# Maximum entropy sequence models

Maximum entropy Markov models (MEMMs) or Conditional Markov models

Christopher Manning

# Sequence problems

- Many problems in NLP have data which is a sequence of characters, words, phrases, lines, or sentences …

- We can think of our task as one of labeling each item

| VBG | NN | IN | DT | NN | IN | NN |
|---|---|---|---|---|---|---|
| Chasing | opportunity | in | an | age | of | upheaval |

**POS tagging**

| PERS | O | O | O | ORG | ORG |
|---|---|---|---|---|---|
| Murdoch | discusses | future | of | News | Corp. |

**Named entity recognition**

| B | B | I | I | B | I | B | I | B | B |
|---|---|---|---|---|---|---|---|---|---|
| 而 | 相 | 对 | 于 | 这 | 些 | 品 | 牌 | 的 | 价 |

**Word segmentation**

Q
A
Q
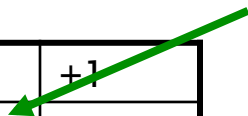A
A
A
A
Q
A

**Text segmentation**

# MEMM inference in systems

- For a Conditional Markov Model (CMM) a.k.a. a Maximum Entropy Markov Model (MEMM), the classifier makes a single decision at a time, conditioned on evidence from observations and previous decisions

- A larger space of sequences is usually explored via search

**Local Context**

**Decision Point**

| -3 | -2 | -1 | 0 | +1 |
|-----|-----|-----|-----|-----|
| DT | NNP | VBD | ??? | ??? |
| The | Dow | fell | 22.6 | % |

(Ratnaparkhi 1996; Toutanova et al. 2003, etc.)

**Features**

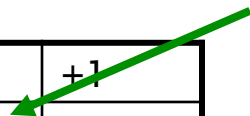| | |
|-----|-----|
| $W_0$ | 22.6 |
| $W_{+1}$ | % |
| $W_{-1}$ | fell |
| $T_{-1}$ | VBD |
| $T_{-1}$-$T_{-2}$ | NNP-VBD |
| hasDigit? | true |
| … | … |

# Example: POS Tagging

- Scoring individual labeling decisions is no more complex than standard classification decisions
  - We have some assumed labels to use for prior positions
  - We use features of those and the observed data (which can include current, previous, and next words) to predict the current label

Local Context

Decision Point

Features

| -3 | -2 | -1 | 0 | +1 |
|-----|-----|-----|------|-----|
| DT | NNP | VBD | ??? | ??? |
| The | Dow | fell | 22.6 | % |

| | |
|----------------------|---------|
| $W_0$ | 22.6 |
| $W_{+1}$ | % |
| $W_{-1}$ | fell |
| $T_{-1}$ | VBD |
| $T_{-1}$-$T_{-2}$ | NNP-VBD |
| hasDigit? | true |
| … | … |

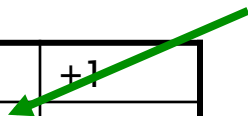(Ratnaparkhi 1996; Toutanova et al. 2003, etc.)

# Example: POS Tagging

- POS tagging Features can include:
  - Current, previous, next words in isolation or together.
  - Previous one, two, three tags.
  - Word-internal features: word types, suffixes, dashes, etc.

**Local Context**

| -3 | -2 | -1 | 0 | +1 |
|----|----|----|----|----|
| DT | NNP | VBD | ??? | ??? |
| The | Dow | fell | 22.6 | % |

**Decision Point**

**Features**

| | |
|----|----|
| $W_0$ | 22.6 |
| $W_{+1}$ | % |
| $W_{-1}$ | fell |
| $T_{-1}$ | VBD |
| $T_{-1}$-$T_{-2}$ | NNP-VBD |
| hasDigit? | true |
| … | … |

(Ratnaparkhi 1996; Toutanova et al. 2003, etc.)

