

# **Large-Margin Methods for Natural Language Learning**

Michael Collins  
AT&T Labs-Research

# MOTIVATION: APPLICATIONS

- Sequential data is **Everywhere**

Biological data

UCUCUGCCCCAGAGCGGAUACACGUGUCCGAACGUCAC  
UAAGUGUCCGCGCUAGGGUGAUGACAGUCGUGCC...

Text

The screenshot shows the front page of The New York Times. The main headline is 'Agency Sues White House to Force Disclosure' by Susan Vassilakis. Other articles include 'Middle East Cease-Fire Talks Resume', 'Angolan Army Kills Rebel Leader', and 'Kidnapping Suspect Details Plot to Attack U.S. Consulate'. The page layout includes a navigation bar, a main content area with multiple columns of text and images, and a sidebar with 'MORE HEADLINES'.

The screenshot shows the CNN.com website. The main headline is 'Neighbor Arrested in DANI VAN BARI CASE'. Other articles include 'Police officers search a vehicle at Wall Street', 'OTHER TOP NEWS', 'WAR AGAINST TERROR', and 'TIME: A casualty of war terror'. The page layout includes a navigation bar, a main content area with multiple columns of text and images, and a sidebar with 'MORE HEADLINES'.

## Estimators for Stochastic "Unification-Based" Grammars\*

Mark Johnson  
Cognitive and Linguistic Sciences  
Brown University

Stuart Geman  
Applied Mathematics  
Brown University

Stephen Canon  
Cognitive and Linguistic Sciences  
Brown University

Zhiyi Chi  
Dept. of Statistics  
The University of Chicago

Stefan Riecker  
Institut für Mathematik  
Universität Stuttgart

### Abstract

Log-linear models provide a statistically sound framework for Stochastic "Unification-Based" Grammars (SUBGs) and stochastic variants of other kinds of grammars. We describe two computationally tractable ways of estimating the parameters of such grammars from a training corpus of syntactic analyses, and apply these to estimate a stochastic version of Lexical-Functional Grammar.

### 1 Introduction

Probabilistic methods have revolutionized computational linguistics. They can provide a systematic treatment of preferences in parsing. Given a suitable estimator procedure, stochastic models can be "trained" to reflect the properties of a corpus. On the other hand, "Unification-Based" Grammars (UBGs) can express a variety of linguistically-important syntactic and semantic constraints. However, deriving Stochastic "Unification-Based" Grammars (SUBGs) has not proved as straightforward as might be hoped.

### 2 Features in SUBGs

We follow the statistical literature in using the term *feature* to refer to the properties that parameters are associated with. We use the word "attribute" here for the attributes of features of a UBG's feature structure. As to the set of all possible grammatical or well-defined analyses, each feature  $f$  maps syntactic analyses to a real value  $f(x)$ . The form of a syntactic analysis depends on the underlying linguistic theory. For example, in a PCFG a tuple consisting of (at least) a constituent, an L-structure and a mapping from structure nodes to feature structure nodes, and for a Chomskian transformational grammar a would be a derivation.

\* This research was supported by the National Science Foundation (IBN-0720010), the Office of Naval Research (N00014-05-1-0402), and Office of Naval Research (N00014-05-1-0402).

# NEED TO INDUCE STRUCTURE

## SEGMENTATIONS

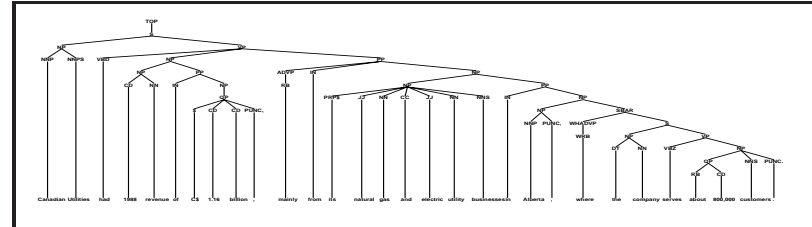
UCUCUGCCCCAGAGCGAUACACG  
UGUCCGAACGUCACUAAGUGUCC  
GCGCUAGGGUGAUGACAGUCGUG



UCUCUGCCCCAGAGCGAUACACG  
UGUCCGAACGUCACUAAGUGUCC  
GCGCUAGGGUGAUGACAGUCGUG

## HIERARCHICAL STRUCTURES (TREES)

Canadian Utilities had 1988 revenue of C\$ 1.16 billion, mainly from its natural gas and electric utility businesses in Alberta, where the company serves about 800,000 customers.



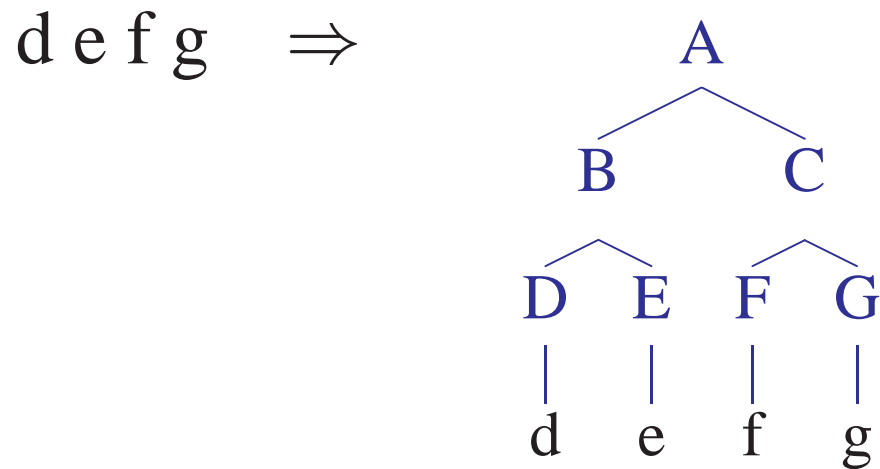
## TWO FUNDAMENTAL PROBLEMS

**TAGGING:** Strings to **Tagged Sequences**

a b e e a f h j  $\Rightarrow$  a/C b/D e/C e/C a/D f/C h/D j/C

**PARSING:** Strings to **Trees**

d e f g  $\Rightarrow$  (A (B (D d) (E e)) (C (F f) (G g)))



# OVERVIEW

- **The basic goal:** Supervised learning for recovering structure
- **A challenge:** Not like “typical” problems in statistics
  - Functions from strings to other discrete structures
- **Some new methods:**
  - How the voted perceptron algorithm can be applied
  - Global features
  - **Kernels** over discrete structures
  - Efficient training of weighted automata
  - Advantages over existing methods:  
new representations, discriminative parameter estimation

# LANGUAGE GENERATION [Walker et al. 2001]

Input = text plan

```
implicit-confirm(orig-city:NEWARK)
implicit-confirm(dest-city:DALLAS)
implicit-confirm(month:9)
implicit-confirm(day-number:1)
request(depart-time)
```

Output = text

What time would you like to travel on September the 1st to Dallas from Newark?

Leaving on September 1st. What time would you like to travel from Newark to Dallas?

Leaving in September. Leaving on the 1st. What time would you, traveling from Newark to Dallas, like to leave?

## A CHALLENGE

- Interesting class of learning problems:  
**Functions from strings to other discrete structures**
- “Typical” problems in statistics:
  - Classification functions  $F : \mathbb{R}^n \rightarrow \{-1, +1\}$
  - Regression functions  $F : \mathbb{R}^n \rightarrow \mathbb{R}$

## COMMON METHODS

- Stochastic automata
  - Probabilistic finite-state machines  
Hidden Markov Models (HMMs)
  - Probabilistic context-free grammars  
PCFGs
- Limitations
  - Representation
  - Strong assumptions in parameter estimation

## NEW METHODS

- “Distribution-free” methods for classification (regression)  
(Support vector machines, Voted Perceptron, Boosting)
- Very different statistical basis from ML/Bayesian methods  
Uniform convergence bounds, VC-dimension,  
**Large-margin bounds**
- Questions:
  - Can these methods be derived for tagging, parsing?
  - Do the statistical guarantees (e.g. about generalization) apply to tagging and parsing?
  - What **representations** can be used?

## MAIN RESULTS

- A variant of the voted perceptron algorithm for parsing, tagging, other NLP problems
- Experiments: Improvements over state-of-the-art baselines, using “global” features
- New kernels for discrete structures such as trees (e.g., the “all subtrees” representation)
- Experiments: Improvements over state-of-the-art baselines, using these kernels
- Efficient discriminative algorithms for weighted automata (Alternative parameter estimation method to Maximum-entropy taggers, Conditional Markov Random Fields)

# LINEAR MODELS FOR TAGGING AND PARSING

- Three components:

**GEN** is a function from a string to a set of **candidates**

$\Phi$  maps a candidate to a feature vector

**W** is a parameter vector

Related to Markov Random Field approaches

[Lafferty et al. 2001]

[Johnson, Geman, Canon, Chi, and Riezler 1999]

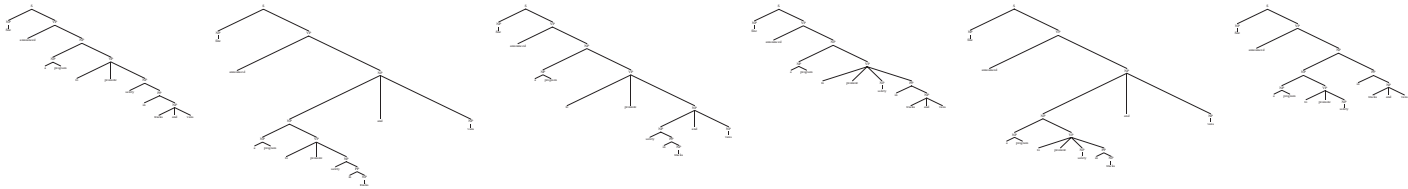
[Abney 1997]

[Della Pietra, Della Pietra, and Lafferty 1997]

[Ratnaparkhi, Roukos and Ward 94]

