Information Extraction

- Introduction
 - Task definition
 - Evaluation
 - IE system architecture
 - Named entity detection --- ML methods
- Relation extraction
 - Manually specified patterns
 - Machine learning approaches

HMMs for NE detection

```
American
Airlines
a
unit
of
AMR
Corp.
immediately
matched
the
move
spokesman
Tim
Wagner
said
```

Tag Set for NER

- IOB tags
 - B-xxx
 - First (i.e. Beginning) token in a NE of type xxx
 - I-xxx
 - Inside an entity of type xxx
 - **–** O
 - Outside all NEs

HMMs for NE detection

| American | $\mathrm{B}_{\mathit{ORG}}$ |
|-------------|-----------------------------|
| Airlines | ${ m I}_{ORG}$ |
| , | O |
| a | O |
| unit | O |
| of | O |
| AMR | $\mathrm{B}_{\mathit{ORG}}$ |
| Corp. | I_{ORG} |
| , | O |
| immediately | O |
| matched | O |
| the | O |
| move | O |
| , | O |
| spokesman | O |
| Tim | B_{PER} |
| Wagner | I_{PER} |
| said | O |
| | O |

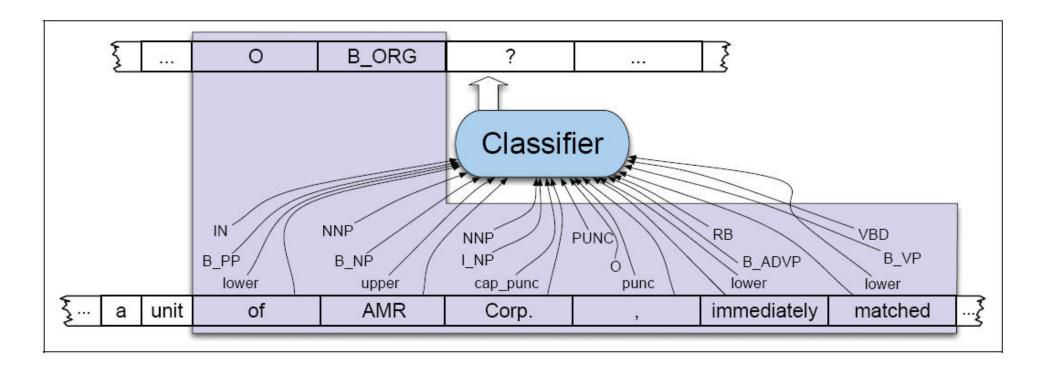
Figure, copyright J&M 2nd ed

HMMs for NE detection

- Just as for POS tagging except
 - States are IOB tags

Window-based Classification

- Fixed-size moving window
- Classify the target token as one of IOB



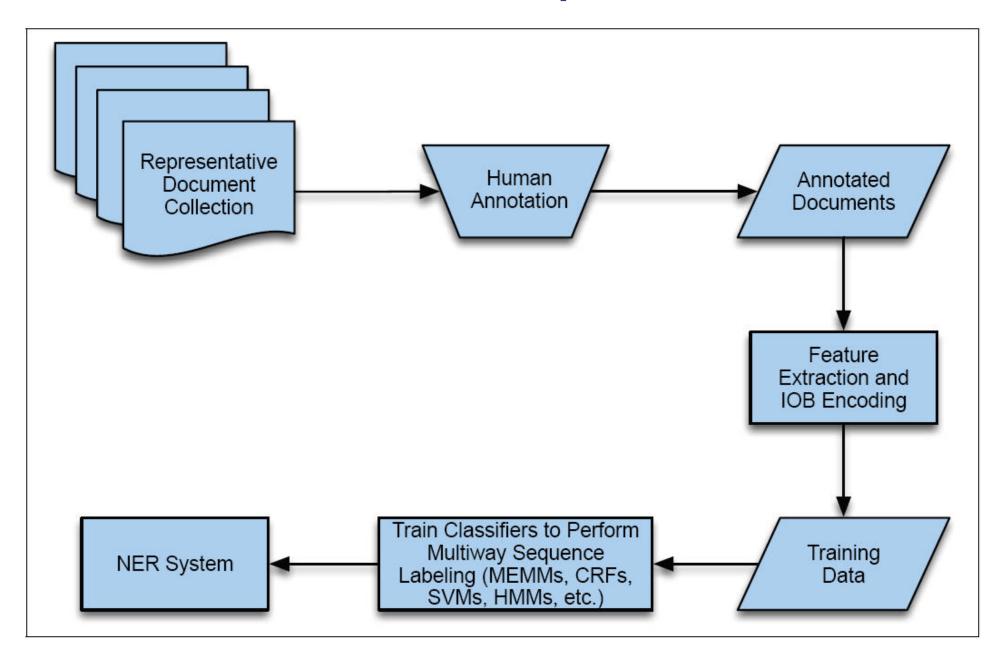
Feature engineering

· Word features can be especially useful

Table 3.1 Word features, examples and intuition behind them.²

| Word Feature | Example Text | Intuition |
|------------------------|---------------|------------------------------------|
| twoDigitNum | 90 | Two-digit year |
| fourDigitNum | 1990 | Four digit year |
| containsDigitAndAlpha | A8956-67 | Product code |
| containsDigitAndDash | 09-96 | Date |
| containsDigitAndSlash | 11/9/89 | Date |
| containsDigitAndComma | 23,000.00 | Monetary amount |
| containsDigitAndPeriod | 1.00 | Monetary amount, percentage |
| otherNum | 456789 | Other number |
| allCaps | BBN | Organization |
| capPeriod | М. | Person name initial |
| firstWord | first word of | No useful capitalization |
| | sentence | information |
| initCap | Sally | Capitalized word |
| lowerCase | can | Uncapitalized word |
| other | , | Punctuation marks, all other words |

End-to-end process



NE Results Using HMM's

Table 5.1 F-measure Scores. This table illustrates IdentiFinder's performance as compared to the best reported scores for each category.

| | Language | Best Rules | HMM |
|-------------|---------------|------------|------|
| Mixed Case | English (WSJ) | 96.4 | 94.9 |