# Inference in Systems

Sequence Level

Sequence Model

Inference

Sequence
Data

Local Level

Local
Data

Feature
Extraction

Label

Features

Classifier Type

Optimization

Smoothing

Label

Features

Maximum Entropy
Models

Conjugate
Gradient

Quadratic
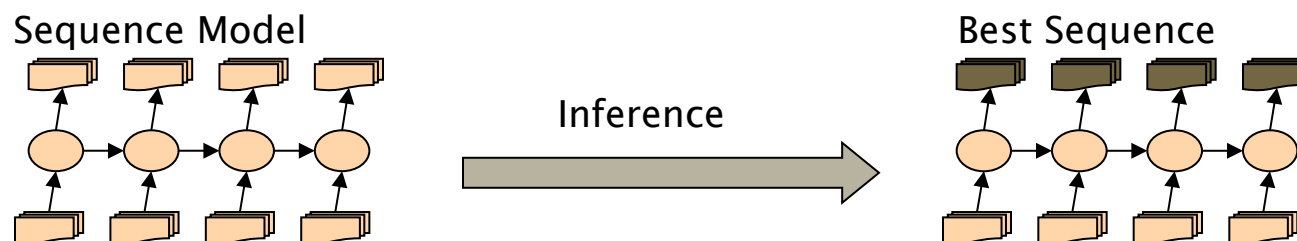Penalties

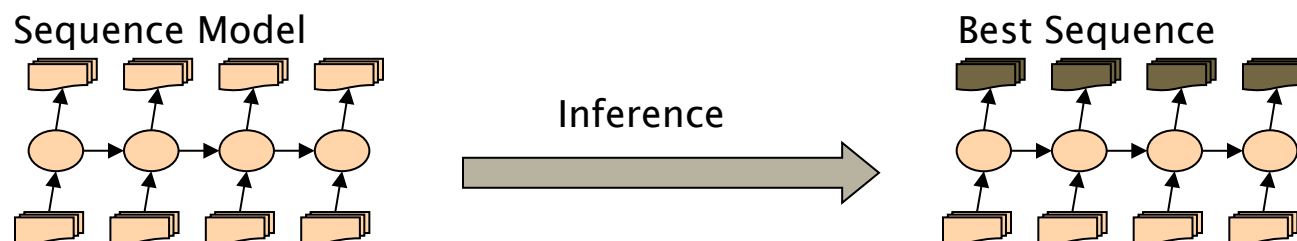# Greedy Inference

Sequence Model                    Best Sequence

Inference

- Greedy inference:
  - We just start at the left, and use our classifier at each position to assign a label
  - The classifier can depend on previous labeling decisions as well as observed data
- Advantages:
  - Fast, no extra memory requirements
  - Very easy to implement
  - With rich features including observations to the right, it may perform quite well
- Disadvantage:
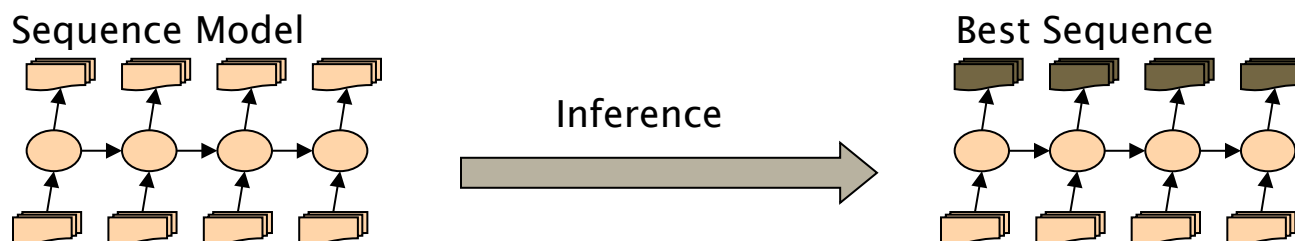  - Greedy. We make commit errors we cannot recover from

# Beam Inference

Sequence Model

Best Sequence

Inference

- Beam inference:
  - At each position keep the top $k$ complete sequences.
  - Extend each sequence in each local way.

  - The extensions compete for the $k$ slots at the next position.
- Advantages:
  - Fast; beam sizes of 3–5 are almost as good as exact inference in many cases.
  - Easy to implement (no dynamic programming required).
- Disadvantage:
  - Inexact: the globally best sequence can fall off the beam.

# Viterbi Inference

Sequence Model                                    Best Sequence

Inference

- Viterbi inference:
  - Dynamic programming or memoization.
  - Requires small window of state influence (e.g., past two states are relevant).
- Advantage:
  - Exact: the global best sequence is returned.
- Disadvantage:
  - Harder to implement long-distance state-state interactions (but beam inference tends not to allow long-distance resurrection of sequences anyway).

# CRFs [Lafferty, Pereira, and McCallum 2001]

- Another sequence model: Conditional Random Fields (CRFs)
- A whole-sequence conditional model rather than a chaining of local models.

$$P(c \mid d, \lambda) = \frac{\exp \sum_i \lambda_i f_i(c,d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c',d)}$$

- The space of $c$'s is now the space of sequences
  - But if the features $f_i$ remain local, the conditional sequence likelihood can be calculated exactly using dynamic programming
- Training is slower, but CRFs avoid causal-competition biases
- These (or a variant using a max margin criterion) are seen as the state-of-the-art these days … but in practice usually work much the same as MEMMs.

# Maximum entropy sequence models

Maximum entropy Markov models (MEMMs) or Conditional Markov models