# Multi-Component Word Sense Disambiguation

## Massimiliano Ciaramita and Mark Johnson

## Brown University

**BLLIP:** http://www.cog.brown.edu/Research/nlp

# Outline

- <span style="color:red">**Pattern classification for WSD**</span>

  - **Features**

  - **Flat multiclass averaged perceptron**

- **Multi-component WSD**

  - **Generating external training data**

  - **Multi-component perceptron**

- **Experiments and results**

# Pattern classification for WSD

**English lexical sample**: **57 test words: 32 verbs, 20 nouns, 5 adjectives. For each word** $w$**:**

1. **compile a training set:** $S(w) = (x_i, y_i)^n$

   - $x_i \in \mathbb{R}^d$ **a vector of features**
   - $y_i \in Y(w)$**, one of the possible senses of** $w$

2. **learn a classifier on** $S(w)$**:** $H : \mathbb{R}^d \rightarrow Y(w)$

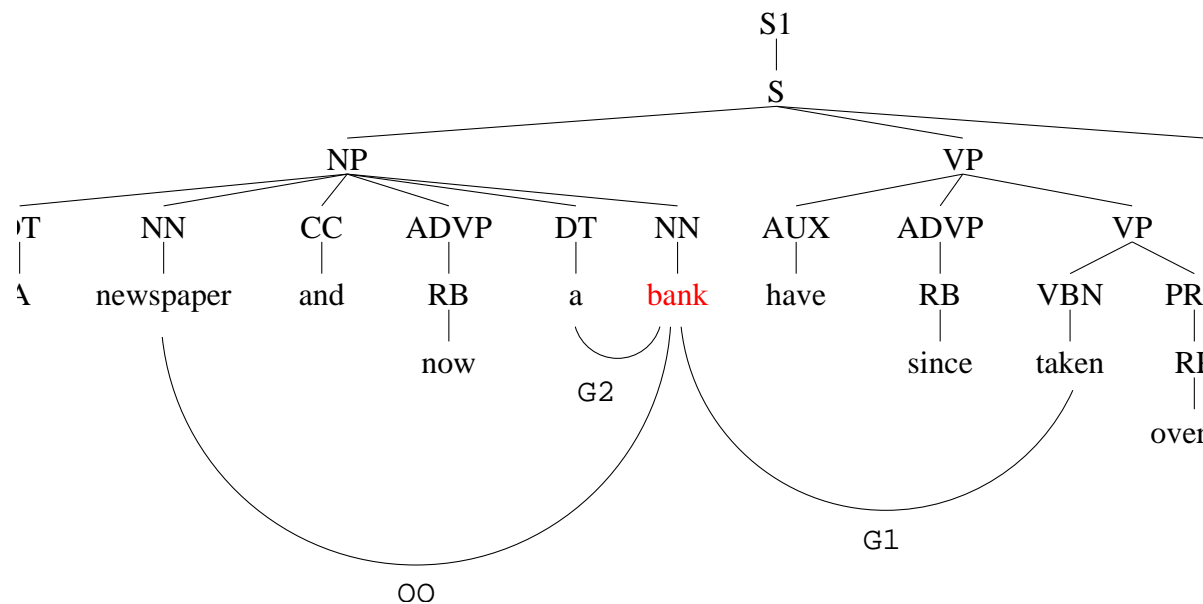3. **use the classifier to disambiguate the unseen test data**

# Features

- **Standard feature set for wsd (derived from (Yoong and Hwee, 2002))**

    - **"A-DT newspaper-NN and-CC now-RB a-DT $bank\text{-}NN$ have-AUX since-RB taken-VBN over-RB"**

- **POS of neighboring words** - $\mathrm{P}_{\mathrm{x,x}\in\{-3,-2,-1,0,+1,+2,+3\}}$**; e.g.,** $\mathrm{P}_{-1} = \mathrm{DT}$, $\mathrm{P}_0 = \mathrm{NN}$, $\mathrm{P}_{+1} = \mathrm{AUX}$, **...**

