COS 484

Sequence Models - 2

Spring 2022



- Lecture slides:
 - uploaded previous evening/night
 - available on class website)
 - download from website)
- message

Logistics

• if you want a head start, you can view slides from Spring 2021 (link

• poll slides may be skipped, updated slides available after class (just re-

• Office hours: if you're unable to make any of the office hours, send us a





- 1. Set of states $S = \{1, 2, ..., K\}$ and set of observations O
- 2. Initial state probability distribution $\pi(s_1)$
- 3. Transition probabilities $P(s_{t+1} | s_t)$ (OR $\theta_{s_t \rightarrow s_{t+1}}$)
- 4. Emission probabilities $P(o_t | s_t)$ (OR $\phi_{s_t \rightarrow o_t}$)

Recap: Hidden Markov Models

Strong assumptions





Markov assumption: 1.

 $P(s_{t+1} | s_1, \dots, s_t) \approx P(s_{t+1} | s_t)$

Output independence: 2.

 $P(o_t | s_1, \ldots, s_t) \approx P(o_t | s_t)$

1) assumes state sequences do not have very strong priors/longrange dependencies 2) assumes neighboring states don't affect current **o**bservation

Recap:Viterbi decoding

.

.

k



M[i, j] stores joint probability of most probable sequence of states ending with state j at time i

$M[i,j] = \max M[i-1,k] P(s_i | s_k) P(o_i | s_j) \quad 1 \le k \le K \quad 1 \le i \le n$

Pick max M[n, k] and backtrack using B





Maximum Entropy Markov Models

Generative vs Discriminative

• HMM is a generative model

• Can we model $P(s_1, \ldots, s_n | o_1, \ldots, o_n)$ directly?

Generative

Text classification

Naive Bayes: P(c)P(d | c)

Sequence prediction

HMM:

 $P(s_1,\ldots,s_n)P(o_1,\ldots,o_n \mid s_1,\ldots,s_n)$

Discriminative

Logistic Regression: $P(c \mid d)$

MEMM: $P(s_1, ..., s_n | o_1, ..., o_n)$

(No factorization)



• Compute the posterior directly:

•
$$P(S|O) \approx \arg\max_{S} \prod_{i} P(s_i|o_i, s_{i-1})$$

• Use features and weights: $P(s_i = s | o_i, s_{i-1}) \propto \exp(\mathbf{w} \cdot \mathbf{f}(s_i = s, o_i, s_{i-1}))$

Maximum Entropy Markov Model



 $O = \langle o_1, o_2, \dots, o_n \rangle$

Fratures

No factorization into transition, emission

~ weights

(Bigram MEMM)

- Use features and weights: $P(s_i = s \mid a_i)$
- Which of the following is the correct way to calculate this probability? A) $P(s_i = s \mid o_i, s_{i-1}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(s_i = s, o_i, s_{i-1}))}{\sum_{s' \in S} \exp(\mathbf{w} \cdot \mathbf{f}(s_i = s, o_i, s_{i-1} = s'))}$ $\mathsf{B}) P(s_i = s \mid o_i, s_{i-1}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(s_i = s, o_i, s_{i-1}))}{\sum_{s' \in S} \exp(\mathbf{w} \cdot \mathbf{f}(s_i = s', o_i, s_{i-1}))}$ C) $P(s_i = s \mid o_i, s_{i-1}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(s_i = s, o_i, s_{i-1}))}{\sum_{o' \in O} \exp(\mathbf{w} \cdot \mathbf{f}(s_i = s, o_i = o', s_{i-1}))}$

$$o_i, s_{i-1}) \propto \exp(\mathbf{w} \cdot \mathbf{f}(s_i = s, o_i, s_{i-1}))$$





• Compute the posterior directly:

•
$$P(S|O) \approx \arg\max_{S} \prod_{i} P(s_i|o_i, s_{i-1})$$

• Use features and weights: $P(s_i = s | o_i, s_{i-1})$

Maximum Entropy Markov Model



 $O = \langle o_1, o_2, \dots, o_n \rangle$

$$\sum_{i=1}^{n} = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(s_i = s, o_i, s_{i-1}))}{\sum_{s' \in S} \exp(\mathbf{w} \cdot \mathbf{f}(s_i = s', o_i, s_{i-1}))}$$

(Bigram MEMM)