| Upper Case | English (WSJ) | 89 | 93.6 |
| Speech Form | English (WSJ) | 74 | 90.7 |
| Mixed Case | Spanish | 93 | 90 |

Information Extraction

Introduction

- Task definition
- Evaluation
- IE system architecture
- Named entity detection
 Relation extraction

Person OUT: "Barry Nelson"

OUT position: "president"

VICTIM: "Jesuit priests"

PERP: "the FMLN"

- Manually specified patterns
- Machine learning approaches

Why bother????

Provide intuition for useful features for the machine learning approaches

Acquiring IE patterns

Goal

- Given a training set of annotated documents
 - answer keys / gold standard
- Develop a set of extraction patterns for each slot type.

Evergreen Information said Barry Nelsen, who had a heart-bypass operation last week, resigned as president and chief executive. The board formally accepted the resignation of Thomas Casey, its former chairman, who stepped down effective Feb. 2.

Martin Bell was named president, CEO, and chairman. Mr. Bell — who has been chief financial officer since the fall — also got voting control of 970,000 shares held by the Evergreen Partnership, a vehicle for the company's three co-founders.

Excluding these shares, Evergreen Informmillion shares or exercisable warrants outstar spokeswoman.

Type:

Person: "Barry Nelsen"

The computer products and services concern has cut its staff to fewer than 10 employees from about 35, and has deferred and reduced managers' salaries. In a press release, it said it believes the company is still viable.

Evergreen Information said Barry Nelsen, who had a heart-bypass operation last week, resigned as president and chief executive. The board formally accepted the resignation of Thomas Casey, its former chairman, who stepped down effective Feb. 2.

Martin Bell was named president, CEO, and chairman. Mr. Bell -- who has been chief financial officer since the fall -- also got voting control of 970,000 shares held by the Evergreen Partnership,

a vehicle for the company's three co-founder

Excluding these shares, Evergreen Information shares or exercisable warrants outstants spokeswoman.

In-out-event

Type: OUT

Position: PRESIDENT,

CHIEF EXECUTIVE

The computer products and services concern has cut its staff to fewer than 10 employees from about 35, and has deferred and reduced managers' salaries. In a press release, it said it believes the company is still viable.