She announced a program to promote safety in trucks and vans

⇓ GEN



⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⟨1, 1, 3, 5⟩

⟨2, 0, 0, 5⟩

⟨1, 0, 1, 5⟩

⟨0, 0, 3, 0⟩

⟨0, 1, 0, 5⟩

⟨0, 0, 1, 5⟩

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

13.6

12.2

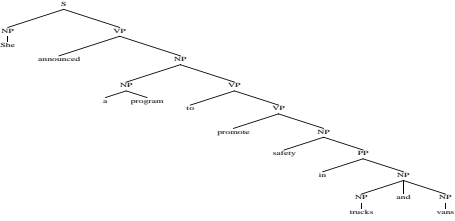
12.1

3.3

9.4

11.1

⇓ arg max



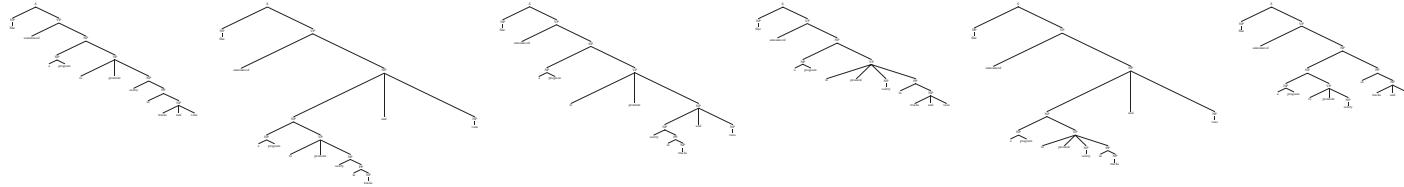
# COMPONENT 1: GEN

- **GEN** enumerates a set of **candidates** for a sentence

---

She announced a program to promote safety in trucks and vans

⇓ **GEN**



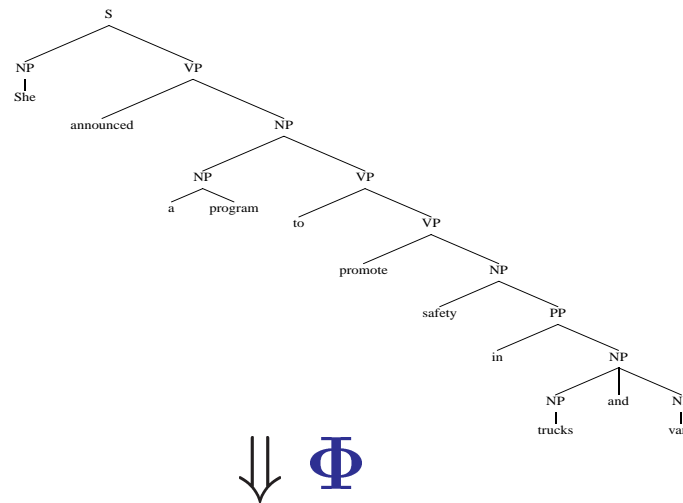
Note: Optimality Theory [[Prince and Smolensky](#)]

## EXAMPLES OF GEN

- A context-free grammar
- A finite-state machine
- Top  $N$  most probable analyses from a probabilistic grammar

## COMPONENT 2: $\Phi$

- $\Phi$  maps a candidate to a **feature vector**  $\in \mathbb{R}^d$
  - $\Phi$  defines the **representation** of a candidate
- 



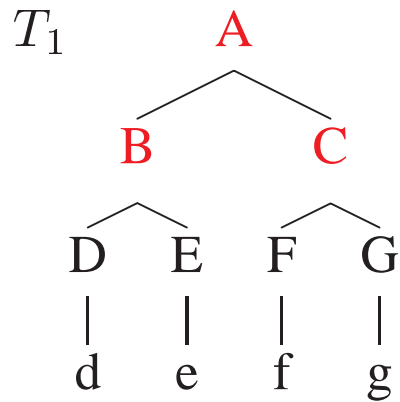
$\langle 1, 0, 2, 0, 0, 15, 5 \rangle$

# FEATURES

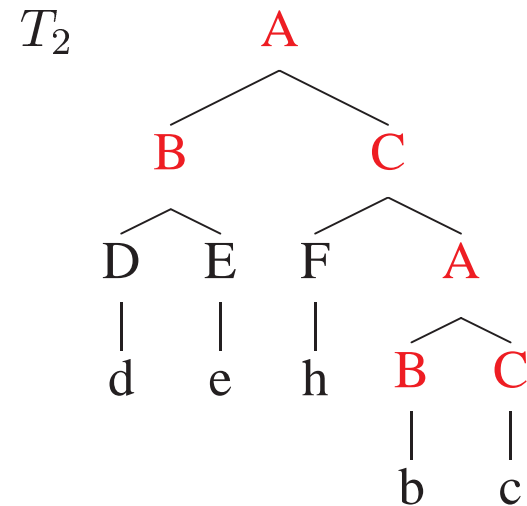
- A “feature” is a function on a structure, e.g.,

$$h(x) = \text{Number of times } \boxed{\begin{array}{c} A \\ \wedge \\ B \quad C \end{array}} \text{ is seen in } x$$

---



$$h(T_1) = 1$$



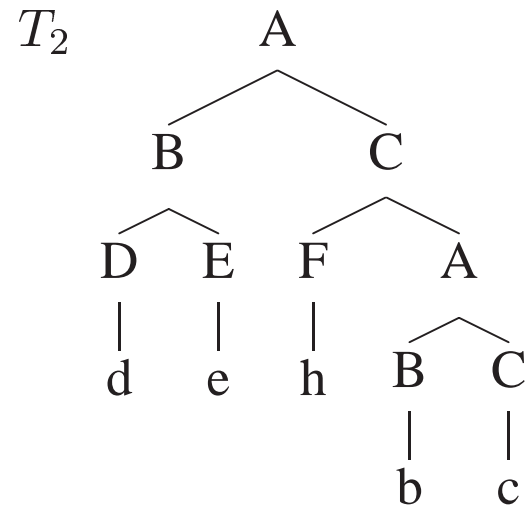
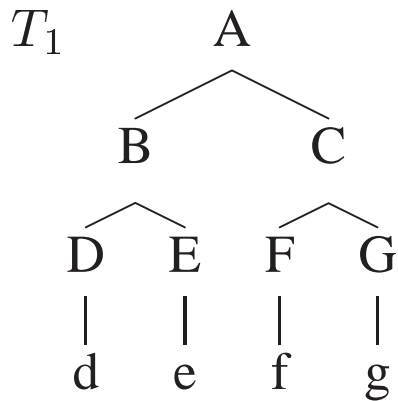
$$h(T_2) = 2$$

# FEATURE VECTORS

- A set of functions  $h_1 \dots h_d$  define a **feature vector**

$$\Phi(x) = \langle h_1(x), h_2(x) \dots h_d(x) \rangle$$

---



$$\Phi(T_1) = \langle 1, 0, 0, 3 \rangle$$

$$\Phi(T_2) = \langle 2, 0, 1, 1 \rangle$$

# A NEED TO DEFINE ARBITRARY FEATURES

**Example 1** Parallelism in coordination [Johnson et. al 1999]

**Constituents with similar structure tend to be coordinated**

$h(x)$  = Number of instances of non-parallel coordination in  $x$

Bars in New York and pubs in London	$h(x) = 0$
vs. Bars in New York and pubs	$h(x) = 1$

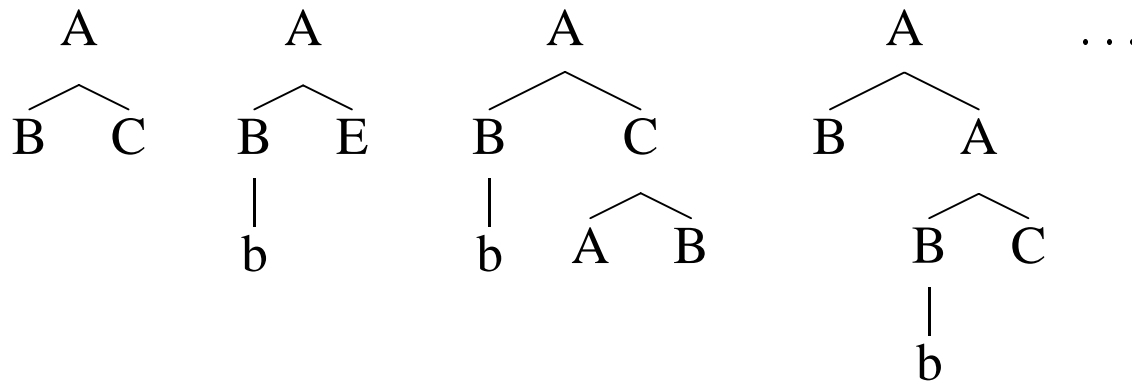
**Example 2** Semantic features

$h(x)$  = Number of times **pour** takes an object with a +liquid feature

pour the <b>cappucino</b>	$h(x) = 1$
vs. pour the <b>book</b>	$h(x) = 0$

# “ALL SUBTREES” REPRESENTATION

- Given: Non-Terminal symbols  $\{A, B, \dots\}$   
Terminal symbols  $\{a, b, c, \dots\}$
- An infinite set of subtrees



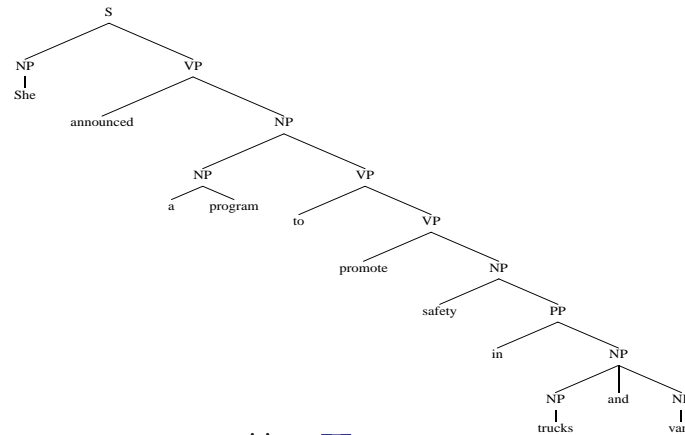
- An infinite set of features, e.g.,

$h_3(x) =$  Number of times  is seen in  $x$

Representation studied by [Bod 98]

## COMPONENT 3: $\mathbf{W}$

- $\mathbf{W}$  is a **parameter vector**  $\in \mathbb{R}^d$
  - $\Phi$  and  $\mathbf{W}$  together map a candidate to a real-valued score
- 



$\Downarrow \Phi$

$\langle 1, 0, 2, 0, 0, 15, 5 \rangle$

$\Downarrow \Phi \cdot \mathbf{W}$

$$\langle 1, 0, 2, 0, 0, 15, 5 \rangle \cdot \langle 1.9, -0.3, 0.2, 1.3, 0, 1.0, -2.3 \rangle = 5.8$$

## PUTTING IT ALL TOGETHER

- $\mathcal{X}$  is set of sentences,  $\mathcal{Y}$  is set of possible outputs (e.g. trees)
- Need to learn a function  $\mathbf{F} : \mathcal{X} \rightarrow \mathcal{Y}$
- $\mathbf{GEN}$ ,  $\Phi$ ,  $\mathbf{W}$  define

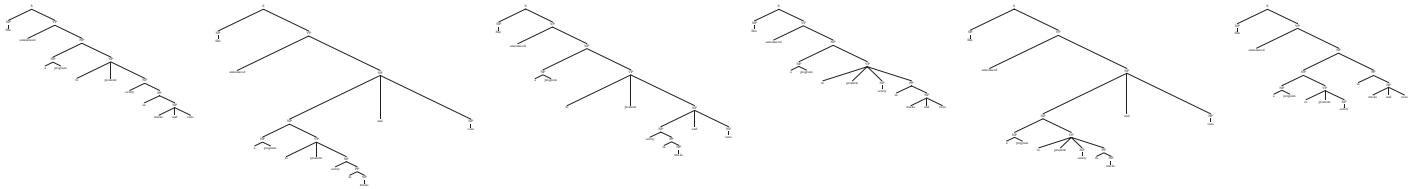
$$\mathbf{F}(x) = \arg \max_{y \in \mathbf{GEN}(x)} \Phi(y) \cdot \mathbf{W}$$

**Choose the highest scoring tree as the most plausible structure**

- Given examples  $(x_i, y_i)$ , how to set  $\mathbf{W}$ ?