• In general, we can use all observations and all previous states:

$$P(S \mid O) = \arg \max_{S} \prod_{i} P(s_i \mid o_n, o_{i-1}, \dots, o_1, s_{i-1}, \dots, s_1)$$

 $P(s_i | s_{i-1}, \ldots, s_1, O)$

MEMM



$$O) \propto \exp(\mathbf{w} \cdot \mathbf{f}(s_i, s_{i-1}, \dots, s_1, O))$$

Q: Why couldn't we do this with HMMs?



$$\langle t_i, w_{i-2} \rangle, \langle t_i, w_{i-1} \rangle, \langle t_i, w_i \rangle, \langle t_i, w_{i+1} \rangle, \langle t_i, w_{i+2} \rangle \langle t_i, t_{i-1} \rangle, \langle t_i, t_{i-2}, t_{i-1} \rangle, \langle t_i, t_{i-1}, w_i \rangle, \langle t_i, w_{i-1}, w_i \rangle \langle t_i, w_i, w_{i+1} \rangle,$$

Feature templates

Features in an MEMM

 $t_i = VB$ and $w_{i-2} = Janet$ $t_i = VB$ and $w_{i-1} = will$ $t_i = VB$ and $w_i = back$ $t_i = VB$ and $w_{i+1} = the$ t_i = VB and w_{i+2} = bill $t_i = VB$ and $t_{i-1} = MD$ t_i = VB and t_{i-1} = MD and t_{i-2} = NNP $t_i = VB$ and $w_i = back$ and $w_{i+1} = the$

Features (binary)

 $t_i = tags (states)$ w_i = words (observations)



DT JJ NN Incorrect

DT NN Correct

The old

Which of these feature templates would help most to tag 'old' correctly? A) $\langle t_i, w_i, w_{i-2}, w_{i-1}, w_{i+1} \rangle$ B) $\langle t_i, w_i, w_{i-1} \rangle$ C) $\langle t_i, w_i, w_{i-1}, w_{i+1} \rangle$ D) $\langle t_i, w_i, w_{i-1}, w_{i+1}, w_{i+2} \rangle$

Features in an MEMM

- DT NN
- VB DT NN
- the man boat

t = tagsw = words





$$\hat{S} = \arg\max_{S} P(S \mid O) = \arg\max_{S} P(S \mid O)$$

(assume bigram MEMM, i.e. features only on previous time step and current obs)

• Greedy decoding:



 $\operatorname{rgmax}_{S} \Pi_{i} P(s_{i} \mid o_{i}, s_{i-1})$

s = argmax P(s|The)

$$\hat{S} = \arg\max_{S} P(S \mid O) = \arg$$

• Greedy decoding:



 $g \max_{G} \prod_{i} P(s_i \mid o_i, s_{i-1})$ S

 $S_2 = \alpha Agmax P(S|cat, DT)$

= NN

$$\hat{S} = \arg\max_{S} P(S \mid O) = \arg$$

• Greedy decoding:



 $g\max_{S} \prod_{i} P(s_i \mid o_i, s_{i-1})$ S

$$\hat{S} = \arg\max_{S} P(S \mid O)$$

Greedy decoding

• Viterbi decoding:

k \bigwedge DP Lattice (Best sequence ending in s_i)

 $) = \arg\max_{S} \prod_{i} P(s_i \mid o_i, s_{i-1})$

 $M[i,j] = \max_{i} M[i-1,k] P(s_j | o_i, s_k) \quad 1 \le k \le K \quad 1 \le i \le n$ # states # timesteps

(or equivalent log form)

Viterbi decoding for MEMMs

.

k



M[i, j] stores joint probability of most probable sequence of states ending with state j at time i

Pick max M[n, k] and backtrack using B





How would you compare the computational complexity of Viterbi decoding for bigram MEMMs compared to decoding for bigram HMMs?

