y_1, y_2, \dots, y_m$$

Note: each $y_i \in \mathcal{Y}_i$ where \mathcal{Y}_i is a **finite** set of possible labels at the i 'th position

- ▶ This is a core problem in natural language processing
- ▶ Examples: part-of-speech tagging, named-entity recognition

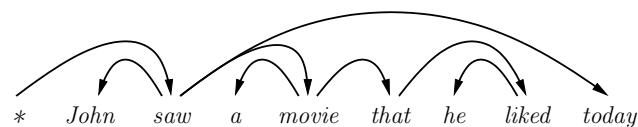


Dependency Parsing

- ▶ The task: learn a function that maps a sentence, e.g.,

John saw a movie that he liked today

to a *dependency structure*,



Machine Translation

- ▶ The task: learn a function that maps a sentence in one language, e.g.,

In wenigen Tagen finden Parlamentswahlen in Slovenian statt
to a sentence in another language,

In a few days elections take place in Slovenia



Topics Covered in the Class

- ▶ Probabilistic models for structured NLP data
 - ▶ e.g., hidden Markov models (HMMs), maximum-entropy Markov models (MEMMs), conditional random fields (CRFs), probabilistic context-free grammars, synchronous context-free grammars, dependency parsing models, etc.
- ▶ Inference algorithms
 - ▶ e.g., dynamic programming, belief propagation, methods based on linear programming and integer linear programming, dual decomposition/Lagrangian relaxation
- ▶ Semi-supervised learning
 - ▶ e.g., deriving lexical representations from unlabeled data, cotraining, entropy regularization, canonical correlation analysis (CCA)



Mapping Sentences to Logical Form

- ▶ The task: learn a function that maps a sentence e.g.,

Show me the latest flight from Boston to Seattle on Friday
to a expression in logical form that represents its meaning,
e.g.,

$$\operatorname{argmax}(\lambda x. \text{flight}(x) \wedge \text{from}(x, \text{BOS}) \wedge \text{to}(x, \text{SEA}) \wedge \text{day}(x, \text{FRI}), \lambda y. \text{time}(y))$$


Admin

- ▶ Background required for the class: a prior class in machine learning and/or natural language processing
- ▶ Evaluation:
 - ▶ Final class project (65%)
 - ▶ 3 homeworks (25%)
 - ▶ Class participation (10%)



Lecture 1, COMS E6998-3: Hidden Markov Models

Michael Collins

January 19, 2011



Markov Sequences

- ▶ Consider a sequence of random variables X_1, X_2, \dots, X_m where m is the length of the sequence
- ▶ Each variable X_i can take any value in $\{1, 2, \dots, k\}$
- ▶ How do we model the joint distribution

$$P(X_1 = x_1, X_2 = x_2, \dots, X_m = x_m)$$

?



Overview

- ▶ Markov models
- ▶ Hidden Markov models



The Markov Assumption

$$\begin{aligned} & P(X_1 = x_1, X_2 = x_2, \dots, X_m = x_m) \\ &= P(X_1 = x_1) \prod_{j=2}^m P(X_j = x_j | X_1 = x_1, \dots, X_{j-1} = x_{j-1}) \\ &= P(X_1 = x_1) \prod_{j=2}^m P(X_j = x_j | X_{j-1} = x_{j-1}) \end{aligned}$$

- ▶ The first equality is exact (by the chain rule).
- ▶ The second equality follows from *the Markov assumption*: for all $j = 2 \dots m$,

$$P(X_j = x_j | X_1 = x_1, \dots, X_{j-1} = x_{j-1}) = P(X_j = x_j | X_{j-1} = x_{j-1})$$



Homogeneous Markov Chains

- ▶ In a *homogeneous* Markov chain, we make an additional assumption, that for $j = 2 \dots m$,

$$P(X_j = x_j | X_{j-1} = x_{j-1}) = q(x_j | x_{j-1})$$

where $q(x' | x)$ is some function

- ▶ Idea behind this assumption: the transition probabilities do not depend on the position in the Markov chain (do not depend on the index j)



A Generative Story for Markov Models

- ▶ A sequence x_1, x_2, \dots, x_m is generated by the following process:

1. Pick x_1 at random from the distribution $q(x)$
2. For $j = 2 \dots m$:
 - ▶ Choose x_j at random from the distribution $q(x | x_{j-1})$



Markov Models

- ▶ Our model is then as follows:

$$p(x_1, x_2, \dots, x_m; \underline{\theta}) = q(x_1) \prod_{j=2}^m q(x_j | x_{j-1})$$

- ▶ Parameters in the model:

- ▶ $q(x)$ for $x = \{1, 2, \dots, k\}$
Constraints: $q(x) \geq 0$ and $\sum_{x=1}^k q(x) = 1$
- ▶ $q(x' | x)$ for $x = \{1, 2, \dots, k\}$ and $x' = \{1, 2, \dots, k\}$
Constraints: $q(x' | x) \geq 0$ and $\sum_{x'=1}^k q(x' | x) = 1$



Today's Lecture

- ▶ Markov models
- ▶ **Hidden Markov models**



Modeling Pairs of Sequences

- ▶ In many applications, we need to model *pairs* of sequences
- ▶ Examples:
 1. Part-of-speech tagging in natural language processing (assign each word in a sentence to one of the categories noun, verb, preposition etc.)
 2. Speech recognition (map acoustic sequences to sequences of words)
 3. Computational biology: recover gene boundaries in DNA sequences



Hidden Markov Models (HMMs)

- ▶ In HMMs, we assume that:

$$\begin{aligned} &P(X_1 = x_1, \dots, X_m = x_m, S_1 = s_1, \dots, S_m = s_m) \\ &= P(S_1 = s_1) \prod_{j=2}^m P(S_j = s_j | S_{j-1} = s_{j-1}) \prod_{j=1}^m P(X_j = x_j | S_j = s_j) \end{aligned}$$



Probabilistic Models for Sequence Pairs

- ▶ We have two sequences of random variables:
 X_1, X_2, \dots, X_m and S_1, S_2, \dots, S_m
- ▶ Intuitively, each X_i corresponds to an “observation” and each S_i corresponds to an underlying “state” that generated the observation. Assume that each S_i is in $\{1, 2, \dots, k\}$, and each X_i is in $\{1, 2, \dots, o\}$
- ▶ How do we model the joint distribution

$$P(X_1 = x_1, \dots, X_m = x_m, S_1 = s_1, \dots, S_m = s_m)$$

?



Independence Assumptions in HMMs

- ▶ By the chain rule, the following equality is exact:

$$\begin{aligned} &P(X_1 = x_1, \dots, X_m = x_m, S_1 = s_1, \dots, S_m = s_m) \\ &= P(S_1 = s_1, \dots, S_m = s_m) \times \\ &P(X_1 = x_1, \dots, X_m = x_m | S_1 = s_1, \dots, S_m = s_m) \end{aligned}$$

- ▶ Assumption 1: the state sequence forms a Markov chain

$$P(S_1 = s_1, \dots, S_m = s_m) = P(S_1 = s_1) \prod_{j=2}^m P(S_j = s_j | S_{j-1} = s_{j-1})$$



Independence Assumptions in HMMs

- ▶ By the chain rule, the following equality is exact:

$$\begin{aligned} & P(X_1 = x_1, \dots, X_m = x_m | S_1 = s_1, \dots, S_m = s_m) \\ = & \prod_{j=1}^m P(X_j = x_j | S_1 = s_1, \dots, S_m = s_m, X_1 = x_1, \dots, X_{j-1} = x_j) \end{aligned}$$

- ▶ Assumption 2: each observation depends only on the underlying state

$$\begin{aligned} & P(X_j = x_j | S_1 = s_1, \dots, S_m = s_m, X_1 = x_1, \dots, X_{j-1} = x_j) \\ = & P(X_j = x_j | S_j = s_j) \end{aligned}$$



A Generative Story for Hidden Markov Models

- ▶ Sequence pairs s_1, s_2, \dots, s_m and x_1, x_2, \dots, x_m are generated by the following process:

1. Pick s_1 at random from the distribution $t(s)$. Pick x_1 from the distribution $e(x|s_1)$
2. For $j = 2 \dots m$:
 - ▶ Choose s_j at random from the distribution $t(s|s_{j-1})$
 - ▶ Choose x_j at random from the distribution $e(x|s_j)$



The Model Form for HMMs

- ▶ The model takes the following form:

$$p(x_1 \dots x_m, s_1 \dots s_m; \theta) = t(s_1) \prod_{j=2}^m t(s_j | s_{j-1}) \prod_{j=1}^m e(x_j | s_j)$$

- ▶ Parameters in the model:

1. Initial state parameters $t(s)$ for $s \in \{1, 2, \dots, k\}$
2. Transition parameters $t(s'|s)$ for $s, s' \in \{1, 2, \dots, k\}$
3. Emission parameters $e(x|s)$ for $s \in \{1, 2, \dots, k\}$ and $x \in \{1, 2, \dots, o\}$



Today's Lecture

- ▶ More on Hidden Markov models:

- ▶ parameter estimation
- ▶ The Viterbi algorithm



Parameter Estimation with Fully Observed Data

- ▶ We'll now discuss parameter estimates in the case of *fully observed data*: for $i = 1 \dots n$, we have pairs of sequences $x_{i,j}$ for $j = 1 \dots m$ and $s_{i,j}$ for $j = 1 \dots m$. (i.e., we have n training examples, each of length m .)