She announced a program to promote safety in trucks and vans

⇓ GEN



⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⟨1, 1, 3, 5⟩

⟨2, 0, 0, 5⟩

⟨1, 0, 1, 5⟩

⟨0, 0, 3, 0⟩

⟨0, 1, 0, 5⟩

⟨0, 0, 1, 5⟩

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

13.6

12.2

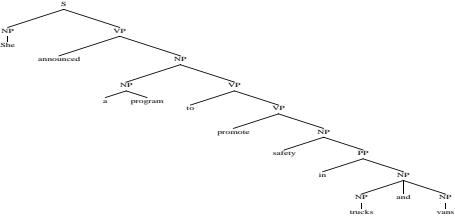
12.1

3.3

9.4

11.1

⇓ arg max



# PREVIOUS WORK ON PARAMETER ESTIMATION

- **Markov Random Fields**

[Johnson et. al 1999],[Lafferty et al. 2001],

[Ratnaparkhi, Roukos and Ward 94]

$$P(y_i | x_i, \mathbf{W}) = \frac{e^{\Phi(y_i) \cdot \mathbf{W}}}{\sum_{y \in \mathbf{GEN}(x_i)} e^{\Phi(y) \cdot \mathbf{W}}}$$

- **Boosting** (ICML 2000 paper):

Minimize: 
$$\sum_i \sum_{z \in \mathbf{GEN}(x_i), z \neq y_i} e^{\mathbf{W} \cdot \Phi(z) - \mathbf{W} \cdot \Phi(y_i)}$$

- Close connections between the two:

[Friedman, Hastie, and Tibshirani 1998, Lafferty 1999]

(COLT 2000 paper with Rob Schapire, Yoram Singer)

# A VARIANT OF THE PERCEPTRON ALGORITHM

**Inputs:** Training set  $(x_i, y_i)$  for  $i = 1 \dots n$

**Initialization:**  $\mathbf{W} = 0$

**Define:**  $\mathbf{F}(x) = \operatorname{argmax}_{y \in \mathbf{GEN}(x)} \Phi(y) \cdot \mathbf{W}$

**Algorithm:** For  $t = 1 \dots T, i = 1 \dots n$   
 $z_i = \mathbf{F}(x_i)$   
If  $(z_i \neq y_i)$   $\mathbf{W} = \mathbf{W} + \Phi(y_i) - \Phi(z_i)$

**Output:** Parameters  $\mathbf{W}$

## THEORY UNDERLYING THE ALGORITHM

- **Definition:**  $\overline{\text{GEN}}(x_i) = \text{GEN}(x_i) - \{y_i\}$
- **Definition:** The training set is **separable with margin  $\delta$** , if there is a vector  $\mathbf{U} \in \mathbb{R}^d$  with  $\|\mathbf{U}\| = 1$  such that

$$\forall i, \forall z \in \overline{\text{GEN}}(x_i) \quad \mathbf{U} \cdot \Phi(y_i) - \mathbf{U} \cdot \Phi(z) \geq \delta$$

## THEORY UNDERLYING THE ALGORITHM

- **Definition:**  $\overline{\mathbf{GEN}}(x_i) = \mathbf{GEN}(x_i) - \{y_i\}$
- **Definition:** The training set is **separable with margin  $\delta$** , if there is a vector  $\mathbf{U} \in \mathbb{R}^d$  with  $\|\mathbf{U}\| = 1$  such that

$$\forall i, \forall z \in \overline{\mathbf{GEN}}(x_i) \quad \mathbf{U} \cdot \Phi(y_i) - \mathbf{U} \cdot \Phi(z) \geq \delta$$

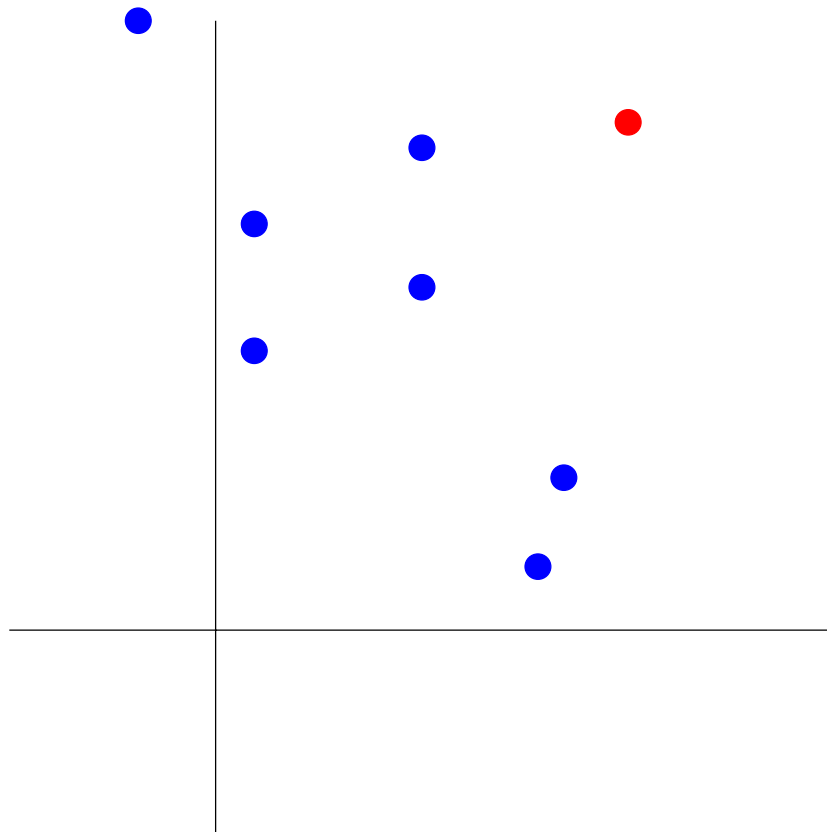
---

**Theorem:** For any training sequence  $(x_i, y_i)$  which is separable with margin  $\delta$ , then for the perceptron algorithm

$$\text{Number of mistakes} \leq \frac{R^2}{\delta^2}$$

where  $R$  is a constant such that  $\forall i, \forall z \in \overline{\mathbf{GEN}}(x_i) \quad \|\Phi(y_i) - \Phi(z)\| \leq R$