Evergreen Information said Barry Nelsen, who had a heart-bypass operation last week, <u>resigned as president and chief executive</u>. The board formally accepted the resignation of Thomas Casey, its former chairman, who stepped down effective Feb. 2.

Martin Bell was named president, CEO, and chairman. Mr. Bell -- who has been chief financial officer since the fall -- also got voting control of 970,000 shares held by the Evergreen Partnership,

a vehicle for the company's three co-founder

Excluding these shares, Evergreen Information shares or exercisable warrants outstants spokeswoman.

In-out-event

Type: OUT

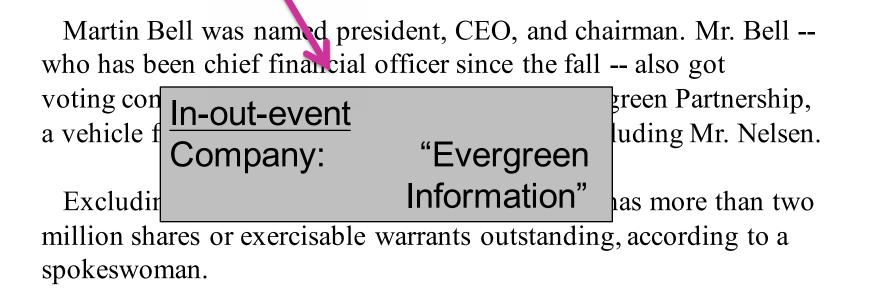
Person: "Barry Nelson"

Position: "president",

"chief executive"

The computer products and services concern has cut his stan to fewer than 10 employees from about 35, and has deferred and reduced managers' salaries. In a press release, it said it believes the company is still viable.

Evergreen Information said Barry Nelsen, who had a heart-bypass operation last week, resigned as president and chief executive. The board formally accepted the resignation of Thomas Casey, its former chairman, who stepped down effective Feb. 2.



The computer products and services concern has cut its staff to fewer than 10 employees from about 35, and has deferred and reduced managers' salaries. In a press release, it said it believes the company is still viable.

Natural disasters

The twister occurred without warning at approximately 7:15p.m. and destroyed *two mobile homes*.



Natural disaster

Type:

TORNADO

Damaged-obj: "two mobile homes"

Syntactico-semantic patterns

The twister occurred without warning at approximately 7:15p.m. and destroyed *two mobile homes*.

Pattern:

Trigger: "destroyed"

condition: active voice verb?

Slot: Damaged-Object

Position: direct-object

Also works here...

High winds destroyed <u>several office buildings</u> in the neighboring town.

Pattern:

Trigger: "destroyed"

condition: active voice verb?

Slot: Damaged-Object

Position: direct-object

And here...uh-oh

The divorce destroyed *his confidence*.

Pattern:

Trigger: "destroyed"

condition: active voice verb?

Slot: Damaged-Object

Position: direct-object

condition: DO is a physical-object?

Information Extraction

- Introduction
 - Task definition
 - Evaluation
 - IE system architecture
- Named entity detection
- Relation extraction
 - Manually specified patterns
 - Machine learning approaches

Why bother????
Provide intuition for useful features for the machine learning approaches

Learning IE patterns from examples

Goal

- Given a training set of annotated documents
 - "answer keys" or "gold standard"
 - text spans annotated with slot type
- Learn extraction patterns for each slot type using an appropriate machine learning algorithm.

Learning IE patterns

- Methods vary with respect to
 - The class of pattern learned (e.g. lexically-based regular expression, syntactic-semantic pattern)
 - Training corpus requirements
 - Amount and type of human feedback required
 - Degree of pre-processing necessary
 - Other resources/knowledge bases presumed

Entity and Relation Extraction

- Learning approaches
 - Weakly supervised methods
 - Mostly unsupervised methods
- Sequence-tagging methods
 - MEMM's
 - ILP for joint extraction of entities and relations
 - Opinion extraction

ML Approaches to IE

The twister occurred without warning at approximately 7:15p.m. and destroyed *two mobile homes*.



