

The Brown et al. Word Clustering Algorithm

Michael Collins, Columbia University

The Brown Clustering Algorithm

- ▶ Input: a (large) corpus of words
- ▶ Output 1: a partition of words into *word clusters*
- ▶ Output 2 (generalization of 1): a hierarchical word clustering

Example Clusters (from Brown et al, 1992)

Friday Monday Thursday Wednesday Tuesday Saturday Sunday weekends Sundays Saturdays
June March July April January December October November September August
people guys folks fellows CEOs chaps doubters commies unfortunates blokes
down backwards ashore sideways southward northward overboard aloft downwards adrift
water gas coal liquid acid sand carbon steam shale iron
great big vast sudden mere sheer gigantic lifelong scant colossal
man woman boy girl lawyer doctor guy farmer teacher citizen
American Indian European Japanese German African Catholic Israeli Italian Arab
pressure temperature permeability density porosity stress velocity viscosity gravity tension
mother wife father son husband brother daughter sister boss uncle
machine device controller processor CPU printer spindle subsystem compiler plotter
John George James Bob Robert Paul William Jim David Mike
anyone someone anybody somebody
feet miles pounds degrees inches barrels tons acres meters bytes
director chief professor commissioner commander treasurer founder superintendent dean cus-
todian

A Sample Hierarchy (from Miller et al., NAACL 2004)

| | |
|-----------------|--------------------------|
| lawyer | 1000001101000 |
| newspaperman | 100000110100100 |
| stewardess | 100000110100101 |
| toxicologist | 10000011010011 |
| slang | 1000001101010 |
| babysitter | 100000110101100 |
| conspirator | 1000001101011010 |
| womanizer | 1000001101011011 |
| mailman | 10000011010111 |
| salesman | 100000110110000 |
| bookkeeper | 1000001101100010 |
| troubleshooter | 10000011011000110 |
| bouncer | 10000011011000111 |
| technician | 1000001101100100 |
| janitor | 1000001101100101 |
| saleswoman | 1000001101100110 |
| ... | |
| Nike | 1011011100100101011100 |
| Maytag | 10110111001001010111010 |
| Generali | 10110111001001010111011 |
| Gap | 1011011100100101011110 |
| Harley-Davidson | 10110111001001010111110 |
| Enfield | 101101110010010101111110 |
| genus | 101101110010010101111111 |
| Microsoft | 10110111001001011000 |
| Ventritex | 101101110010010110010 |
| Tractebel | 1011011100100101100110 |
| Synopsys | 1011011100100101100111 |
| WordPerfect | 1011011100100101101000 |
| ... | |
| John | 101110010000000000 |
| Consuelo | 101110010000000001 |
| Jeffrey | 101110010000000010 |
| Kenneth | 10111001000000001100 |
| Phillip | 101110010000000011010 |
| WILLIAM | 101110010000000011011 |
| Timothy | 10111001000000001110 |
| Terrence | 101110010000000011110 |

The Intuition

- ▶ Similar words appear in similar contexts
- ▶ More precisely: similar words have similar distributions of words to their immediate left and right

The Formulation

- ▶ \mathcal{V} is the set of all words seen in the corpus w_1, w_2, \dots, w_n
- ▶ Say $C : \mathcal{V} \rightarrow \{1, 2, \dots, k\}$ is a *partition* of the vocabulary into k classes
- ▶ The model:

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n e(w_i | C(w_i)) q(C(w_i) | C(w_{i-1}))$$

(note: $C(w_0)$ is a special start state)

An Example

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n e(w_i | C(w_i)) q(C(w_i) | C(w_{i-1}))$$

An Example

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n e(w_i | C(w_i)) q(C(w_i) | C(w_{i-1}))$$

$$C(\text{the}) = 1, \quad C(\text{dog}) = C(\text{cat}) = 2, \quad C(\text{saw}) = 3$$

An Example

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n e(w_i | C(w_i)) q(C(w_i) | C(w_{i-1}))$$

$$C(\text{the}) = 1, \quad C(\text{dog}) = C(\text{cat}) = 2, \quad C(\text{saw}) = 3$$

$$e(\text{the}|1) = 1, \quad e(\text{cat}|2) = e(\text{dog}|2) = 0.5, \quad e(\text{saw}|3) = 1$$

An Example

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n e(w_i | C(w_i)) q(C(w_i) | C(w_{i-1}))$$

$$C(\text{the}) = 1, \quad C(\text{dog}) = C(\text{cat}) = 2, \quad C(\text{saw}) = 3$$

$$e(\text{the}|1) = 1, \quad e(\text{cat}|2) = e(\text{dog}|2) = 0.5, \quad e(\text{saw}|3) = 1$$

$$q(1|0) = 0.2, \quad q(2|1) = 0.4, \quad q(3|2) = 0.3, \quad q(1|3) = 0.6$$

An Example

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n e(w_i | C(w_i)) q(C(w_i) | C(w_{i-1}))$$

$$C(\text{the}) = 1, \quad C(\text{dog}) = C(\text{cat}) = 2, \quad C(\text{saw}) = 3$$

$$e(\text{the}|1) = 1, \quad e(\text{cat}|2) = e(\text{dog}|2) = 0.5, \quad e(\text{saw}|3) = 1$$

$$q(1|0) = 0.2, \quad q(2|1) = 0.4, \quad q(3|2) = 0.3, \quad q(1|3) = 0.6$$

$$p(\text{the dog saw the cat}) =$$

The Brown Clustering Model

A Brown clustering model consists of:

- ▶ A vocabulary \mathcal{V}
- ▶ A function $C : \mathcal{V} \rightarrow \{1, 2, \dots, k\}$ defining a *partition* of the vocabulary into k classes
- ▶ A parameter $e(v|c)$ for every $v \in \mathcal{V}$, $c \in \{1 \dots k\}$
- ▶ A parameter $q(c'|c)$ for every $c', c \in \{1 \dots k\}$

Measuring the Quality of C

- ▶ How do we measure the quality of a partition C ?

$$\begin{aligned}\text{Quality}(C) &= \sum_{i=1}^n \log e(w_i | C(w_i)) q(C(w_i) | C(w_{i-1})) \\ &= \sum_{c=1}^k \sum_{c'=1}^k p(c, c') \log \frac{p(c, c')}{p(c)p(c')} + G\end{aligned}$$

where G is a constant

- ▶ Here

$$p(c, c') = \frac{n(c, c')}{\sum_{c, c'} n(c, c')} \quad p(c) = \frac{n(c)}{\sum_c n(c)}$$

where $n(c)$ is the number of times class c occurs in the corpus, $n(c, c')$ is the number of times c' is seen following c , under the function C

A First Algorithm

- ▶ We start with $|\mathcal{V}|$ clusters: each word gets its own cluster
- ▶ Our aim is to find k final clusters
- ▶ We run $|\mathcal{V}| - k$ merge steps:
 - ▶ At each merge step we pick two clusters c_i and c_j , and merge them into a single cluster
 - ▶ We greedily pick merges such that

Quality(C)

for the clustering C after the merge step is maximized at each stage

- ▶ Cost? Naive = $O(|\mathcal{V}|^5)$. Improved algorithm gives $O(|\mathcal{V}|^3)$: still too slow for realistic values of $|\mathcal{V}|$

A Second Algorithm

- ▶ Parameter of the approach is m (e.g., $m = 1000$)
- ▶ Take the top m most frequent words, put each into its own cluster, c_1, c_2, \dots, c_m
- ▶ For $i = (m + 1) \dots |\mathcal{V}|$
 - ▶ Create a new cluster, c_{m+1} , for the i 'th most frequent word. We now have $m + 1$ clusters
 - ▶ Choose two clusters from $c_1 \dots c_{m+1}$ to be merged: pick the merge that gives a maximum value for $\text{Quality}(C)$. We're now back to m clusters
- ▶ Carry out $(m - 1)$ final merges, to create a full hierarchy

Running time: $O(|\mathcal{V}|m^2 + n)$ where n is corpus length

Miller et al, NAACL 2004

Name Tagging with Word Clusters and Discriminative Training

Scott Miller, Jethran Guinness, Alex Zamanian

BBN Technologies

10 Moulton Street

Cambridge, MA 02138

szmiller@bbn.com

At a recent meeting, we presented name-tagging technology to a potential user. The technology had performed well in formal evaluations, had been applied successfully by several research groups, and required only annotated training examples to configure for new name classes. Nevertheless, it did not meet the user's needs.

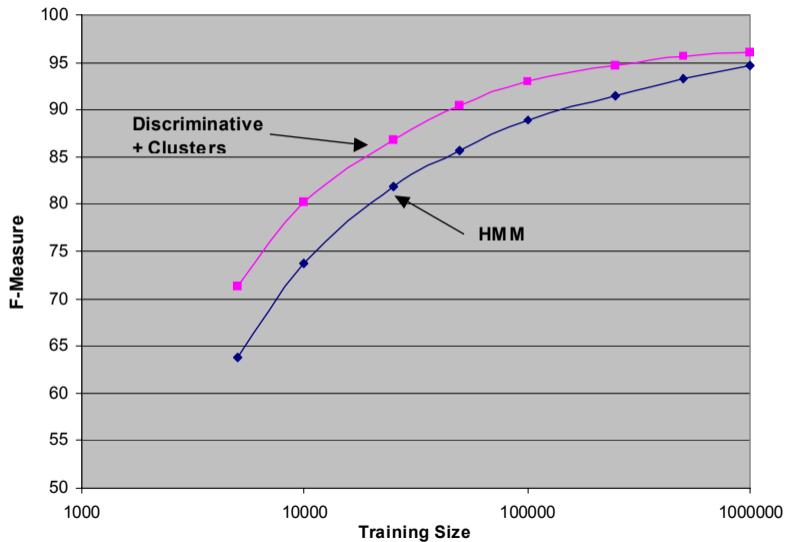
Miller et al, NAACL 2004

To achieve reasonable performance, the HMM-based technology we presented required roughly 150,000 words of annotated examples, and over a million words to achieve peak accuracy. Given a typical annotation rate of 5,000 words per hour, we estimated that setting up a name finder for a new problem would take four person days of annotation work – a period we considered reasonable. However, this user's problems were too dynamic for that much setup time. To be useful, the system would have to be trainable in minutes or hours, not days or weeks.

Miller et al, NAACL 2004

1. Tag + PrevTag
2. Tag + CurWord
3. Tag + CapAndNumFeatureOfCurWord
4. ReducedTag + CurWord
 //collapse start and continue tags
5. Tag + PrevWord
6. Tag + NextWord
7. Tag + DownCaseCurWord
8. Tag + Pref8ofCurrWord
9. Tag + Pref12ofCurrWord
10. Tag + Pref16ofCurrWord
11. Tag + Pref20ofCurrWord
12. Tag + Pref8ofPrevWord
13. Tag + Pref12ofPrevWord
14. Tag + Pref16ofPrevWord
15. Tag + Pref20ofPrevWord
16. Tag + Pref8ofNextWord
17. Tag + Pref12ofNextWord
18. Tag + Pref16ofNextWord
19. Tag + Pref20ofNextWord

Miller et al, NAACL 2004



Miller et al, NAACL 2004