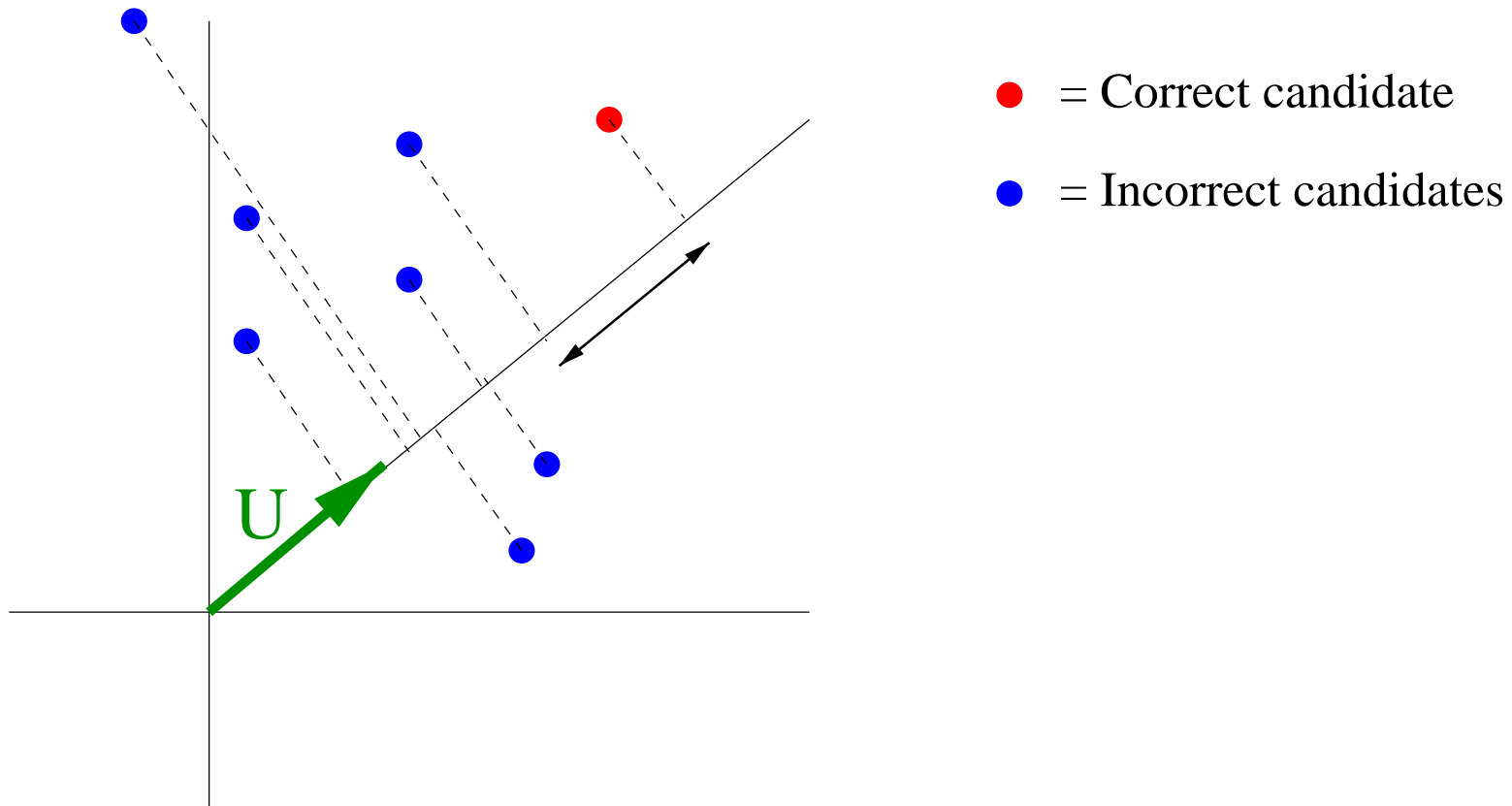
# GEOMETRIC INTUITION BEHIND SEPARATION



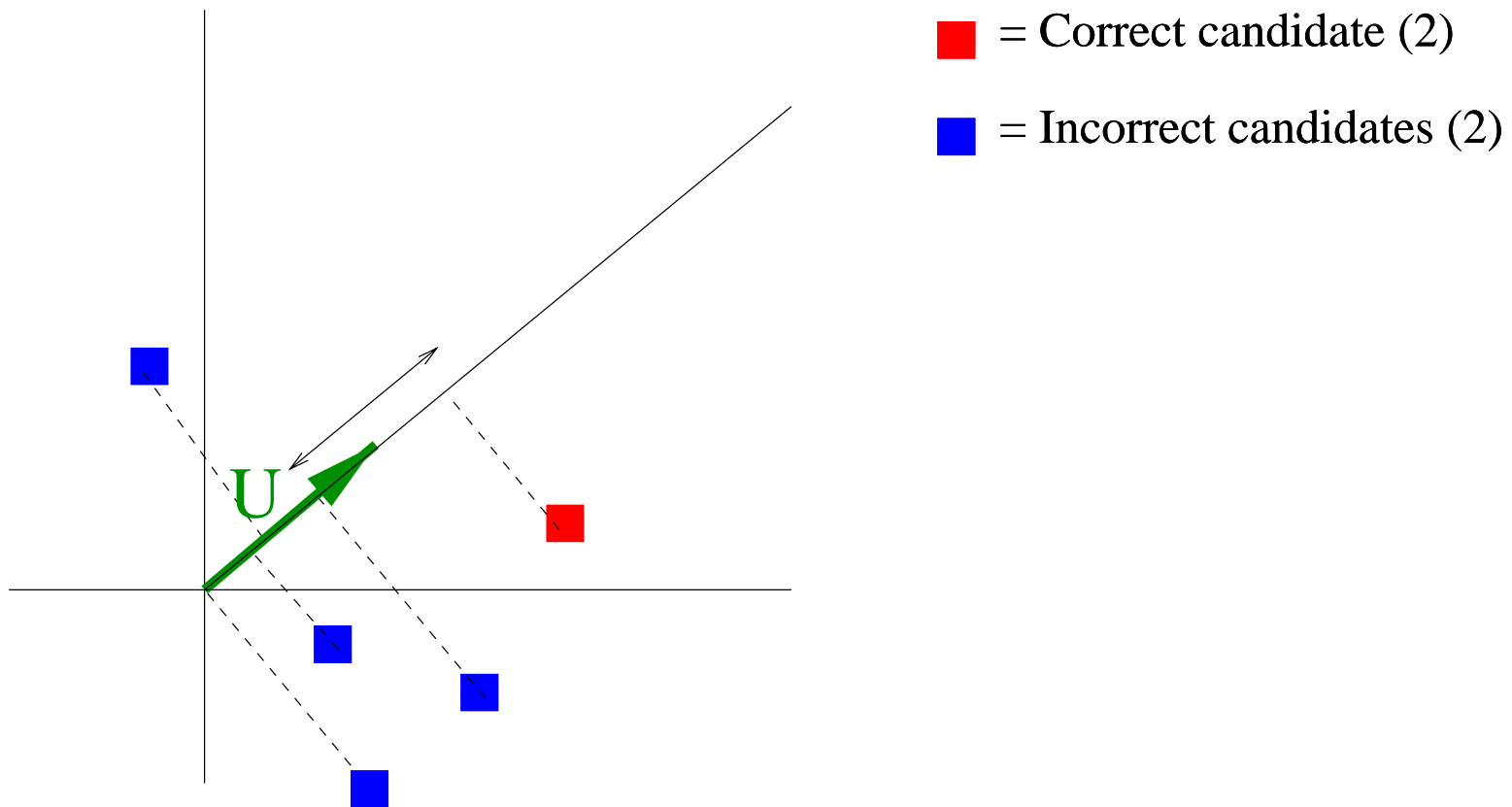
● = Correct candidate

● = Incorrect candidates

# GEOMETRIC INTUITION BEHIND SEPARATION



# ALL EXAMPLES ARE SEPARATED



# THEORY UNDERLYING THE ALGORITHM

---

## Old definition

- Classification case: examples  $(x_i, y_i)$ , where  $y_i \in \{-1, +1\}$
- **Definition:** The training set is **separable with margin  $\delta$** , if there is a vector  $\mathbf{U} \in \mathbb{R}^d$  with  $\|\mathbf{U}\| = 1$  such that

$$\forall i, \quad \mathbf{U} \cdot \Phi(x_i) \times y_i \geq \delta$$

---

## New definitions

- **Definition:**  $\overline{\text{GEN}}(x_i) = \text{GEN}(x_i) - \{y_i\}$
- **Definition:** The training set is **separable with margin  $\delta$** , if there is a vector  $\mathbf{U} \in \mathbb{R}^d$  with  $\|\mathbf{U}\| = 1$  such that

$$\forall i, \forall z \in \overline{\text{GEN}}(x_i) \quad \mathbf{U} \cdot \Phi(y_i) - \mathbf{U} \cdot \Phi(z) \geq \delta$$

## THEORY UNDERLYING THE ALGORITHM

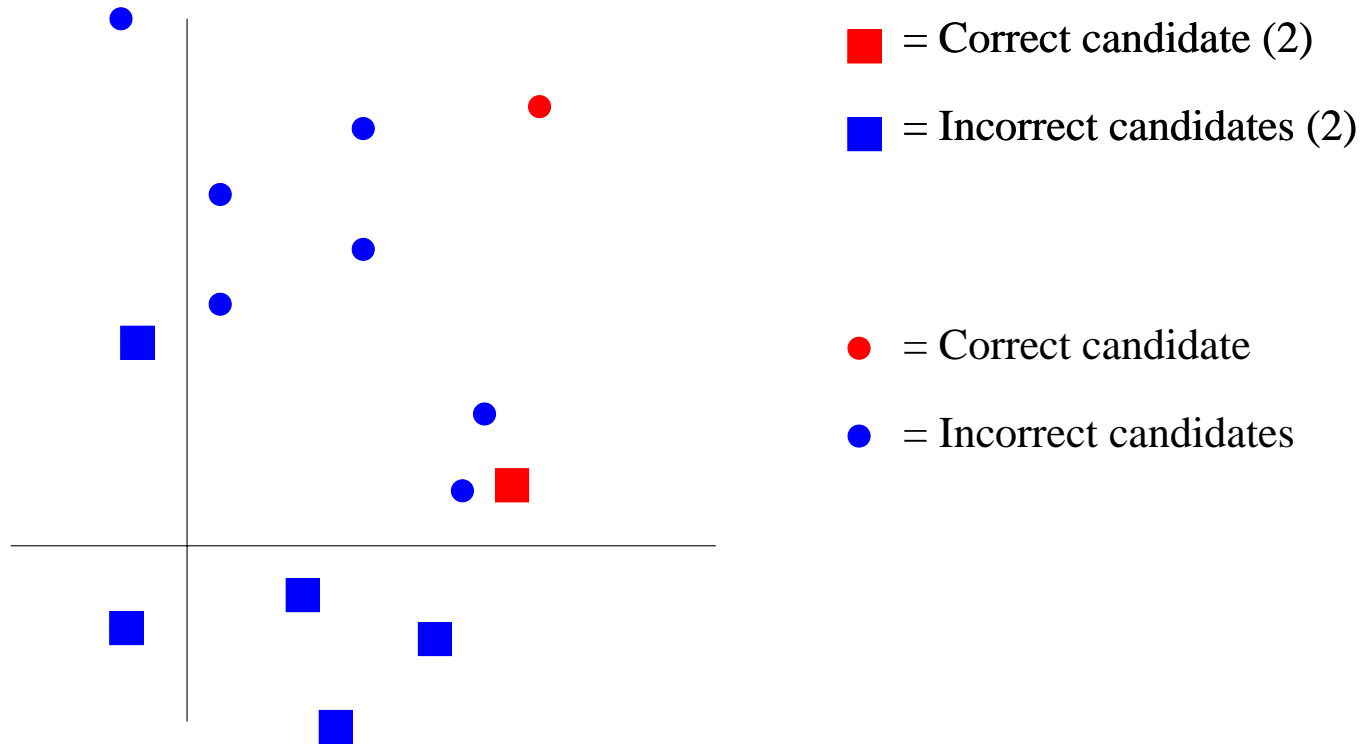
**Theorem:** For any training sequence  $(x_i, y_i)$  which is separable with margin  $\delta$ , then for the perceptron algorithm

$$\text{Number of mistakes} \leq \frac{R^2}{\delta^2}$$

where  $R$  is a constant such that  $\forall i, \forall z \in \overline{\text{GEN}}(x_i) \quad \|\Phi(y_i) - \Phi(z)\| \leq R$

**Proof:** Direct modification of the proof for the classification case.

# WHAT WE'RE NOT DOING



- We aim to separate the correct structure from its competitors
- We're **not** aiming to do **classification** of good vs. bad parses  
This is a harder (and unnecessary) problem

## Proof:

Let  $\mathbf{W}^k$  be the weights before the  $k$ 'th mistake.  $\mathbf{W}^1 = 0$

If the  $k$ 'th mistake is made at  $i$ 'th example,

and  $z_i = \operatorname{argmax}_{y \in \mathbf{GEN}(x_i)} \Phi(y) \cdot \mathbf{W}^k$ , then

$$\begin{aligned}\mathbf{W}^{k+1} &= \mathbf{W}^k + \Phi(y_i) - \Phi(z_i) \\ \Rightarrow \mathbf{U} \cdot \mathbf{W}^{k+1} &= \mathbf{U} \cdot \mathbf{W}^k + \mathbf{U} \cdot \Phi(y_i) - \mathbf{U} \cdot \Phi(z_i) \\ &\geq \mathbf{U} \cdot \mathbf{W}^k + \delta \\ &\geq k\delta \\ \Rightarrow \|\mathbf{W}^{k+1}\| &\geq k\delta\end{aligned}$$

Also,

$$\begin{aligned}\|\mathbf{W}^{k+1}\|^2 &= \|\mathbf{W}^k\|^2 + \|\Phi(y_i) - \Phi(z_i)\|^2 + 2\mathbf{W}^k \cdot (\Phi(y_i) - \Phi(z_i)) \\ &\leq \|\mathbf{W}^k\|^2 + R^2 \\ \Rightarrow \|\mathbf{W}^{k+1}\|^2 &\leq kR^2 \\ \Rightarrow k^2\delta^2 &\leq \|\mathbf{W}^{k+1}\|^2 \leq kR^2 \\ \Rightarrow k &\leq R^2/\delta^2\end{aligned}$$

## MORE THEORY FOR THE PERCEPTRON ALGORITHM

- **Question 1: what if the data is not separable?**  
[Freund and Schapire 99] give a modified theorem for this case
- **Question 2: performance on training data is all very well, but what about performance on new test examples?**

Assume some distribution  $P(x, y)$  underlying examples

**Theorem** [Helmbold and Warmuth 95]: For any distribution  $P(x, y)$  generating examples, if  $e =$  expected number of mistakes of an online algorithm on a sequence of  $m + 1$  examples, then a randomized algorithm trained on  $m$  samples will have probability  $\frac{e}{m+1}$  of making an error on a newly drawn example from  $P$ .

[Freund and Schapire 99] use this to define the **Voted Perceptron**

## QUESTION 1: WHAT IF THE DATA IS NOT SEPARABLE?

- A modified theorem follows from [Freund and Schapire 99]

For any  $\mathbf{U}$ ,  $\delta$  pair with  $\|\mathbf{U}\| = 1$  define

- $m_i = \mathbf{U} \cdot \Phi(y_i) - \max_{z \in \overline{\text{GEN}}(x_i)} \mathbf{U} \cdot \Phi(z)$

- $\epsilon_i = \begin{cases} \delta - m_i & \text{if } m_i < \delta \\ 0 & \text{otherwise} \end{cases}$

- $D(\mathbf{U}, \delta) = \sqrt{\sum_i \epsilon_i^2}$

Then it can be shown that

$$\text{Number of mistakes} \leq \min_{\mathbf{U}, \delta} \frac{(R + D(\mathbf{U}, \delta))^2}{\delta^2}$$

## QUESTION 2: PERFORMANCE ON NEW EXAMPLES?

**Theorem** [[Helmbold and Warmuth 95](#)]: For any distribution  $P(x, y)$  generating examples, if  $e =$  expected number of mistakes of an online algorithm on a sequence of  $m + 1$  examples, then a randomized algorithm trained on  $m$  samples will have probability  $\frac{e}{m+1}$  of making an error on a newly drawn test example.