Parameter Estimation: Emission Parameters

- ▶ Assume we have fully observed data: for $i = 1 \dots n$, we have pairs of sequences $x_{i,j}$ for $j = 1 \dots m$ and $s_{i,j}$ for $j = 1 \dots m$
- ▶ Define $\text{count}(i, s \rightsquigarrow x)$ to be the number of times state s is paired with emission x . More formally:

$$\text{count}(i, s \rightsquigarrow x) = \sum_{j=1}^m [[s_{i,j} = s \wedge x_{i,j} = x]]$$

- ▶ The maximum-likelihood estimates of emission probabilities are then

$$e(x|s) = \frac{\sum_{i=1}^n \text{count}(i, s \rightsquigarrow x)}{\sum_{i=1}^n \sum_x \text{count}(i, s \rightsquigarrow x)}$$



Parameter Estimation: Transition Parameters

- ▶ Assume we have fully observed data: for $i = 1 \dots n$, we have pairs of sequences $x_{i,j}$ for $j = 1 \dots m$ and $s_{i,j}$ for $j = 1 \dots m$
- ▶ Define $\text{count}(i, s \rightarrow s')$ to be the number of times state s' follows state s in the i 'th training example. More formally:

$$\text{count}(i, s \rightarrow s') = \sum_{j=1}^{m-1} [[s_{i,j} = s \wedge s_{i,j+1} = s']]$$

(We define $[[\pi]]$ to be 1 if π is true, 0 otherwise.)

- ▶ The maximum-likelihood estimates of transition probabilities are then

$$t(s'|s) = \frac{\sum_{i=1}^n \text{count}(i, s \rightarrow s')}{\sum_{i=1}^n \sum_{s'} \text{count}(i, s \rightarrow s')}$$



Parameter Estimation: Initial State Parameters

- ▶ Assume we have fully observed data: for $i = 1 \dots n$, we have pairs of sequences $x_{i,j}$ for $j = 1 \dots m$ and $s_{i,j}$ for $j = 1 \dots m$
- ▶ Define $\text{count}(i, s)$ to be 1 if state s is the initial state in the sequence, and 0 otherwise:

$$\text{count}(i, s) = [[s_{i,1} = s]]$$

- ▶ The maximum-likelihood estimates of initial state probabilities are:

$$t(s) = \frac{\sum_{i=1}^n \text{count}(i, s)}{n}$$



Today's Lecture

- ▶ Hidden Markov models:
 - ▶ parameter estimation
 - ▶ the Viterbi algorithm



The Viterbi Algorithm

- ▶ Goal: for a given input sequence x_1, \dots, x_m , find

$$\arg \max_{s_1, \dots, s_m} p(x_1 \dots x_m, s_1 \dots s_m; \underline{\theta})$$

- ▶ The *Viterbi algorithm* is a dynamic programming algorithm. Basic data structure:

$$\pi[j, s]$$

will be a table entry that stores the maximum probability for any state sequence ending in state s at position j . More formally: $\pi[1, s] = t(s)e(x_1|s)$, and for $j > 1$,

$$\pi[j, s] = \max_{s_1 \dots s_{j-1}} \left[t(s_1)e(x_1|s_1) \left(\prod_{k=2}^{j-1} t(s_k|s_{k-1})e(x_k|s_k) \right) t(s|s_{j-1})e(x_j|s) \right]$$



The Viterbi Algorithm

- ▶ Goal: for a given input sequence x_1, \dots, x_m , find

$$\arg \max_{s_1, \dots, s_m} p(x_1 \dots x_m, s_1 \dots s_m; \underline{\theta})$$

- ▶ This is the most likely state sequence $s_1 \dots s_m$ for the given input sequence $x_1 \dots x_m$



The Viterbi Algorithm

- ▶ Initialization: for $s = 1 \dots k$

$$\pi[1, s] = t(s)e(x_1|s)$$

- ▶ For $j = 2 \dots m$, $s = 1 \dots k$:

$$\pi[j, s] = \max_{s' \in \{1 \dots k\}} [\pi[j-1, s'] \times t(s|s') \times e(x_j|s)]$$

- ▶ We then have

$$\max_{s_1 \dots s_m} p(x_1 \dots x_m, s_1 \dots s_m; \underline{\theta}) = \max_s \pi[m, s]$$

- ▶ The algorithm runs in $O(mk^2)$ time



