Sequence Learning Statistical Methods in NLP? ISCL-BA-08

Çağrı Çöltekin ccoltekin@sfs.uni-tuebingen.de

Summer Semester 2025

Some (typical) machine learning applications Spam detection

Credit scoring

x (input) y (output) document spam or not Sentiment analysis product review sentiment Medical diagnosis patient data diagnosis financial history loan decision

The cases (input-output) pairs are assumed to be independent and identically distributed (i.i.d.).

Structured prediction

In many applications, the i.i.d. assumption is wrong

x (input) POS tagging word sequence Parsing word sequence OCR image (array of pixels) Gene prediction

POS sequence parse tree sequences of letters

y (output)

Structured/sequence learning is prevalent in NLP

Sequence learning - a demonstration of the problem





Recap: (conditional) independence





- · Individual predictions depend on each other
- . Can we treat the whole sequence as a single label?

Recap: chain rule

We rewrite the relation between the joint and the conditional probability as

P(X,Y) = P(X|Y)P(Y)

 $P(X,Y) = P(Y \mid X)P(X)$

In general, for any number of random variables, we can write

 $P(X_1, X_2, ..., X_n) = P(X_1 | X_2, ..., X_n)P(X_2, ..., X_n)$

If two variables X and Y are independent,

 $P(X \,|\, Y) = P(X) \quad and \quad P(X,Y) = P(X)P(Y)$

If two variables X and Y are independent given another variable Z, P(X|Y|Z) = P(X|Z)P(Y|Z)

An example: probability of a sentence

P(The old man the boats) = ?

- We cannot estimate this probability by counting all occurrences of the sentence, and dividing it to the total number of sentences in English
 We can (potentially) calculate its probability based on the probabilities of the words. Using chain rule
- $P(S) = P(boats \mid The \ old \ man \ the) P(The \ old \ man \ the)$
 - P(boats | The old man the)P(the | The old man)P(The old)
 - $= P(boats \mid The \ old \ man \ the) P(the \mid The \ old \ man) P(man \mid The \ old) P(old \mid The) P(The)$
- . Did we solve the problem of probability estimation?

Markov chains Given a sequence of events (or states), $q_1, q_2, \dots q_t$

 \star In a first-order Markov chain, the probability of an event q_t is

 $P(q_t|q_1,...,q_{t-1}) = P(q_t|q_{t-1})$

* In higher order chains, the dependence of history is extended, e.g., second-order Markov chain

 $P(q_t|q_t,\dots,q_{t-1}) = P(q_t|q_{t-2},q_{t-1})$

The conditional independence properties simplify the probability distributions.

Markov chains

A Markov model is defined by.

- * A set of states $Q = \{q_1, \dots, q_n\}$
- A special start state qo A transition probability matrix

 - $\mathbf{A} = \begin{bmatrix} a_{01} & a_{02} & \dots & a_{0n} \\ a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ \end{bmatrix}$ a_, a_, ... a_

where an is the probability of transition from state i to state i

Back to sentence probability example

- · With a first-order Markov assumption,
 - $P(S) = P(boats \,|\, The \,old \,man \,the) P(The \,old \,man \,the)$ = P(boats | the)P(the | man)P(man | old)P(old | The)P(The)
- · Now the probabilities are easier to estimate
- * The above approach is an example of n-gram language models that we will get back to very soon

Hidden/latent variables

- In many machine learning problems we want to account for unobserved/unobservable latent or hidden variables
- Some examples

 - 'personality' in many psychological data
 'topic' of a text
 - text tomic class' of a speaker
- Accounting for latent variables improve the accuracy of the models Since we cannot observe them, latent variables make learning algorithms difficult
- Learning with hidden variables
 - The EM algorithm (or its variants) is used in many machine learning models with latent/hidden variables 1. Randomly initialize the parameters
 - Iterate until convergence:
 - E-stop compute likelihood of the data, given the parameters
 M-step re-estimate the parameters using the predictions based on the E-step

Hidden Markov models (HMM)

HMMs are like Markov chains: probability of a state depends only a limited history of previous states

$$P(q_t|q_1,\ldots,q_{t-1}) = P(q_t|q_{t-1})$$

- . Unlike Markov chains, state sequence is hidden, they are not the observations
- + At every state q_t , an HMM emits an output, o_t , whose probability depends only on the associated hidden state
- \ast Given a state sequence $q=q_1,\dots,q_T,$ and the corresponding observation sequence $o = o_1, ..., o_T$,

$$P(\mathbf{o}, \mathbf{q}) = p(q_1) \left[\prod_{2}^{T} P(q_t | q_{t-1}) \right] \prod_{1}^{T} P(o_t | q_t)$$

HMMs: formal definition An HMM is defined by

- - A set of states $Q = \{q_1, \dots, q_n\}$ * The set of possible observations $O = \{o_1, \dots, o_m\}$
 - · A transition probability matrix
 - [a11 a12 ... a1n] a_{ij} is the probability of from state q_i to state q_j
 - Initial probability distribut $\pi = \{P(q_1), ..., P(q_n)\}\$
 - · Probability distributions of [b₁₁ b₁₂ ... b_{1n}]
 - b_{ij} is the probability of emitting output σ_i at state q_j
 - b_{m1} b_{m2} ... b_{mn}

HMM transition diagram



HMMs: three problems

P(a | o:M)

P(o | M)

Given observation sequences, a set of states, and (sometimes) corresponding state sequences, estimate the parameters (π , A, B) of

Calculating likelihood of a given sequence

Assigning probabilities to observation sequences

 Start with calculating all forward probabilities for t = 1 $\alpha_{1,i} = \pi_i P(o_1|q_i)$ for $1 \leqslant i \leqslant |Q|$

store the α values For t > 1.

 $\alpha_{t,t} = \sum \alpha_{t-1,j} P(q_t|q_j) P(o_t|q_t) \quad \text{for } 1\leqslant t\leqslant |Q|, 2\leqslant t\leqslant n$

. Likelihood of the observation is the sum of the forward probabilities of the

$$P(o|M) = \sum_{j=1}^{|Q|} \alpha_{n,j}$$

Determining best sequence of latent variables

- + We often want to know the hidden state sequence given an observation sequence, $P(\mathbf{q}\mid\mathbf{o};M)$ For example, given a sequence of tokens, find the most likely POS tag sequence.
- The problem (also the solution, the Viterbi algorithm) is very similar to the forward algorithm

- Two major differences

 we store maximum likelihood leading to each node on the lattice

 we also store backlinks, the previous state that leads to the maxim

Example: HMMs for POS tagging



- Probability of a tag depends on the previous tag
- . Probability of a word at a given state depends only on the current tag

. Four possible observations: a, b, c , d

A simple example

. Three states: N, V, D

Unfolding the states



Assigning probabilities to observation sequences

$$P(o \mid M) = \sum_{\mathbf{q}} P(o, \mathbf{q} \mid M)$$

We need to sum over an exponential number of hidden state sequer The solution is using a dynamic programming algorithm
 for each node of the lattice, store forward probabilities

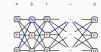
$$\alpha_{t,i} = \sum_{j}^{N} \alpha_{t-1,j} P(q_i|q_j) P(o_i|q_i)$$

Forward algorithm



 $\alpha_{2,2}=\alpha_{1,1}\alpha_{\mathrm{NV}}b_{\mathrm{yV}}+\alpha_{1,2}\alpha_{\mathrm{VV}}b_{\mathrm{yV}}+\alpha_{1,3}\alpha_{\mathrm{DV}}b_{\mathrm{yV}}$

HMM decoding problem



Learning the parameters of an HMM . We want to estimate π. A. B . If we have both the observ sequence, MLE estimate is
$$\begin{split} \pi_t &= \frac{C(q_0 \rightarrow q_1)}{\sum_k C(q_0 \rightarrow q_k)} \\ \alpha_{tj} &= \frac{C(q_1 \rightarrow q_j)}{\sum_k C(q_1 \rightarrow q_k)} \\ b_{ti} &= \frac{C(q_1 \rightarrow q_j)}{\sum_k C(q_1 \rightarrow q_k)} \end{split}$$

s, e.g., p(glo) can be a non

+ The HMMs we discussed so far are called ergodic HMMs: all α_{ij} are non-zero For some applications, it is common to use HMMs with additional restrictions A well known variant (Bakis HMM) allows only forward transitions

· Graphical models define models involving multiple random variables

. It is generally more intuitive (compared to corresponding mathematical equations) to work with graphical models · In a graphical model, by convention, the observed variables are shaded Graphs can also be undirected, which are also called Markov random fields.

HMM variations

dietribution

Graphical models

MaxEnt HMMs (MEMM)

Conditional random fields

+ In HMMs, we model $P(q,\sigma) = P(q)P(\sigma\,|\,q)$ * In many applications, we are only interested in $P(q \mid o)$, which we can calculate using the Bayes theorem \bullet But we can also model $P(q \mid o)$ directly using a maximum entropy model $P(q_t \mid q_{t-1}, o_t) = \frac{1}{Z} e^{\sum w_t f_t(o_t, q_t)}$

 f_{\downarrow} are features – can be any useful feature Z normalizes the probability distribution

Learning the parameters of an HMM

 Given a training set with of want to find θ = (π, A, B) observation sequence(s) o and state sequence q, we

 $\underset{\mathbf{q}}{\operatorname{arg\,max}} P(\mathbf{o} \mid \mathbf{q}, \boldsymbol{\theta})$

Typically solved using EM

sypto.amy Servetu tought Essi

1. Initialize 0

2. Repeat until convergence

E-step given 0, estimate the hidden state sequence

M-step given the estimated hidden states, use 'expected counts' to update 0

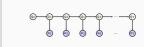
An efficient implementation of EM algorithm is called Baum-Welch algorithm or forward-backward algorithm

Directed graphical models: a brief divergence

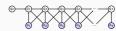
We saw earlier that joint distributions of multiple random variables factorized different ways

 Graphical models display this relations in graphs, variables are denoted by nodes, - the dependence between the var

HMM as a graphical model



MEMMs as graphical models



odel the joint distribu you can generate the output using HMMs
 MEMMs and CRFs are discriminative models they model the conditional

Generative vs. discriminative models

HMMs are generative models, they m

Additional reading, references, credits

. A related model used in NLP is conditional random field (CRF) . CRFs are undirected models

• CRFs also model P(q | o) directly

 $P(q \mid \mathbf{o}) = \frac{1}{Z} \prod f(q_{t-1}, q_t) g(q_t, o_t)$

Summary

- In many problems, e.g., POS tagging, i.i.d. assumption is wrong
- We need models that are aware of the effects of the sequence (or structure in general) in the data · HMMs are generative sequence models:
- Markov assumption between the hidden states (POS tags)
 Observations (words) are conditioned on the state (tag) · There are other sequence learning methods
- Briefly mentioned: MEMM, CRF
 Coming soon: recurrent neural n
- Reading: Jurafsky and Martin (2025, Chapter 17)
- Next · Recurrent and convolutional networks
- * Reading: Jurafsky and Martin (2025, Chapter 8)

It is easier to add arbitrary features on discriminative models

. In general: HMMs work well when the state sequence, P(q), can be modeled