[[Freund and Schapire 99](#)]: For the **voted perceptron**

$$\text{Probability of error} \leq \frac{2}{m+1} E \left[ \min_{\mathbf{U}, \delta} \frac{(R + D(\mathbf{U}, \delta))^2}{\delta^2} \right]$$

where  $E[\dots]$  is the expectation taken over a sample of  $(m + 1)$  examples drawn from  $P(x, y)$

## QUESTION 3: ARE LINEAR SEPARATORS GOOD ENOUGH?

Example: Hand-written digit recognition

- **Kernels** give non-linear separators

Linear model + no voting  $\Rightarrow$  13.5% error rate

Degree 3 polynomial kernel + voting  $\Rightarrow$  1.6% error rate

- **Voting** helps in non-separable cases

Linear model + voting  $\Rightarrow$  8.1% error rate

## THREE SCENARIOS

- Reranking output from a probabilistic model
- Reranking output using **kernels**
- Training a tagging model

## RERANKING APPROACHES

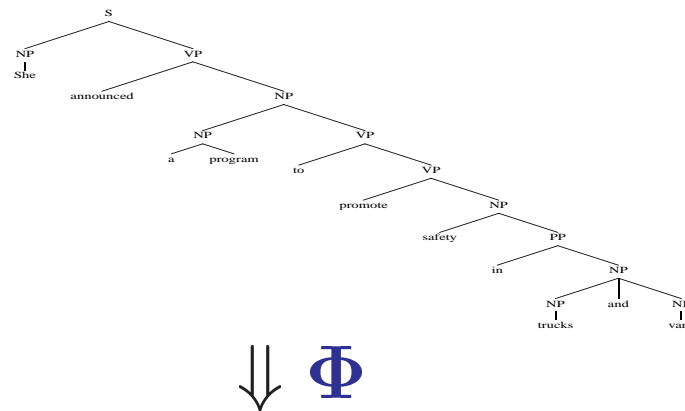
- **GEN** is the top  $n$  most probable candidates from a **base model**
  - Parsing: a lexicalized probabilistic context-free grammar
  - Tagging: “maximum entropy” tagger
  - Speech recognition: existing recogniser

# PARSING EXPERIMENTS

**GEN** Beam search used to parse training and test sentences:  
around 27 parses for each sentence

$\Phi = \langle L(x), h_1(x) \dots h_m(x) \rangle$ , where  $L(x) = \log$ -likelihood from first-pass parser,  $h_1 \dots h_m$  are  $\approx 500,000$  indicator functions

$$e.g., \quad h_1(x) = \begin{cases} 1 & \text{if } x \text{ contains } \langle S \rightarrow NP \ VP \rangle \\ 0 & \text{otherwise} \end{cases}$$



$\langle -15.65, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, \dots, 1, 0, 0 \rangle$

# NAMED ENTITIES

**GEN** Top 20 segmentations from a “maximum-entropy” tagger

$$\Phi = \langle L(x), h_1(x) \dots h_m(x) \rangle,$$

$$e.g., \quad h_1(x) = \begin{cases} 1 & \text{if } x \text{ contains a boundary} = \boxed{\text{“[The} \\ 0 & \text{otherwise} \end{cases}$$

---

Whether you’re an aging flower child or a clueless **[Gen-Xer]**, “**[The Day They Shot John Lennon]**,” playing at the **[Dougherty Arts Center]**, entertains the imagination.

⇓  $\Phi$

$$\langle -3.17, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, \dots 0, 1, 1 \rangle$$

Whether you're an aging flower child or a clueless  
[Gen-Xer], “[The Day They Shot John Lennon],” playing at the  
[Dougherty Arts Center], entertains the imagination.

⇓  $\Phi$

$\langle -3.17, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, \dots 0, 1, 1 \rangle$

---

Whether you're an aging flower child or a clueless  
Gen-Xer, “The Day [They Shot John Lennon],” playing at the  
[Dougherty Arts Center], entertains the imagination.

⇓  $\Phi$

$\langle -3.51, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, \dots 0, 1, 0 \rangle$

---

Whether you're an aging flower child or a clueless  
[Gen-Xer], “The Day [They Shot John Lennon],” playing at the  
[Dougherty Arts Center], entertains the imagination.

⇓  $\Phi$

$\langle -2.87, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, \dots 0, 1, 0 \rangle$

# EXPERIMENTS

## **Parsing Wall Street Journal Treebank**

Training set = 40,000 sentences, test = 2,416 sentences

State-of-the-art parser: 88.2% F-measure

Reranked model: 89.5% F-measure (**11% relative error reduction**)

Boosting: 89.7% F-measure (**13% relative error reduction**)

## **Recovering Named-Entities in Web Data**

Training data = 53,609 sentences (1,047,491 words),

test data = 14,717 sentences (291,898 words)

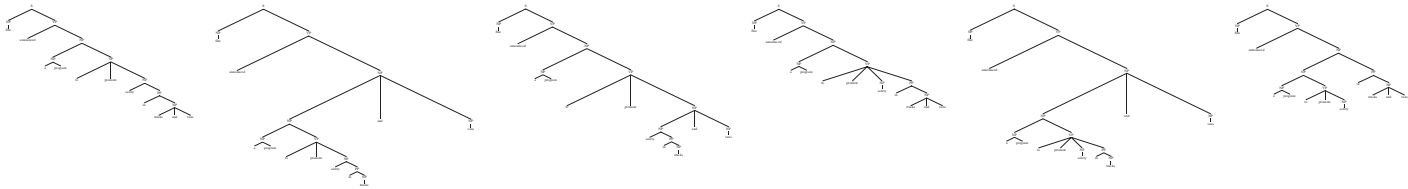
State-of-the-art tagger: 85.3% F-measure

Reranked model: 87.9% F-measure (**17.7% relative error reduction**)

Boosting: 87.6% F-measure (**15.6% relative error reduction**)

She announced a program to promote safety in trucks and vans

⇓ GEN



⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⇓ Φ

⟨1, 1, 3, 5⟩

⟨2, 0, 0, 5⟩

⟨1, 0, 1, 5⟩

⟨0, 0, 3, 0⟩

⟨0, 1, 0, 5⟩

⟨0, 0, 1, 5⟩

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

⇓ Φ · W

13.6

12.2

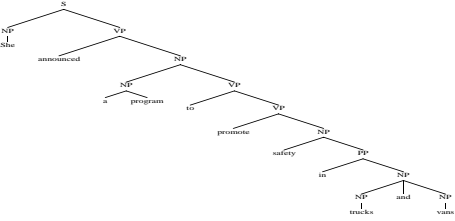
12.1

3.3

9.4

11.1

⇓ arg max



## THREE SCENARIOS

- Reranking output from a probabilistic model
- Reranking output using **kernels**  
(with Nigel Duffy)
- Training a weighted grammar

# A VARIANT OF THE PERCEPTRON ALGORITHM

**Inputs:** Training set  $(x_i, y_i)$  for  $i = 1 \dots n$

**Initialization:**  $\mathbf{W} = 0$

**Define:**  $\mathbf{F}(x) = \operatorname{argmax}_{y \in \mathbf{GEN}(x)} \Phi(y) \cdot \mathbf{W}$

**Algorithm:** For  $t = 1 \dots T, i = 1 \dots n$   
 $z_i = \mathbf{F}(x_i)$   
If  $(z_i \neq y_i)$   $\mathbf{W} = \mathbf{W} + \Phi(y_i) - \Phi(z_i)$

**Output:** Parameters  $\mathbf{W}$

# “KERNEL” FORM OF THE PERCEPTRON ALGORITHM

- Manipulate sets  $\mathcal{A}$  and  $\mathcal{B}$  of training examples such that

$$\mathbf{W} = \sum_{a \in \mathcal{A}} \Phi(a) - \sum_{b \in \mathcal{B}} \Phi(b)$$

$$\mathbf{W} \cdot \Phi(y) = \left( \sum_{a \in \mathcal{A}} \Phi(y) \cdot \Phi(a) - \sum_{b \in \mathcal{B}} \Phi(y) \cdot \Phi(b) \right)$$

- Intuition:
  - $\mathcal{A}$  is a set of “good parses”
  - $\mathcal{B}$  is a set of “bad parses”
  - $\Phi(y) \cdot \Phi(x)$  is a similarity measure

# “KERNEL” FORM OF THE PERCEPTRON ALGORITHM

**Define:**

$$\mathbf{F}(x) = \arg \max_{y \in \mathbf{GEN}(x)} \left( \sum_{a \in \mathcal{A}} \Phi(y) \cdot \Phi(a) - \sum_{b \in \mathcal{B}} \Phi(y) \cdot \Phi(b) \right)$$

**Initialization:**

$$\mathcal{A} = \{\}, \mathcal{B} = \{\}$$

**Algorithm:**

For  $t = 1 \dots T, i = 1 \dots n$

$$z_i = \mathbf{F}(x_i)$$

$$\text{If } (z_i \neq y_i) \quad \begin{aligned} \mathcal{A} &= \mathcal{A} \cup \{y_i\}, \\ \mathcal{B} &= \mathcal{B} \cup \{z_i\} \end{aligned}$$

**Output:**

Sets  $\mathcal{A}$  and  $\mathcal{B}$

---

Equivalence:  $\mathbf{W} = \sum_{a \in \mathcal{A}} \Phi(a) - \sum_{b \in \mathcal{B}} \Phi(b)$

$$\text{Equivalence: } \mathbf{W} = \sum_{a \in \mathcal{A}} \Phi(a) - \sum_{b \in \mathcal{B}} \Phi(b)$$

---

## Original Form

**Initialization:**  $\mathbf{W} = 0$

**Define:**  $\mathbf{F}(x) = \operatorname{argmax}_{y \in \mathbf{GEN}(x)} \Phi(y) \cdot \mathbf{W}$

**Algorithm:** For  $t = 1 \dots T, i = 1 \dots n$   
 $z_i = \mathbf{F}(x_i)$   
If  $(z_i \neq y_i)$   $\mathbf{W} = \mathbf{W} + \Phi(y_i) - \Phi(z_i)$

---

## Dual Form

**Initialization:**  $\mathcal{A} = \{\}, \mathcal{B} = \{\}$

**Define:**  $\mathbf{F}(x) = \operatorname{argmax}_{y \in \mathbf{GEN}(x)} \left( \sum_{a \in \mathcal{A}} \Phi(y) \cdot \Phi(a) - \sum_{b \in \mathcal{B}} \Phi(y) \cdot \Phi(b) \right)$

**Algorithm:** For  $t = 1 \dots T, i = 1 \dots n$   
 $z_i = \mathbf{F}(x_i)$   
If  $(z_i \neq y_i)$   $\mathcal{A} = \mathcal{A} \cup \{y_i\}, \mathcal{B} = \mathcal{B} \cup \{z_i\}$