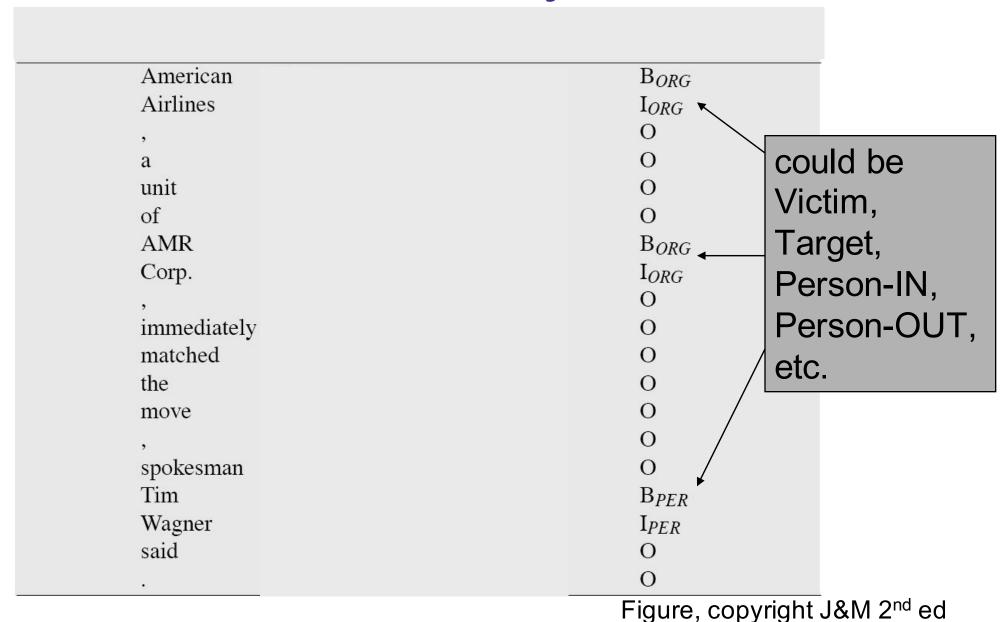
Natural disaster

Type:

TORNADO

Damaged-obj: "two mobile homes"

HMMs for entity detection

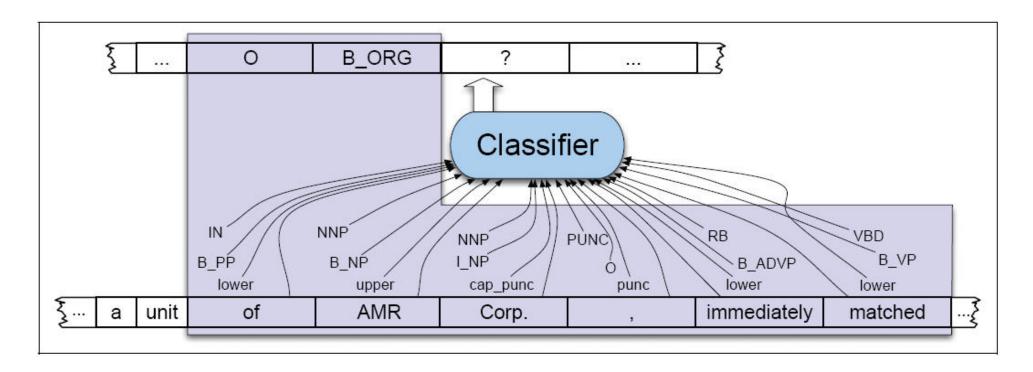


IOB

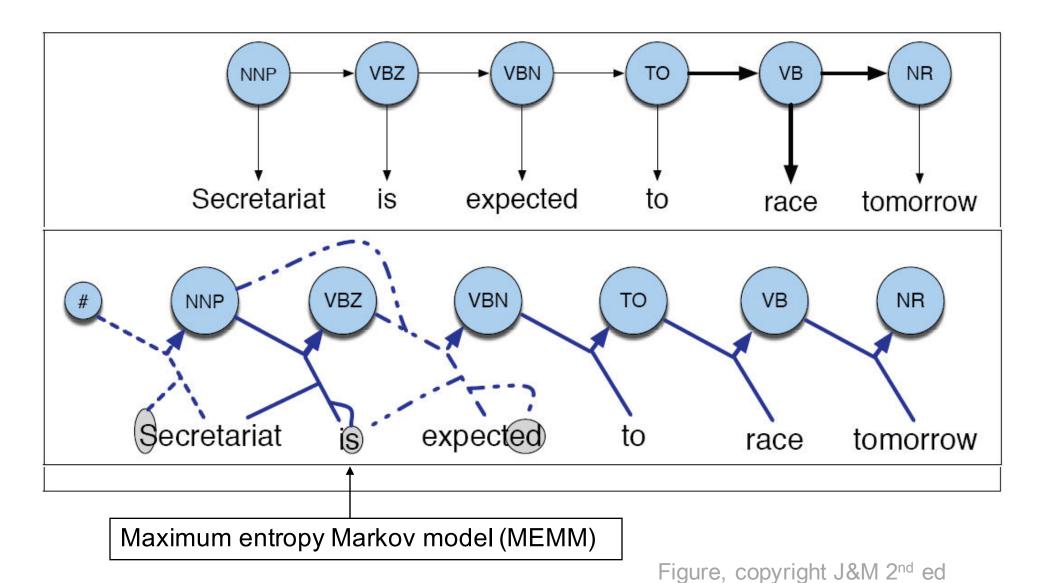
| $oldsymbol{B}_{tornado}$ | tornado | 0 | 0 | (|) (| 0 | |
|---|---------------------|--------|----------|-------------------|--------------------|--------------------|---|
| The twister occurred without warning at approximately | | | | | | | |
| B_{tim} | e I _{time} | 0 | 0 | $B_{d	ext{-}obj}$ | I _{d-obj} | I _{d-obj} | 0 |
| 7:15 | p.m. | and de | estroyed | two m | obile | homes | |

Feature extraction

 We'd like to be able to include lots of features as in classification-based approaches (e.g. SVMs, dtrees)

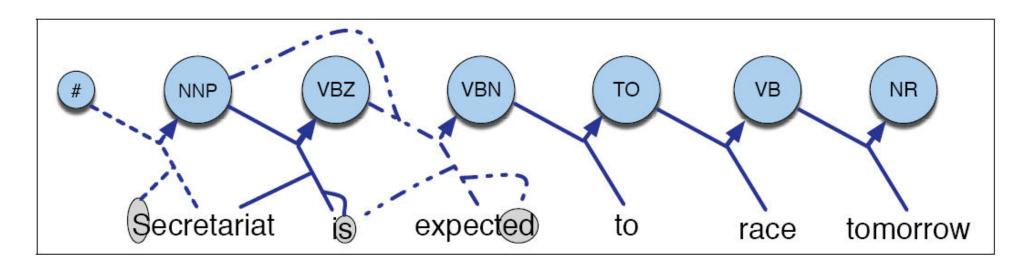


Not possible with HMMs



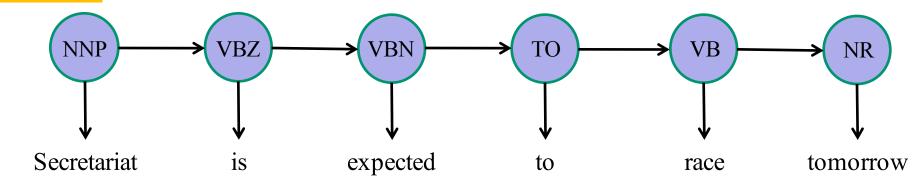
MEMM for p-o-s tagging

- Can condition on many features of the input
 - Capitalization
 - Morphology
 - Earlier words
 - Earlier tags

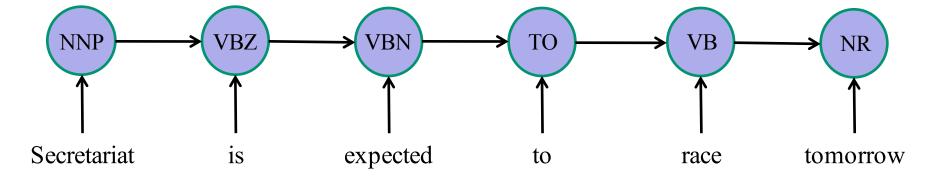


HMM vs. MEMM

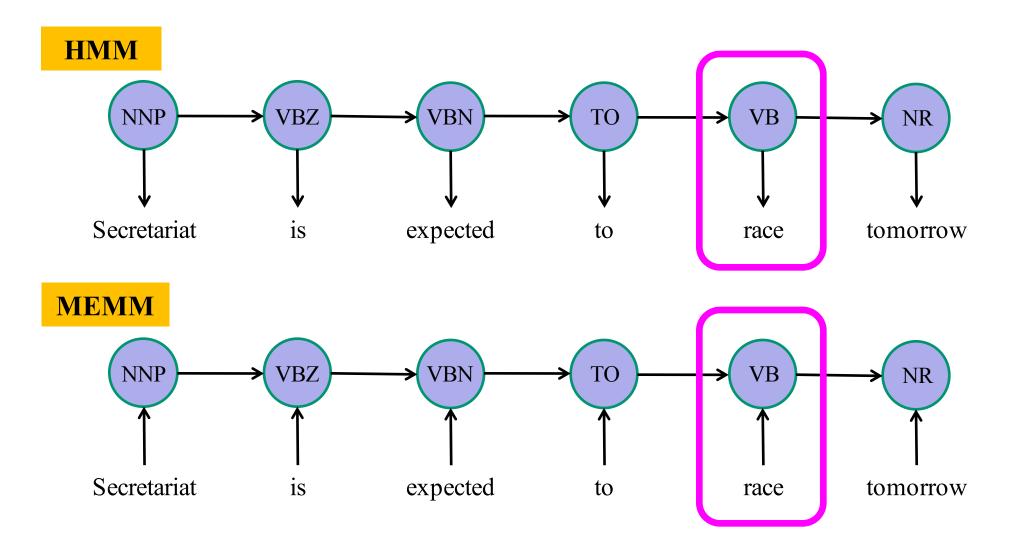
HMM

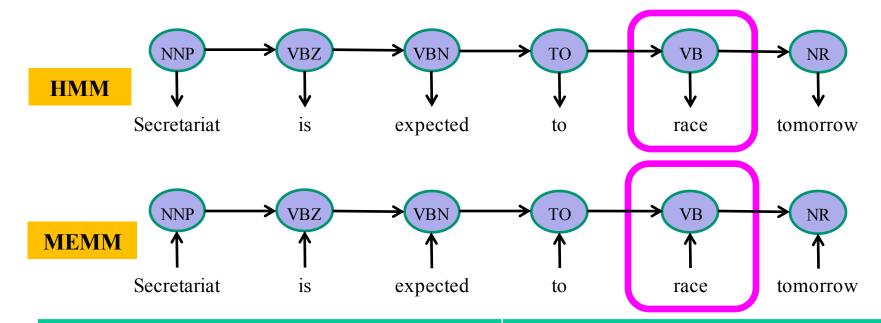


MEMM



HMM vs. MEMM





HMM

"Generative" models

- → joint probability p(words, tags)
- → "generate" input (in addition to tags)
- → but we need to predict tags, not words!

Probability of each slice = transition * emission = p(tag_i | tag_i-1) * p(word_i | tag_i)

→ Cannot incorporate long distance features

MEMM

"Discriminative" or "Conditional" models

- → conditional probability p(tags | words)
- → "condition" on input
- → Focusing only on predicting tags

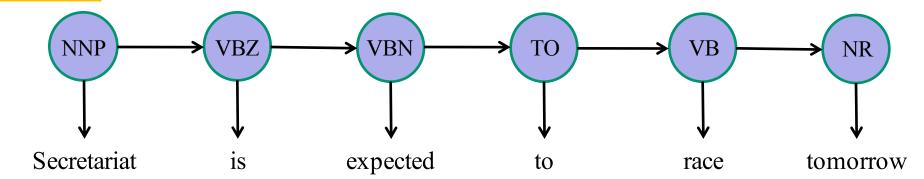
Probability of each slice =
p(tag_i | tag_i-1, word_i)
or
p(tag_i | tag_i-1, all words)

via a probabilistic classifier: NB, MaxEnt

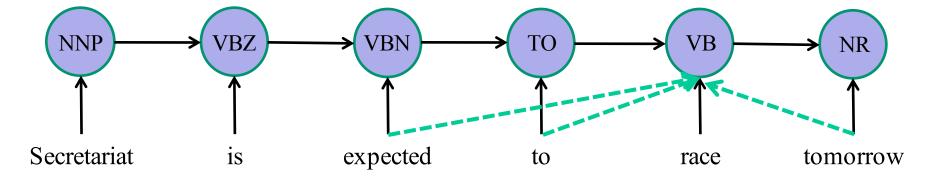
→ Can incorporate long distance features

HMM v.s. MEMM

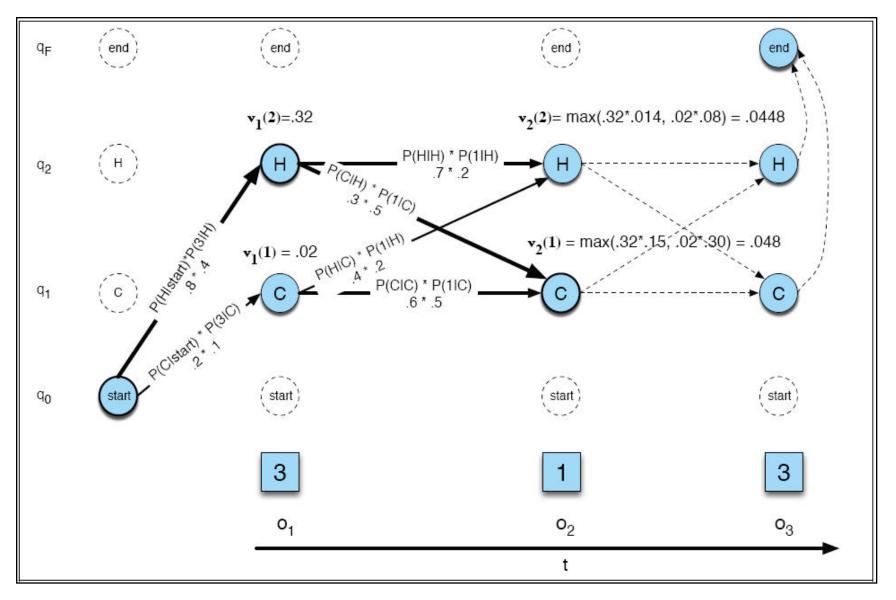
HMM



MEMM

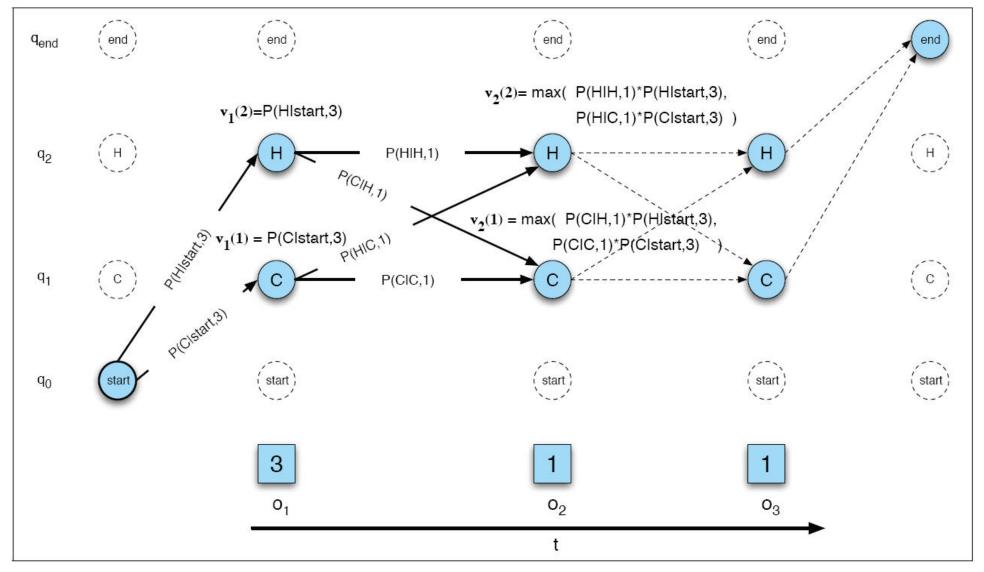


Decoding/inference in HMMs



Figure, copyright J&M 2nd ed

Decoding/inference in MEMMs



Figure, copyright J&M 2nd ed

Next...

Maximum Entropy classification