- A) More operations in MEMM
- B) More operations in HMM
- C) Equal

D) Depends on number of features in MEMM

$$M[i,j] = \max_{k} M[i-1,k] P$$

$$M[i,j] = \max_{k} M[i-1,k] P$$





MEMM: Learning

• Gradient descent: similar to logistic regression!

$$P(s_i = s \mid s_1, \dots, s_{i-1}, O) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(s_1, \dots, s_{i-1}, s_i = s, O))}{\sum_{s'} \exp(\mathbf{w} \cdot \mathbf{f}(s_1, \dots, s_{i-1}, s_i = s', O))}$$

• Given: annotated pairs of (

Loss for one sequence, L

$$(S, O)$$
 where each $S = \langle s_1, s_2, \ldots, s_n \rangle$

$$= -\sum_{i=1}^{n} \log P(s_i | s_1, \dots, s_{i-1}, O)$$

• Compute gradients with respect to weights w and update

MEMM vs HMM

- HMM models the joint P(S, O) while MEMM models the required prediction $P(S \mid O)$
- MEMM has more expressivity
 - observation sequence
 - allows for more flexible features
- HMM may hold an advantage if the dataset is small

• accounts for dependencies between neighboring states and entire

Label bias



$$P(JJ \mid DT) P(\text{old} \mid JJ) P(N)$$

$$P(NN \mid DT) P(\text{old} \mid NN) P(N)$$

Low entropy transitions (e.g. JJ -> NN) between labels may override the effect of observations

The/? old/? man/? the/? boat/?



Stanford Parser

Please enter a sentence to be parsed:



Your query

The old man the boat

Tagging

The/DT old/JJ man/NN the/DT boat/NN

	/
	4.5
lantanaa	
Senience	Parce
	r al se

Solution? **Conditional Random Fields**



Conditional Random Fields

Conditional Random Field

- Model $P(s_1, \ldots, s_n | o_1, \ldots, o_n)$ directly
- No Markov assumption
 - global feature vector
- Normalize over entire sequences

Map entire sequence of states S and observations O to a

Conditional Random Field



 $P(S \mid O) = -\frac{1}{\Sigma}$

 $\mathbf{O}(\mathbf{C} \otimes \mathbf{O})$

 $\exp(\mathbf{w} \cdot \mathbf{f}(S, O))$

Z(O)

$$\exp(\mathbf{w} \cdot \mathbf{f}(S, O))$$
$$\mathbf{E}_{S'}\exp(\mathbf{w} \cdot \mathbf{f}(S', O))$$

$\exp(\mathbf{w} \cdot \mathbf{f}(S, O))$ Z(O)

$P(S \mid O) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(S, O))}{\sum_{S'} \exp(\mathbf{w} \cdot \mathbf{f}(S', O))}$

Features

- Each F_k in **f** is a **global** feature function
- Can be computed as a combination of local features (linear chain CRF)

•
$$F_k = \sum_{i=1}^n f_k(s_{i-1}, s_i, O, i)$$

 each local feature only depends on previous and current states (outputs)

$1{x_i = the, y_i = DET}$ $1{y_i = \text{PROPN}, x_{i+1} = \text{Street}, y_{i-1} = \text{NUM}}$ $\mathbb{1}\{y_i = \text{VERB}, y_{i-1} = \text{AUX}\}$

Features

- Each F_k in **f** is a **global** feature function
- Can be computed as a combination of local features (linear chain CRF)

•
$$F_k = \sum_{i=1}^n f_k(s_{i-1}, s_i, O, i)$$

• Training: Optimize weights using supervised data similar to logistic regression

Inference

• $\hat{S} = \arg\max_{S} P(S|O) = \arg\max_{S} \frac{\exp(\mathbf{w} \cdot \mathbf{f}(S,O))}{Z(O)}$

S

 $= \arg \max$ S

(for linear chain CRF)

• Use Viterbi similar to HMM and MEMM

 $= \arg\max_{\mathbf{c}} \exp(\mathbf{w} \cdot \mathbf{f}(S, O))$

$$\exp(\sum_{k=1}^{K}\sum_{i=1}^{n}w_{k}f_{k}(s_{i-1}, s_{i}, O, i))$$

CRF vs MEMM

- MEMM models the required prediction $P(S \mid O)$ using the Markov assumption, while the CRF does not
- CRF uses global features while MEMM features are localized
- Feature design is flexible in both models
- CRF is computationally more complex

History of CRFs

Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data

John Lafferty^{†*} LAFFERTY@CS.CMU.EDU Andrew McCallum^{*†} MCCALLUM@WHIZBANG.COM Fernando Pereira^{*‡} FPEREIRA@WHIZBANG.COM *WhizBang! Labs–Research, 4616 Henry Street, Pittsburgh, PA 15213 USA [†]School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213 USA [‡]Department of Computer and Information Science, University of Pennsylvania, Philadelphia, PA 19104 USA