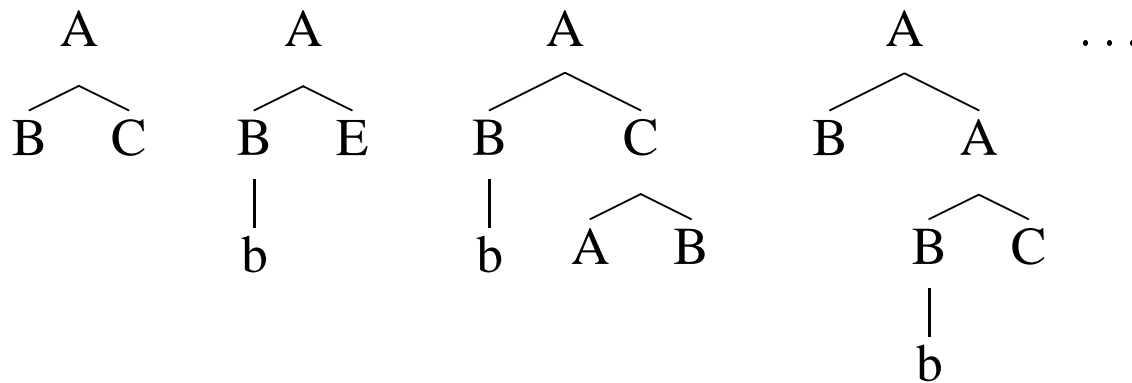
## MOTIVATION FOR THE KERNEL FORM

- Let  $N = \sum_i |\mathbf{GEN}(x_i)|$   
(i.e., size of the training set)
- One pass over training set takes  $O(Nd)$  time for regular form
- Say it takes  $I$  time to calculate  $\Phi(x) \cdot \Phi(y)$ , perceptron makes  $M$  mistakes, then dual form takes  $O(NMI)$  time
- For some representations,  $MI \ll d$

**If we can compute  $\Phi(x) \cdot \Phi(y)$  efficiently  
we can learn efficiently with the representation  $\Phi$**

## “ALL SUBTREES” REPRESENTATION

- Given: Non-Terminal symbols  $\{A, B, \dots\}$   
Terminal symbols  $\{a, b, c, \dots\}$
- An infinite set of subtrees



- **Step 1:**

Choose an (arbitrary) mapping from subtrees to integers

$h_i(x)$  = Number of times subtree  $i$  is seen in  $x$

$\Phi(x) = \langle h_1(x), h_2(x), h_3(x), \dots \rangle$

## ALL SUBTREES REPRESENTATION

- $\Phi$  is now huge
- **But** inner product  $\Phi(T_1) \cdot \Phi(T_2)$  can be computed efficiently using dynamic programming.

## COMPUTING THE INNER PRODUCT

Define –  $N_1$  and  $N_2$  are sets of nodes in  $T_1$  and  $T_2$  respectively.

$$- I_i(x) = \begin{cases} 1 & \text{if } i\text{'th subtree is rooted at } x. \\ 0 & \text{otherwise.} \end{cases}$$

Follows that:

$$h_i(T_1) = \sum_{n_1 \in N_1} I_i(n_1) \quad \text{and} \quad h_i(T_2) = \sum_{n_2 \in N_2} I_i(n_2)$$

$$\begin{aligned} \Phi(T_1) \cdot \Phi(T_2) &= \sum_i h_i(T_1) h_i(T_2) = \sum_i \left( \sum_{n_1 \in N_1} I_i(n_1) \right) \left( \sum_{n_2 \in N_2} I_i(n_2) \right) \\ &= \sum_{n_1 \in N_1} \sum_{n_2 \in N_2} \sum_i I_i(n_1) I_i(n_2) \\ &= \sum_{n_1 \in N_1} \sum_{n_2 \in N_2} \Delta(n_1, n_2) \end{aligned}$$

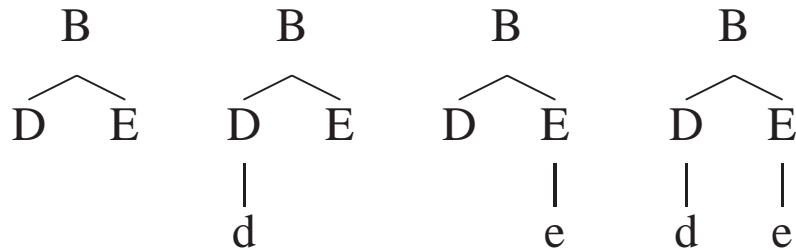
where  $\Delta(n_1, n_2) = \sum_i I_i(n_1) I_i(n_2)$  is the **number of common subtrees at  $n_1, n_2$**

# AN EXAMPLE



$$\Phi(T_1) \cdot \Phi(T_2) = \Delta(A, A) + \Delta(A, B) \dots + \Delta(B, A) + \Delta(B, B) \dots + \Delta(G, G)$$

- Most of these terms are 0 (e.g.  $\Delta(A, B)$ ).
- Some are non-zero, e.g.  $\Delta(B, B) = 4$



## RECURSIVE DEFINITION OF $\Delta(n_1, n_2)$

- If the productions at  $n_1$  and  $n_2$  are different

$$\Delta(n_1, n_2) = 0$$

- Else if  $n_1, n_2$  are pre-terminals,

$$\Delta(n_1, n_2) = 1$$

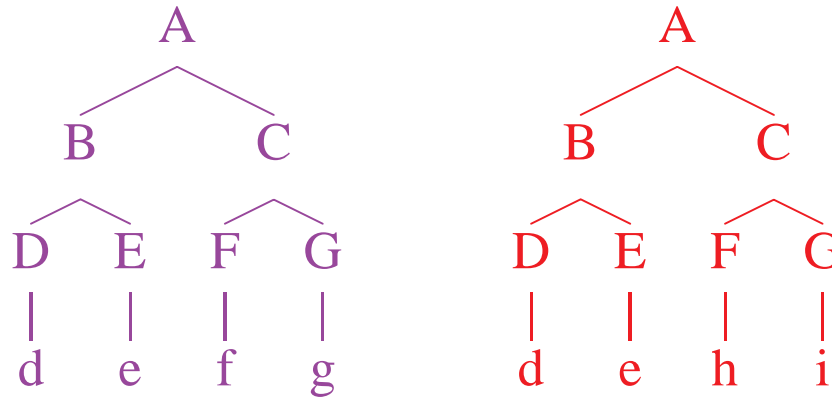
- Else

$$\Delta(n_1, n_2) = \prod_{j=1}^{nc(n_1)} (1 + \Delta(ch(n_1, j), ch(n_2, j)))$$

$nc(n_1)$  is number of children of node  $n_1$ ;

$ch(n_1, j)$  is the  $j$ 'th child of  $n_1$ .

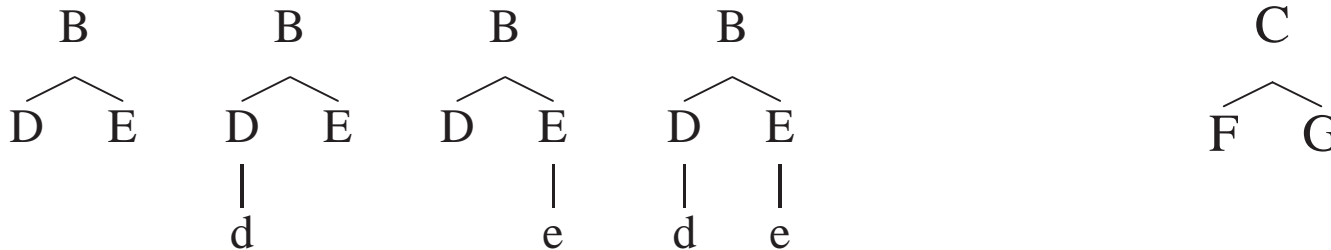
# ILLUSTRATION OF THE RECURSION



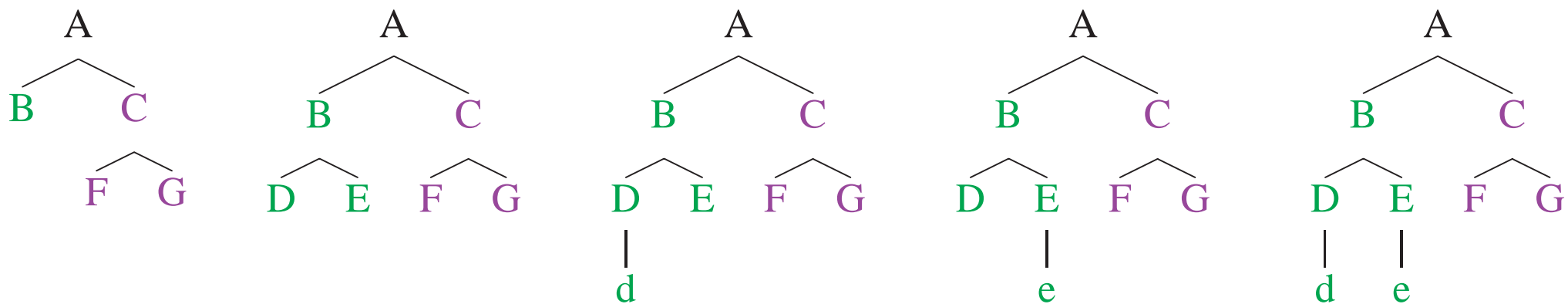
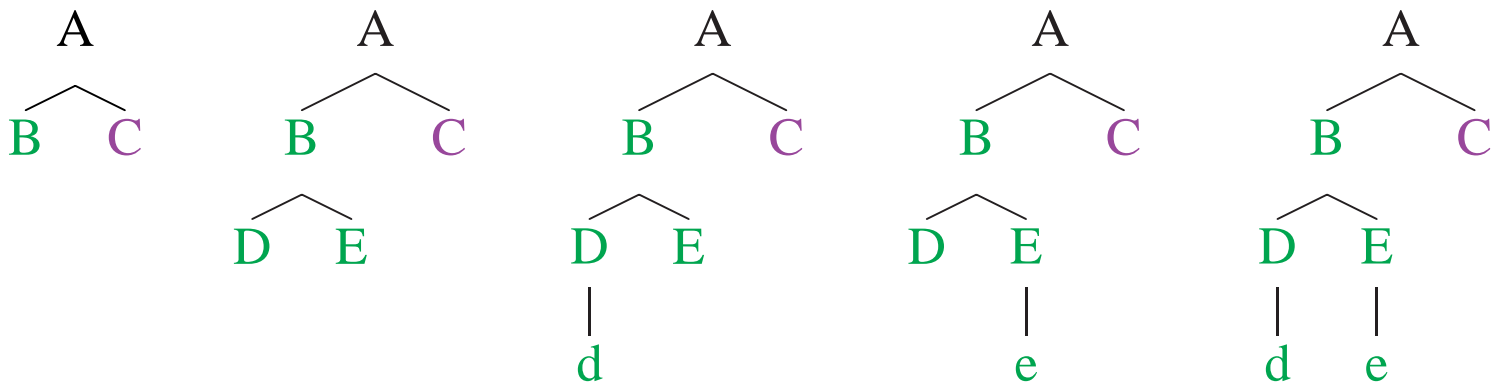
How many subtrees do nodes  $A$  and  $A$  have in common? i.e., What is  $\Delta(A, A)$ ?

$$\Delta(B, B) = 4$$

$$\Delta(C, C) = 1$$

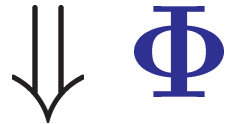


$$\Delta(A, A) = (\Delta(B, B) + 1) \times (\Delta(C, C) + 1) = 10$$



# SIMILAR KERNELS CAN BE DERIVED FOR TAGGED SEQUENCES

Whether you're an aging flower child or a clueless [Gen-Xer], "[The Day They Shot John Lennon]," playing at the [Dougherty Arts Center], entertains the imagination.



Whether [Gen-Xer], Day They John Lennon],” playing

Whether you're an aging flower child or a clueless [Gen

...

# EXPERIMENTS

## **Parsing Wall Street Journal Treebank**

Training set = 40,000 sentences, test = 2,416 sentences

State-of-the-art parser: 88.5% F-measure

Reranked model: 89.1% F-measure

(5% relative error reduction)

## **Recovering Named-Entities in Web Data**

Training data = 53,609 sentences (1,047,491 words),

test data = 14,717 sentences (291,898 words)

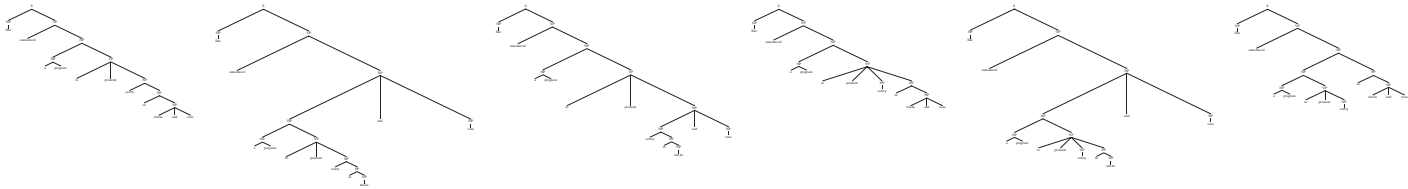
State-of-the-art tagger: 85.3% F-measure

Reranked model: 87.6% F-measure

(15.6% relative error reduction)

She announced a program to promote safety in trucks and vans

⇓ GEN



⇓ Φ

⟨1, 1, 3, 5⟩

⇓ Φ

⟨2, 0, 0, 5⟩

⇓ Φ

⟨1, 0, 1, 5⟩

⇓ Φ

⟨0, 0, 3, 0⟩

⇓ Φ

⟨0, 1, 0, 5⟩

⇓ Φ

⟨0, 0, 1, 5⟩

⇓ Φ · W

13.6

⇓ Φ · W

12.2

⇓ Φ · W

12.1

⇓ Φ · W

3.3

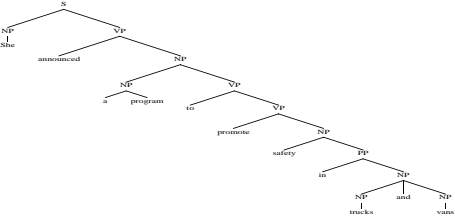
⇓ Φ · W

9.4

⇓ Φ · W

11.1

⇓ arg max



## THREE SCENARIOS

- Reranking output from a probabilistic model
- Reranking output using **kernels**
- Training a weighted grammar

# TRAINING WEIGHTED AUTOMATA

- Main computational expense:

$$\mathbf{F}(x) = \arg \max_{y \in \mathbf{GEN}(x)} \Phi(y) \cdot \mathbf{W}$$

- Worst case:  $O(|\mathbf{GEN}(x)|)$  time
- Can be better for some  $\mathbf{GEN}$ ,  $\Phi$  combinations

## THE TAGGING PROBLEM

Hispaniola/**NNP** quickly/**RB** became/**VB** an/**DT**  
important/**JJ** base/**??** from which Spain expanded  
its empire into the rest of the Western Hemisphere .

- There are many possible tags in the position **??**
- Need to learn a function from (context, tag) pairs to a real value indicating the “plausibility” of the tag in this context

## REPRESENTATION: HISTORIES

- A **history** is a 4-tuple  $\langle t_{-1}, t_{-2}, w_{[1:n]}, i \rangle$
  - $t_{-1}, t_{-2}$  are the previous two tags.
  - $w_{[1:n]}$  are the  $n$  words in the input sentence.
  - $i$  is the index of the word being tagged
- 

Hispaniola/**NNP** quickly/**RB** became/**VB** an/**DT** important/**JJ**  
base/**??** from which Spain expanded its empire into the rest of the  
Western Hemisphere .

- $t_{-1}, t_{-2} = \text{DT, JJ}$
- $w_{[1:n]} = \langle \text{Hispaniola, quickly, became, } \dots, \text{ Hemisphere, .} \rangle$
- $i = 6$

## FEATURE-VECTOR REPRESENTATIONS

- Take a history/tag pair  $(h, t)$ .
  - $\phi_s(h, t)$  for  $s = 1 \dots d$  are **features** representing tagging decision  $t$  in context  $h$ .
- 

Example: POS Tagging [Ratnaparkhi 96]

- **Word/tag features**

$$\phi_{100}(h, t) = \begin{cases} 1 & \text{if current word } w_i \text{ is base and } t = \text{VB} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_{101}(h, t) = \begin{cases} 1 & \text{if current word } w_i \text{ ends in ing and } t = \text{VBG} \\ 0 & \text{otherwise} \end{cases}$$

- **Contextual Features**

$$\phi_{103}(h, t) = \begin{cases} 1 & \text{if } \langle t_{-2}, t_{-1}, t \rangle = \langle \text{DT}, \text{JJ}, \text{VB} \rangle \\ 0 & \text{otherwise} \end{cases}$$

## GLOBAL FEATURES FROM LOCAL FEATURES

- **Global features**  $\Phi_s(w_{[1:n]}, t_{[1:n]})$  for  $s = 1 \dots d$  map an entire  $(w_{[1:n]}, t_{[1:n]})$  pair to a feature vector
- Global features can be derived from local features

$$\Phi_s(w_{[1:n]}, t_{[1:n]}) = \sum_i \phi_s(h_i, t_i)$$

where  $h_i = \langle t_{i-2}, t_{i-1}, w_{[1:n]}, i \rangle$

- Local features are indicator functions  $\Rightarrow$   
Global features are **counts**

$$\phi_{103}(h, t) = \begin{cases} 1 & \text{if } \langle t_{-1}, t \rangle = \langle \text{DT}, \text{NN} \rangle \\ 0 & \text{otherwise} \end{cases}$$

$$\Phi_{103}(w_{[1:n]}, t_{[1:n]}) = \text{Number of times } \langle \text{DT}, \text{NN} \rangle \text{ is seen in } t_{[1:n]}$$

e.g.,  $\Phi_{103}(\text{the/DT man/NN bit/VBD the/DT dog/NN}) = 2$

## MAXIMUM ENTROPY MODELS

- $\phi_s(h, t)$  for  $s = 1 \dots d$  are **features**  
 $\mathbf{W}_s$  for  $s = 1 \dots d$  are **parameters**

- Conditional distribution:

$$P(t|h) = \frac{e^{\sum_s \mathbf{W}_s \phi_s(h,t)}}{Z(h, \mathbf{W})}$$

where  $Z(h, \mathbf{W}) = \sum_{t' \in \mathcal{T}} e^{\sum_s \mathbf{W}_s \phi_s(h,t')}$

- Parameters estimated using maximum-likelihood  
e.g., iterative scaling, gradient descent

# MAXIMUM ENTROPY MODELS

- Word sequence  $w_{[1:n]} = [w_1, w_2 \dots w_n]$
- Tag sequence  $t_{[1:n]} = [t_1, t_2 \dots t_n]$
- Histories  $h_i = \langle t_{i-1}, t_{i-2}, w_{[1:n]}, i \rangle$

$$\log P(t_{[1:n]} \mid w_{[1:n]})$$

$$= \sum_{i=1}^n \log P(t_i \mid h_i)$$

$$= \underbrace{\sum_{i=1}^n \sum_s \mathbf{W}_s \phi_s(h_i, t_i)}_{\text{Linear Score}} - \underbrace{\sum_{i=1}^n \log Z(h_i, \mathbf{W})}_{\text{Local Normalization Terms}}$$

Linear Score

Local Normalization  
Terms

# PROBLEMS WITH LOCALLY NORMALIZED MODELS

- “Label bias” problem [[Lafferty et al. 2001](#)]  
See also [[Klein and Manning 2002](#)]

- Example of a conditional distribution that locally normalized models can't capture (under bigram tag representation):

$$\mathbf{a\ b\ c} \Rightarrow \begin{array}{c} \mathbf{A} \text{ --- } \mathbf{B} \text{ --- } \mathbf{C} \\ | \qquad | \qquad | \\ \mathbf{a} \qquad \mathbf{b} \qquad \mathbf{c} \end{array} \quad \text{with } P(\mathbf{A\ B\ C} \mid \mathbf{a\ b\ c}) = 1$$

$$\mathbf{a\ b\ e} \Rightarrow \begin{array}{c} \mathbf{A} \text{ --- } \mathbf{D} \text{ --- } \mathbf{E} \\ | \qquad | \qquad | \\ \mathbf{a} \qquad \mathbf{b} \qquad \mathbf{e} \end{array} \quad \text{with } P(\mathbf{A\ D\ E} \mid \mathbf{a\ b\ e}) = 1$$

- Impossible to find parameters that satisfy

$$P(A \mid a) \times P(B \mid b, A) \times P(C \mid c, B) = 1$$

$$P(A \mid a) \times P(D \mid b, A) \times P(E \mid e, D) = 1$$

# CONDITIONAL MARKOV RANDOM FIELDS

[Lafferty et al. 2001],[Johnson et. al 1999]

- “Global” normalization

$$\log P(t_{[1:n]} \mid w_{[1:n]}) = \underbrace{\sum_i \sum_s \mathbf{W}_s \phi_s(h_i, t_i)}_{\text{Linear Score}} - \underbrace{\log Z(\mathbf{W}, w_{[1:n]})}_{\text{Global Normalization}}$$

where

$$Z(\mathbf{W}, w_{[1:n]}) = \sum_{t_{[1:n]} \in \mathcal{T}^n} e^{\sum_i \sum_s \mathbf{W}_s \phi_s(h_i, t_i)}$$

- Computational issues
  - Viterbi algorithm finds most likely tag sequence
  - Gradients w.r.t parameters calculated using Baum-Welch!

## A PERCEPTRON ALGORITHM

- Score for a  $(w_{[1:n]}, t_{[1:n]})$  pair is

$$\begin{aligned} F(w_{[1:n]}, t_{[1:n]}) &= \sum_i \sum_s \mathbf{W}_s \phi_s(h_i, t_i) \\ &= \sum_s \mathbf{W}_s \Phi_s(t_{[1:n]}, w_{[1:n]}) \end{aligned}$$

- Note: no normalization terms
- Note:  $F(w_{[1:n]}, t_{[1:n]})$  is not a log probability
- Viterbi algorithm for

$$\arg \max_{t_{[1:n]} \in \mathcal{T}^n} F(w_{[1:n]}, t_{[1:n]})$$

# TRAINING THE PARAMETERS

**Inputs:** Training set  $(w_{[1:n_i]}^i, t_{[1:n_i]}^i)$  for  $i = 1 \dots n$ .

**Initialization:**  $\mathbf{W} = 0$

**Algorithm:** For  $t = 1 \dots T, i = 1 \dots n$

$$z_{[1:n_i]} = \arg \max_{u_{[1:n_i]} \in \mathcal{T}^{n_i}} \sum_s \mathbf{W}_s \Phi_s(w_{[1:n_i]}^i, u_{[1:n_i]})$$

$z_{[1:n_i]}$  is output on  $i$ 'th sentence with current parameters

If  $z_{[1:n_i]} \neq t_{[1:n_i]}^i$  then

$$\mathbf{W}_s = \mathbf{W}_s + \underbrace{\Phi_s(w_{[1:n_i]}^i, t_{[1:n_i]}^i)}_{\text{Correct tags' feature value}} - \underbrace{\Phi_s(w_{[1:n_i]}^i, z_{[1:n_i]})}_{\text{Incorrect tags' feature value}}$$

**Output:** Parameter vector  $\mathbf{W}$ .

## AN EXAMPLE

Say the correct tags for  $i$ 'th sentence are

the/**DT** man/**NN** bit/**VBD** the/**DT** dog/**NN**

Under current parameters, output is

the/**DT** man/**NN** bit/**NN** the/**DT** dog/**NN**

---

Assume also that features track: (1) all bigrams; (2) word/tag pairs

Parameters incremented:

$\langle \text{NN}, \text{VBD} \rangle$ ,  $\langle \text{VBD}, \text{DT} \rangle$ ,  $\langle \text{VBD} \rightarrow \text{bit} \rangle$

Parameters decremented:

$\langle \text{NN}, \text{NN} \rangle$ ,  $\langle \text{NN}, \text{DT} \rangle$ ,  $\langle \text{NN} \rightarrow \text{bit} \rangle$

## A REFINEMENT: AVERAGED PARAMETERS

$\mathbf{W}_{t,i}$  = weight vector at  $i$ 'th example on  $t$ 'th pass

$$\mathbf{W}_{AVG} = \sum_{t=1 \dots T, i=1 \dots n} \frac{\mathbf{W}_{t,i}}{Tn}$$

---

**Algorithm:** For  $t = 1 \dots T, i = 1 \dots n$

- $z_{[1:n_i]} = \arg \max_{u_{[1:n_i]} \in \mathcal{T}^{n_i}} \sum_s \mathbf{W}_s \Phi_s(w_{[1:n_i]}^i, u_{[1:n_i]})$
- If  $z_{[1:n_i]} \neq t_{[1:n_i]}^i$  then  
 $\mathbf{W}_s = \mathbf{W}_s + \Phi_s(w_{[1:n_i]}^i, t_{[1:n_i]}^i) - \Phi_s(w_{[1:n_i]}^i, z_{[1:n_i]})$
- $\mathbf{W}_{AVG} = \mathbf{W}_{AVG} + \mathbf{W} / Tn$

## EXPERIMENTS

- Wall Street Journal part-of-speech tagging data

≈ One million words of training data

300,000 words of development data

300,000 words of test data

- [Ramshaw and Marcus 95] NP chunking data

≈ 150,000 words of training data

50,000 words of development data

50,000 words of test data

## RESULTS ON THE POS TAGGING TASK

### Development Data:

Method	Error rate/%	Numits
Perc, avg, cc=0	<b>2.93</b>	10
Perc, noavg, cc=0	3.68	20
Perc, avg, cc=5	3.03	6
Perc, noavg, cc=5	4.04	17
ME, cc=0	3.4	100
ME, cc=5	<b>3.28</b>	200

*cc* = 0 All features included

*cc* = 5 Only features occurring 5 times or more in training

**avg** Averaged perceptron, **noavg** Regular perceptron

**Test Data:** **Perceptron = 2.89%**, **Max-ent = 3.28%**  
(11.9% relative error reduction)

## RESULTS ON THE CHUNKING TASK

### Development Data:

Method	F-Measure	Numits
Perceptron, avg, cc=0	<b>93.53</b>	13
Perceptron, noavg, cc=0	93.04	35
Perceptron, avg, cc=5	93.33	9
Perceptron, noavg, cc=5	91.88	39
Max-ent, cc=0	92.34	900
Max-ent, cc=5	<b>92.65</b>	200

*cc* = 0 All features included

*cc* = 5 Only features occurring 5 times or more in training

**avg** Averaged perceptron, **noavg** Regular perceptron

**Test Data:** **Perceptron = 93.63%**, **Max-ent = 93.29%**  
(5.1% relative error reduction)

## A FINAL TWIST: TRAINING USING BEAM SEARCH

- Even if representation  $\Phi$  contains features with potentially large “span”, can **approximate** search for best hypothesis by keeping top  $N$  highest scoring prefixes at each stage
- Can also use Viterbi algorithm, beam-search for parsing
- **Main point:**  
if we can (approximately) find the highest scoring hypothesis under a representation  $\Phi$ , we can train using the perceptron algorithm
- **Second point:**  
Convergence rate, generalization ability of the algorithm depends on **separability**, not on the size of **GEN**

## SUMMARY

- Margin–based algorithm for NLP: the Voted perceptron
- Three alternatives within the framework:
  - “Arbitrary” features
  - Kernels
  - Training weighted grammars
- Methods are applicable to many other tasks:
  - Language generation, Speech recognition
- Flexibility in representation:
  - Freedom to define global features
  - The tree kernel generalizes to other recursive structures
- Future questions:
  - How to do away with first pass altogether?

# MARKOV RANDOM FIELDS

Parameters  $\mathbf{W}$  define a conditional distribution over candidates:

$$P(y_i \mid x_i, \mathbf{W}) = \frac{e^{\Phi(y_i) \cdot \mathbf{W}}}{\sum_{y \in \mathbf{GEN}(x_i)} e^{\Phi(y) \cdot \mathbf{W}}}$$

Gaussian prior:  $\log P(\mathbf{W}) \sim -C \|\mathbf{W}\|^2 / 2$

MAP parameter estimates maximise

$$\begin{aligned} & \sum_i \log \frac{e^{\Phi(y_i) \cdot \mathbf{W}}}{\sum_{y \in \mathbf{GEN}(x_i)} e^{\Phi(y) \cdot \mathbf{W}}} - C \frac{\|\mathbf{W}\|^2}{2} \\ &= - \left( \sum_i \log \left( 1 + \sum_{y \in \overline{\mathbf{GEN}}(x_i)} e^{\Phi(y) \cdot \mathbf{W} - \Phi(y_i) \cdot \mathbf{W}} \right) + C \frac{\|\mathbf{W}\|^2}{2} \right) \end{aligned}$$

# References

- [[Abney 1997](#)] Abney
- [[ABR64](#)] Aizerman, M., Braverman, E., and Rozonoer, L. (1964). Theoretical Foundations of the Potential Function Method in Pattern Recognition Learning. *Automation and Remote Control*, 25:821–837.
- [[Bod 98](#)] Bod, R. (1998). *Beyond Grammar: An Experience-Based Theory of Language*. CSLI Publications/Cambridge University Press.
- [[CHA97](#)] Charniak, E. (1997). Statistical techniques for natural language parsing. In *AI Magazine*, Vol. 18, No. 4.
- [[COL00](#)] Collins, M. (2000). Discriminative Reranking for Natural Language Parsing. *Proceedings of the Seventeenth International Conference on Machine Learning*. San Francisco: Morgan Kaufmann.
- [[CD01](#)] Collins, M. and Duffy, N. (2001). Parsing with a Single Neuron: Convolution Kernels for Natural Language Problems. Technical report UCSC-CRL-01-01, University of California at Santa Cruz.
- [[CV95](#)] Cortes, C. and Vapnik, V. (1995). Support–Vector Networks. *Machine Learning*, 20(3):273–297.
- [[Work with Rob Schapire, Yoram Singer at COLT 2000](#)]
- [[Friedman, Hastie, and Tibshirani 1998](#)] Friedman.
- [[Lafferty 1999](#)] lafferty.
- [[Della Pietra, Della Pietra, and Lafferty 1997](#)] Lafferty
- [[Freund and Schapire 99](#)] Freund, Y. and Schapire, R. (1999). Large Margin Classification using the Perceptron Algorithm. In *Machine Learning*, 37(3):277–296.
- [[Lafferty et al. 2001](#)] John Lafferty, Andrew McCallum, and Fernando Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In Proceedings of ICML-01, pages 282-289, 2001.
- [[FISS98](#)] Freund, Y., Iyer, R., Schapire, R.E., & Singer, Y. (1998). An efficient boosting algorithm for combining preferences. In *Machine Learning: Proceedings of the Fifteenth International Conference*. San Francisco: Morgan Kaufmann.
- [[HAU99](#)] Haussler, D. (1999). *Convolution Kernels on Discrete Structures*. Technical report, University of Santa Cruz.

[[Helmbold and Warmuth 95](#)] Helmbold, D., and Warmuth, M. On Weak Learning.

[[LCSW01](#)] Lodhi, H., Cristianini, N., Shawe-Taylor, J., and Watkins, C. (2001). Text Classification using String Kernels. To appear in *Advances in Neural Information Processing Systems 13*, MIT Press.

[[Johnson et. al 1999](#)] Johnson, M., Geman, S., Canon, S., Chi, S., & Riezler, S. (1999). Estimators for stochastic ‘unification-based’ grammars. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*. San Francisco: Morgan Kaufmann.

[[Ramshaw and Marcus 95](#)]

[[Johnson, Geman, Canon, Chi, and Riezler 1999](#)] Johnson, M., Geman, S., Canon, S., Chi, S., & Riezler, S. (1999). Estimators for stochastic ‘unification-based’ grammars. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*. San Francisco: Morgan Kaufmann.

[[MSM93](#)] Marcus, M., Santorini, B., & Marcinkiewicz, M. (1993). Building a large annotated corpus of english: The Penn treebank. *Computational Linguistics*, 19, 313-330.

[[Prince and Smolensky](#)] Prince and Smolensky. Optimality theory.

[[Ratnaparkhi 96](#)]

[[Ratnaparkhi, Roukos and Ward 94](#)] Adwait Ratnaparkhi, Salim Roukos, and R. Todd Ward. A Maximum Entropy Model for Parsing. In Proceedings of the International Conference on Spoken Language Processing, pages 803-806. September, 1994. Yokohama, Japan.

[[SSM99](#)] Scholkopf, B., Smola, A., and Muller, K.-R. (1999). Kernel principal component analysis. In B. Scholkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods – SV Learning*, pages 327-352. MIT Press, Cambridge, MA.

[[Walker et al. 2001](#)] Walker, M., Rambow, O., and Rogati, M. (2001). SPoT: a trainable sentence planner. In *Proceedings of the 2nd Meeting of the North American Chapter of the Association for Computational Linguistics (NAACL 2001)*.

[[WAT00](#)] Watkins, C. (2000). Dynamic alignment kernels. In A.J. Smola, P.L. Bartlett, B. Scholkopf, and D. Schuurmans, editors, *Advances in Large Margin Classifiers*, pages 39-50, MIT Press.