- **Surrounding words** - $\mathrm{WS}$**; e.g.,** $\mathrm{WS} = \mathrm{take_v}$, $\mathrm{WS} = \mathrm{over_r}$, $\mathrm{WS} = \mathrm{newspaper_n}$

- **N-grams:**

    - $\mathrm{NG}_{\mathrm{x,x}\in\{-2,-1,+1,+2\}}$**; e.g.,** $\mathrm{NG}_{-2} = \mathrm{now}$, $\mathrm{NG}_{+1} = \mathrm{have}$, $\mathrm{NG}_{+2} = \mathrm{take}$

    - $\mathrm{NG}_{\mathrm{x,y:(x,y)}\in\{(-2,-1),(-1,+1),(+1,+2)\}}$**; e.g.,** $\mathrm{NG}_{-2,-1} = \mathrm{now\_a}$, $\mathrm{NG}_{+1,+2} = \mathrm{have\_since}$

# Syntactic features (Charniak,2000)

- **Governing elements under a phrase** - $G_1$; **e.g.**, $G_1 = \texttt{take\_S}$

- **Governed elements under a phrase** - $G_2$; **e.g.**, $G_2 = \texttt{a\_NP}$, $G_2 = \texttt{now\_NP}$

- **Coordinates** - $\texttt{OO}$; **e.g.**, $\texttt{OO} = \texttt{newspaper}$

S1
|
S

NP                    VP

)T    NN      CC    ADVP    DT    NN    AUX    ADVP    VP    .

A    newspaper    and    RB    a    bank    have    RB    VBN    PR

now                      since    taken    RF

G2                                         over

OO                       G1

# Multiclass Perceptron (Crammer and Singer, 2003)

- **Discriminant function:** $\mathrm{H}(\mathbf{x}; \mathbf{V}) = \arg\max_{r=1}^{k} \langle v_r, x \rangle$

- **Input:** $\mathbf{V} \in \mathbb{R}^{|Y(w)| \times d}$, $d \approx 200,000$, **initialized as** $\mathbf{V} = 0$

- **Repeat** $\mathrm{T}$ **times - passes over training data or** *epochs*

$\mathtt{Multiclass\_Perceptron}((x,y)^n, \mathbf{V})$

**1**   **for** $i = 1$ **to** $i = n$

**2**   **do** $E = \{r : \langle v_r, x_i \rangle > \langle v_y, x_i \rangle\}$

**3**      **if** $|E| > 0$

**4**         **then** 1. $\tau_r = 1$ for $r = y$

**5**                  2. $\tau_r = 0$ for $r \notin E \cup \{y\}$

**6**                  3. $\tau_r = -\frac{1}{|E|}$ for $r \in E$

**7**                  **for** $r = 1$ **to** $r = k$

**8**                  **do** $v_r \leftarrow v_r + \tau_r x_i;$

# Averaged perceptron classifier

- **Perceptron's output:** $V^{(0)}, \dots, V^{(n)}$

- $V^{(i)}$ **is the weight matrix after the first** `i` **training items**

- **Final model:** $V = V^{(n)}$

- **Averaged perceptron: (Collins, 2002)**

  - **final model:** $V = \frac{1}{n} \sum_{i=1}^{n} V^{(i)}$
  - **reduces the effect of over-training**