- Lafferty, McCallum, Pereria (2001): introduced CRFs for sequence modeling
 - Mitigates the label bias problem (in HMMs/MEMMs)
 - Better empirical performance compared to HMMs/MEMMs
 - Parameter estimation not straightforward

History of CRFs

- Very popular in the 2000s
- Wide variety of applications:
 - Information extraction
 - Summarization
 - Image labeling/segmentation

Information extraction from research papers using conditional random fields 🖈

Fuchun Peng ^a [∧] [∞], Andrew McCallum ^b [∞]

Multiscale conditional random fields for image labeling

Publisher: IEEE

Cite This

🔎 PDF

Xuming He; R.S. Zemel; M.A. Carreira-Perpinan All Authors

Document Summarization using Conditional Random Fields

Dou Shen¹, Jian-Tao Sun², Hua Li², Qiang Yang¹, Zheng Chen² ¹Department of Computer Science and Engineering Hong Kong University of Science and Technology, Hong Kong {dshen, qyang}@cse.ust.hk ²Microsoft Research Asia, 49 Zhichun Road, China {jtsun, huli, zhengc}@microsoft.com

History of CRFs

- Very popular in the 2000s
- Wide variety of applications:
 - Information extraction
 - Summarization
 - Image labeling/segmentation

Software [edit]

This is a partial list of software that implement generic CRF tools.

- RNNSharp
 CRFs based on recurrent neural networks (C#, .NET)
- CRF-ADF & Linear-chain CRFs with fast online ADF training (C#, .NET)
- CRFSharp & Linear-chain CRFs (C#, .NET)
- GCO I CRFs with submodular energy functions (C++, Matlab)
- DGM I General CRFs (C++)
- GRMM & General CRFs (Java)
- CRFall & General CRFs (Matlab)
- Sarawagi's CRF I Linear-chain CRFs (Java)
- Accord.NET & Linear-chain CRF, HCRF and HMMs (C#, .NET)

- FlexCRFs & First-order and second-order Markov CRFs (C++)
- crf-chain1 & First-order, linear-chain CRFs (Haskell)
- imageCRF
 Graphic CRF for segmenting images and image volumes (C++)
- MALLET & Linear-chain for sequence tagging (Java)

Empirical performance

Model	F score
SVM combination	94.39%
(Kudo and Matsumoto, 2001)	
CRF	94.38%
Generalized winnow	93.89%
(Zhang et al., 2002)	
Voted perceptron	94.09%
MEMM	93.70%

Table 2: NP chunking F scores

null hypothesis	p-value
CRF vs. SVM	0.469
CRF vs. MEMM	0.00109
CRF vs. voted perceptron	0.116
MEMM vs. voted perceptron	0.0734

Table 4: McNemar's tests on labeling disagreements

Sha and Pereira (2003)

CRFs in deep learning era

Conditional Random Fields as Recurrent **Neural Networks**

Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, Philip H. S. Torr; Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2015, pp. 1529-1537

Neural Architectures for Named Entity Recognition

Guillaume Lample Miguel Ballesteros Sandeep Subramanian[♠] Kazuya Kawakami[♠] Chris Dyer[♠] Carnegie Mellon University
NLP Group, Pompeu Fabra University {glample, sandeeps, kkawakam, cdyer}@cs.cmu.edu, miguel.ballesteros@upf.edu

Bidirectional LSTM-CRF Models for Sequence Tagging

Zhiheng Huang Wei Xu Baidu research Baidu research huangzhiheng@baidu.com xuwei06@baidu.com

Kai Yu Baidu research yukai@baidu.com

- Use CRFs on top of neural representations (instead of features and weights)
- Joint sequence prediction without the need for defining features!
- Recent architectures such as seq2seq w/ attention or Transformer may implicitly do the job

Preview: Recurrent neural networks (RNNs)

How can we model sequences using **neural networks**?

- Recurrent neural networks = A class of neural networks used to model sequences, allowing to handle variable length inputs
- variable-length, sequential inputs

• Very crucial in NLP problems (different from images) because sentences/paragraphs are

Preview: Recurrent neural networks (RNNs)

A family of neural networks allowing to handle variable length inputs

A function: $y = \text{RNN}(\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n) \in \mathbb{R}^h$ where $\mathbf{x}_1, ..., \mathbf{x}_n \in \mathbb{R}^d$

Preview: Recurrent neural networks (RNNs)

Proven to be an highly effective approach to language modeling, sequence tagging as well as text classification tasks:

