

	x (input)	y (output)
Spam detection	document	spam or not
Sentiment analysis	product review	sentiment
Medical diagnosis	patient data	diagnosis
Credit scoring	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)	y (output)
POS tagging	word sequence	POS sequence
Parsing	word sequence	parse tree
OCR	image (array of pixels)	sequences of letters
Gene prediction	genome	genes

Structured/sequence learning is prevalent in NLP.

## In this lecture ...

- Hidden Markov models (HMMs)
- A short note on graphical probabilistic models
- Alternatives to HMMs (briefly): HMEM / CRF

... and soon

- Recurrent neural networks

## Recap: chain rule

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

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

We can also write the same quantity as,

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

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

$$P(X_1, X_2, \dots, X_n) = P(X_1 | X_2, \dots, X_n)P(X_2, \dots, X_n)$$

## Recap: (conditional) independence

If two variables X and Y are independent,

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

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

$$P(\text{It's a beautiful day}) = ?$$

- We cannot just count all occurrences of the sentence, and divide it to the total number of sentences in English
- But we can calculate its probability based on the probabilities of the words. Using chain rule

$$\begin{aligned} P(\text{It's a beautiful day}) &= P(\text{day} | \text{It's a beautiful})P(\text{It's a beautiful}) \\ &= P(\text{day} | \text{It's a beautiful})P(\text{beautiful} | \text{It's a})P(\text{It's a}) \\ &= P(\text{day} | \text{It's a beautiful})P(\text{beautiful} | \text{It's a})P(a | \text{It's})P(\text{It's}) \end{aligned}$$

- Did we solve the problem?

## Markov chains

calculating probabilities

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

- In a *first-order Markov chain*, the probability of an event  $q_t$  is

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

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

$$P(q_1 | q_1, \dots, q_{t-1}) = P(q_1 | q_{t-2}, q_{t-1})$$

- The conditional independence properties simplify the probability distributions

## Markov chains

definition

A Markov model is defined by,

- A set of states  $Q = \{q_1, \dots, q_n\}$
- A special start state  $q_0$
- A transition probability matrix

$$A = \begin{bmatrix} a_{01} & a_{02} & \dots & a_{0n} \\ a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

where  $a_{ij}$  is the probability of transition from state  $i$  to state  $j$

## Back to sentence probability example

- With a first-order Markov assumption,

$$\begin{aligned} P(\text{It's a beautiful day}) &= P(\text{day} | \text{It's a beautiful})P(\text{beautiful} | \text{It's a})P(a | \text{It's})P(\text{It's}) \\ &= P(\text{day} | \text{beautiful})P(\text{beautiful} | a)P(a | \text{It's})P(\text{It's} | (S)) \end{aligned}$$

- Now the probabilities are easier to calculate
- 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
  - “socio-economic class” of a speaker
- Latent variables make learning difficult: since we cannot observe them, how do we set the parameters?

## Learning with hidden variables

(Another) informal/quick introduction to the EM algorithm

- The EM algorithm (or its variants) is used in many machine learning models with latent/hidden variables

1. Randomly initialize the parameters
2. Iterate until convergence:
  - E-step: compute likelihood of the data, given the parameters
  - M-step: re-estimate the parameters using the predictions based on the E-step



## Learning the parameters of an HMM

supervised case

- We want to estimate  $\pi, A, B$
- If we have both the observation sequence  $\mathbf{o}$  and the corresponding state sequence, MLE estimate is

$$\pi_i = \frac{C(q_0 \rightarrow q_i)}{\sum_k C(q_0 \rightarrow q_k)}$$
$$a_{ij} = \frac{C(q_i \rightarrow q_j)}{\sum_k C(q_i \rightarrow q_k)}$$
$$b_{ij} = \frac{C(q_i \rightarrow o_j)}{\sum_k C(q_i \rightarrow o_k)}$$

## HMM variations

- The HMMs we discussed so far are called *ergodic* HMMs: all  $a_{ij}$  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



- The emission probabilities can also be continuous, e.g.,  $p(q|o)$  can be a normal distribution

## Graphical models

- 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*

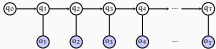
## MaxEnt HMMs (MEMM)

- In HMMs, we model  $P(\mathbf{q}, \mathbf{o}) = P(\mathbf{q})P(\mathbf{o} | \mathbf{q})$
- In many applications, we are only interested in  $P(\mathbf{q} | \mathbf{o})$ , which we can calculate using the Bayes theorem
- But we can also model  $P(\mathbf{q} | \mathbf{o})$  directly using a *maximum entropy model*

$$P(q_i | q_{i-1}, o_i) = \frac{1}{Z} e^{\sum w_j f_j(o_i, q_i)}$$

$f_i$  are features – can be any useful feature  
 $Z$  normalizes the probability distribution

## Conditional random fields



- A related model used in NLP is *conditional random field* (CRF)
- CRFs are *undirected* models
- CRFs also model  $P(\mathbf{q} | \mathbf{o})$  directly

$$P(\mathbf{q} | \mathbf{o}) = \frac{1}{Z} \prod_i f(q_{i-1}, q_i | g(q_i, o_i))$$

## Summary

- In many problems, e.g., POS tagging, I.L.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 networks

Next

- Recurrent and convolutional networks

## Learning the parameters of an HMM

- Given a training set with observation sequence(s)  $\mathbf{o}$  and state sequence  $\mathbf{q}$ , we want to find  $\theta = (\pi, A, B)$

$$\arg \max_{\theta} P(\mathbf{o} | \mathbf{q}, \theta)$$

- Typically solved using EM
  1. Initialize  $\theta$
  2. Repeat until convergence
    - E-step: given  $\theta$ , estimate the hidden state sequence
    - M-step: given the estimated hidden states, use 'expected counts' to update  $\theta$
- An efficient implementation of EM algorithm is called *Baum-Welch algorithm*, or *forward-backward algorithm*

## Directed graphical models: a brief divergence

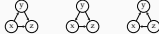
### Bayesian networks

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

$$P(\mathbf{x}, \mathbf{y}, \mathbf{z}) = P(\mathbf{x})P(\mathbf{y} | \mathbf{x})P(\mathbf{z} | \mathbf{x}, \mathbf{y}) = P(\mathbf{y})P(\mathbf{x} | \mathbf{y})P(\mathbf{z} | \mathbf{x}, \mathbf{y}) = P(\mathbf{z})P(\mathbf{x} | \mathbf{z})P(\mathbf{y} | \mathbf{x}, \mathbf{z})$$

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

- Bayesian networks are directed acyclic graphs



- A variable (node) depends only on its parents

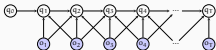
## HMM as a graphical model



## MEMMs as graphical models



We can also have other dependencies as features, for example



## Generative vs. discriminative models

- HMMs are *generative* models, they model the joint distribution
  - you can *generate* the output using HMMs
- MEMMs and CRFs are *discriminative* models they model the conditional probability directly
- It is easier to add arbitrary features on discriminative models
- In general: HMMs work well when the state sequence,  $P(\mathbf{q})$ , can be modeled well

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1

C. Gillekin

SR | University of Tübingen

Summer Semester 2021

AS2

Week 1