# Outline

- **Pattern classification for WSD**

  - **Features**

  - **Flat multiclass perceptron**

- **Multi-component WSD**

  - **Generating external training data**

  - **Multi-component perceptron**
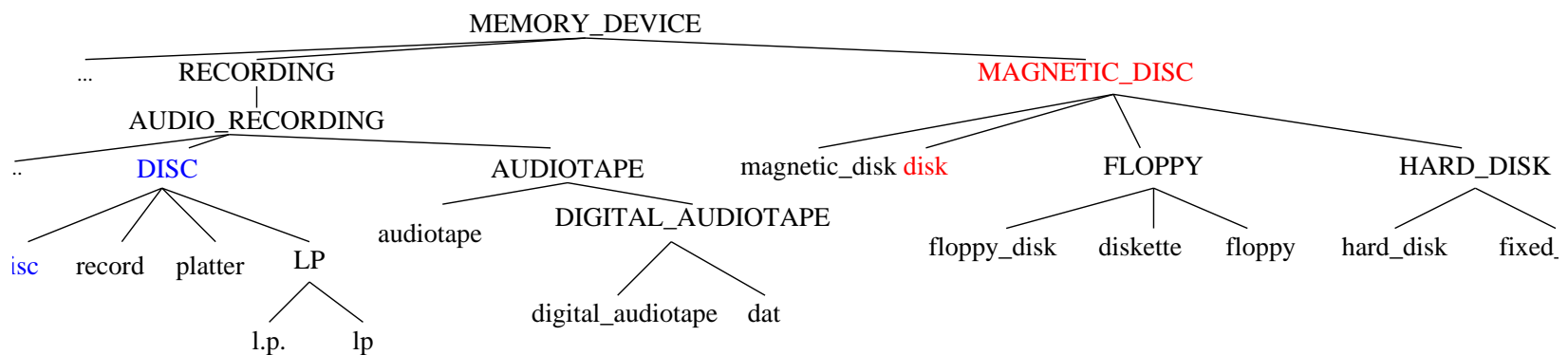
- **Experiments and results**

# Sparse data problem in WSD

- **Thousands of word senses** - **120,000 in Wordnet 2.0**

- **Very specific classes** - **50% of noun synsets contain one noun**

- **Problem:** **training instances often too few for fine-grained semantic distinctions**

- **Solution:**

    1. **use the hierarchy of Wordnet to find similar word senses and generate external training data for these senses**

    2. **integrate task-specific and external data with perceptron**

- **Intuition** - **to classify an instance of the noun disk additional knowledge about concepts such as other "audio" or "computer memory" devices could be helpful**
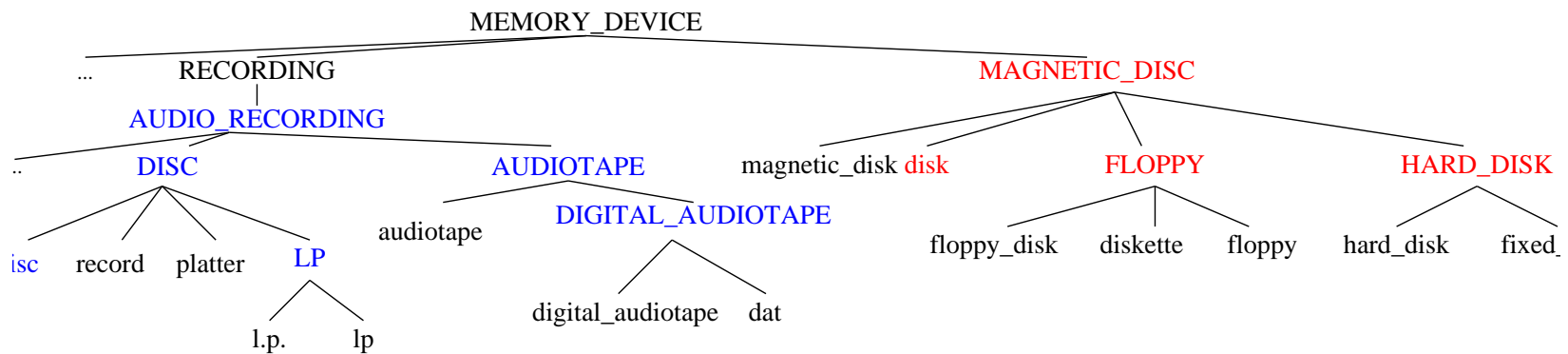
# Finding neighbor senses

- **disc$_1$ = memory device for information storing**
- **disc$_2$ = phonograph record**

```
                                    MEMORY_DEVICE
        ...        RECORDING                                              MAGNETIC_DISC
              AUDIO_RECORDING
       ..           DISC                    AUDIOTAPE      magnetic_disk disk        FLOPPY                  HARD_DISK
                                                    DIGITAL_AUDIOTAPE
     isc   record  platter   LP      audiotape                               floppy_disk  diskette  floppy   hard_disk   fixed_
                         l.p.    lp          digital_audiotape   dat
```

# Finding neighbor senses

- **neighbors(disc$_1$) = floppy disk, hard disk, ...**

- **neighbors(disc$_2$) = audio recording, lp, soundtrack, audiotape, talking book, digital audio tape, ...**

MEMORY_DEVICE

... RECORDING

AUDIO_RECORDING

.. DISC          AUDIOTAPE          magnetic_disk disk          FLOPPY          HARD_DISK

isc   record   platter   LP          audiotape          DIGITAL_AUDIOTAPE          floppy_disk   diskette   floppy          hard_disk   fixed_

l.p.   lp          digital_audiotape   dat

# External training data

- **Find neighbors**: for each sense $y$ of a noun or verb in the task a set $\hat{y}$ of $k = 100$ neighbor senses is generated from the Wordnet hierarchy

- **Generate new instances**: for each synset in $\hat{y}$ a training instance $(x_i, \hat{y}_i)$ is compiled from the corresponding Wordnet glosses (definitions/example sentences) using the same set of features

- **Result**: for each noun/verb

  1. task-specific training data $(x_i, y_i)^n$
  2. external training data $(x_i, \hat{y}_i)^m$

# Multi-component perceptron

- **Simplification of hierarchical perceptron (Ciaramita et al., 2003)**

- **A weight matrix $\mathbf{V}$ is trained on the task-specific data**

- **A weight matrix $\mathbf{M}$ is trained on the external data**

- **Discriminant function:**

$$\mathtt{H}(\mathbf{x}; \mathbf{V}, \mathbf{M}) = \arg \max_{\mathtt{y} \in \mathtt{Y}(\mathtt{w})} \lambda_{\mathtt{y}} \langle \mathtt{v}_{\mathtt{y}}, \mathbf{x} \rangle + \lambda_{\hat{\mathtt{y}}} \langle \mathtt{m}_{\hat{\mathtt{y}}}, \mathbf{x} \rangle$$

  $-$ $\lambda_{\mathtt{y}}$ **is an adjustable parameter that weights each component's contribution:** $\lambda_{\hat{\mathtt{y}}} = 1 - \lambda_{\mathtt{y}}$

# Multi-Component Perceptron

- **The algorithm learns $\mathbf{V}$ and $\mathbf{M}$ independently**

$\texttt{Multi-Component\_Perceptron}((\mathbf{x_i}, \mathbf{y_i})^n, (\mathbf{x_i}, \hat{\mathbf{y_i}})^m, \mathbf{V}, \mathbf{M})$

**1**   $\mathbf{V} \leftarrow 0$

**2**   $\mathbf{M} \leftarrow 0$

**3**   for $\mathbf{t} = 1$ to $\mathbf{i} = \mathbf{T}$

**4**   do $\texttt{Multiclass\_Perceptron}((\mathbf{x_i}, \mathbf{y_i})^n, \mathbf{V})$

**5**       $\texttt{Multiclass\_Perceptron}((\mathbf{x_i}, \mathbf{y_i})^n, \mathbf{M})$

**6**       $\texttt{Multiclass\_Perceptron}((\mathbf{x_i}, \mathbf{y_i})^m, \mathbf{M})$

# Outline

- **Pattern classification for WSD**

  - **Features**

  - **Flat multiclass averaged perceptron**

- **Multi-component WSD**

  - **Generating external training data**

  - **Multi-component perceptron**

- **Experiments and results**

# Experiments and results

- **One classifier trained for each test word**

- **Adjectives**: standard perceptron, only set T

- **Verbs/nouns**: multicomponent perceptron, set T and $\lambda_y$

- **Cross-validation** experiments on the training data for each test word:

  1. choose the value for $\lambda_y$; $\lambda_y = 1$ use only the "flat" perceptron, or $\lambda_y = 0.5$ use both component equally weighted
  2. choose the number of iterations T

- Average T value $= 13.9$

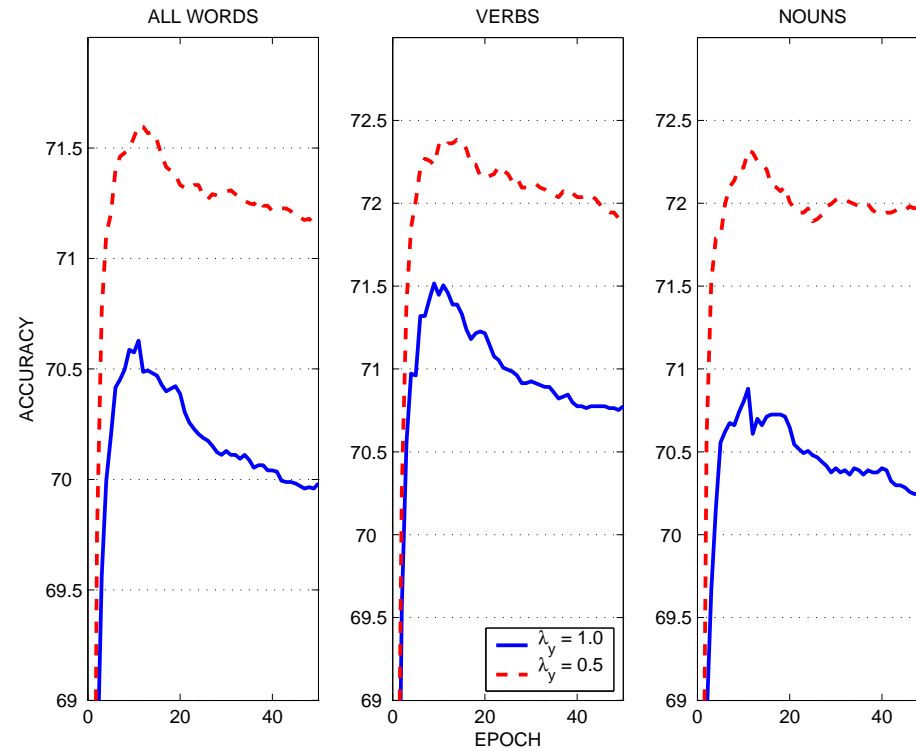- For 37 out of 52 nouns/verbs $\lambda_y = 0.5$; the two-component model is more accurate than the flat perceptron

# English Lexical Sample Results

| Measure | Precision | Recall | Attempted % |
|---|---|---|---|
| Fine all POS | 71.1 | 71.1 | 100 |
| Coarse all POS | 78.1 | 78.1 | 100 |
| Fine verbs | 72.5 | 72.5 | 100 |
| Coarse verbs | 80.0 | 80.0 | 100 |
| Fine nouns | 71.3 | 71.3 | 100 |
| Coarse nouns | 77.4 | 77.4 | 100 |
| Fine adjectives | 49.7 | 49.7 | 100 |
| Coarse adjectives | 63.5 | 63.5 | 100 |

# Flat vs. Multi-component: cross validation on train

# Conclusion

- **Advantages** of the multi-component perceptron trained on neighbors' data

  - **Neighbors**: one "supersense" for each sense, same amount of additional data per sense

  - **Simpler model**: smaller variance more homogeneous external data

  - **Efficiency**: fast and efficient training

  - **Architecture**: simple, easy to add any number of (weighted) "components"