

Feature-based Discriminative Models

Yoav Goldberg

Bar Ilan University

Sentence-boundary-detection revisited

Georg "Mr. George" Svendsen (19 March 1894 – 1966) was a Norwegian journalist and crime novelist.

He was born in Eidanger, and started his journalistic career in Bratsberg-Demokraten before moving on to Demokraten where he was a subeditor. In 1921 he was hired in Fremtiden and replaced in Demokraten by Evald O. Solbakken. In 1931 he was hired in Arbeiderbladet. Under the pen name "Mr. George" he became known for his humorous articles in the newspaper. At his death he was also called "the last of the three great criminal and police reporters in Oslo", together with Fridtjof Knutsen and Axel Kielland. He was also known for practising journalism as a trade in itself, and not as a part of a political career. He retired in 1964, and died in 1966.

He released the criminal novels *Mannen med ljåen* (1942), *Ridderne av øksen* (1945) and *Den hvite streken* (1946), and translated the book *S.S. Murder by Quentin Patrick as Mord ombord* in 1945. He released several historical books: *Rørleggernes Fagforenings historie gjennom 50 år* (1934), *Telefonmontørenes Forening Oslo, gjennom 50 år* (1939), *Norsk nærings- og nydelsmiddelarbeiderforbund: 25-års beretning* (1948), *De tause vitner: av rettskjemiker Ch. Bruffs memoarer* (1949, with Fridtjof Knudsen) and *Elektriske montørers fagforening gjennom 50 år* (1949).

Sentence-boundary-detection revisited

- ▶ $P(\text{boundary} \mid \text{subeditor} . \text{In})$
- ▶ $P(\text{boundary} \mid \text{O} . \text{Solbakken})$
- ▶ $P(\text{boundary} \mid \text{Solbakken} . \text{In})$
- ▶ $P(\text{boundary} \mid \text{Arbeiderbladet} . \text{Under})$
- ▶ $P(\text{boundary} \mid \text{Mr} . \text{George})$
- ▶ $P(\text{boundary} \mid 1945 . \text{He})$

Sentence-boundary-detection revisited

- ▶ $P(\text{boundary} \mid L=\text{subeditor} . R=\text{In})$
- ▶ $P(\text{boundary} \mid L=\text{O} . R=\text{Solbakken})$
- ▶ $P(\text{boundary} \mid L=\text{Solbakken} . R=\text{In})$
- ▶ $P(\text{boundary} \mid L=\text{Arbeiderbladet} . R=\text{Under})$
- ▶ $P(\text{boundary} \mid L=\text{Mr} . R=\text{George})$
- ▶ $P(\text{boundary} \mid L=1945 . R=\text{He})$

Sentence-boundary-detection revisited

- ▶ $P(\text{boundary} \mid L=\text{subeditor} . R=\text{In})$
- ▶ $P(\text{boundary} \mid L=\text{O} . R=\text{Solbakken})$
- ▶ $P(\text{boundary} \mid L=\text{Solbakken} . R=\text{In})$
- ▶ $P(\text{boundary} \mid L=\text{Arbeiderbladet} . R=\text{Under})$
- ▶ $P(\text{boundary} \mid L=\text{Mr} . R=\text{George})$
- ▶ $P(\text{boundary} \mid L=1945 . R=\text{He})$

Useful indicators

- ▶ Identity of L
- ▶ Identity of R
- ▶ Length of L
- ▶ Is R capitalized
- ▶ ...

Sentence-boundary-detection revisited

- ▶ $P(\text{boundary} \mid L=\text{subeditor} . R=\text{In})$
- ▶ $P(\text{boundary} \mid L=\text{O} . R=\text{Solbakken})$
- ▶ $P(\text{boundary} \mid L=\text{Solbakken} . R=\text{In})$
- ▶ $P(\text{boundary} \mid L=\text{Arbeiderbladet} . R=\text{Under})$
- ▶ $P(\text{boundary} \mid L=\text{Mr} . R=\text{George})$
- ▶ $P(\text{boundary} \mid L=1945 . R=\text{He})$

Useful indicators

- ▶ **Identity of L** Highly correlated.
- ▶ Identity of R Do we need both?
- ▶ **Length of L** How do we integrate information?
- ▶ Is R capitalized
- ▶ ...

Sentence-boundary-detection revisited

$P(\text{boundary} | L=\text{subeditor}, R=\text{In}, \text{len}(L)=9, \text{isCap}(R)=\text{True},$
 $\text{isCap}(L)=\text{False}, \dots)$

- ▶ How do we compute the probability model?
- ▶ Do we really need a probability model here?

Probabilities → Scores

- ▶ Instead of $P(y|x)$, compute $\text{score}(x, y)$
- ▶ Return $\arg \max_y \text{score}(x, y)$
 - ▶ Or, if we just have two classes, return class 1 iff $\text{score}(x) > 0$
- ▶ But how do we compute the score?

Feature Representation

Each indicator will be a feature.

Feature

- ▶ A **feature** is a function $\phi_i(x)$ looking at a particular property of x and returning a value.
 - ▶ Generally, the value can be any number, $\phi_i(x) \in \mathbb{R}$
 - ▶ Often, the value is binary, $\phi_i(x) \in \{0, 1\}$.

Feature Extractor / Feature Vector

- ▶ A **feature extractor** is a collection of features
$$\phi = \phi_1, \dots, \phi_m$$
- ▶ A **feature vector** is the result of ϕ applied on a particular x
 - ▶ $\phi(x) = \phi_1(x), \dots, \phi_m(x)$
 - ▶ For binary features: $\phi(x) \in \{0, 1\}^m$

Feature Representation

Feature Examples

$$\phi_1(L.R) = \begin{cases} 1 & \text{if } L = \text{'subeditor'}$$

$$\phi_2(L.R) = \begin{cases} 1 & \text{if } L = \text{'O'}$$

$$\phi_{57}(L.R) = \begin{cases} 1 & \text{if } R = \text{'George'}$$

$$\phi_{1039}(L.R) = \begin{cases} 1 & \text{if } \text{len}(L) = 1$$

$$\phi_{1040}(L.R) = \begin{cases} 1 & \text{if } \text{len}(L) = 2$$

$$\phi_{1045}(L.R) = \begin{cases} 1 & \text{if } \text{len}(L) > 4$$

$$\phi_{1092}(L.R) = \begin{cases} 1 & \text{if } \text{isCap}(R)$$

$$\phi_{1093}(L.R) = \begin{cases} 1 & \text{if } \text{isNumber}(L)$$

$$\phi_{2045}(L.R) = \begin{cases} 1 & \text{if } L = \text{'US' and } \text{isCap}(R)$$

$$\phi_{3066}(L.R) = \begin{cases} 1 & \text{if } \text{len}(L) = 1 \text{ and } \text{isCap}(R)$$

$$\phi_{5032}(L.R) = \begin{cases} 1 & \text{if } L = \text{'Ca' and } R = \text{'The'}$$

$$\phi_{5033}(L.R) = \begin{cases} 1 & \text{if } L = \text{'Ca' and } R = \text{'snow'}$$

Feature Representation

Feature Vector

$\phi(\text{L=subeditor . R=In}) = 1, 0, 0, 0, 0, \dots, 1, 0, 0, 0, 1 \dots$

$\phi(\text{L=Mr . R=George}) = 0, 0, 0, 1, 0, \dots, 0, 0, 1, 0, 1 \dots$

- ▶ Each example is mapped to a very long vector of 0 and 1.
- ▶ Most values are 0.

Sparse Representation

We don't need to write the zeros:

$\phi(\text{L=subeditor . R=In}) = 1:1 \ 1045:1 \ 1092:1 \ \dots$

Machine Learning

We had a set of example:

- ▶ subeditor . In \rightarrow boundary
- ▶ O . Solbakken \rightarrow no-boundary
- ▶ Mr . George \rightarrow no-boundary
- ▶ ...

We can map them to a list of vectors

- ▶ $\phi(L=\text{subeditor} . R=\text{In}) = 1:1 \ 1045:1 \ 1092:1 \dots$
- ▶ $\phi(L=O . R=\text{Solbakken}) = 2:1 \ 1039:1 \ 1092:1 \dots$
- ▶ $\phi(L=\text{Mr} . R=\text{George}) = 57:1 \ 1040:1 \ 1092:1 \dots$
- ▶ ...

And a list of answers:

- ▶ $+1, -1, -1, \dots$

Feature Representation

Why did we do this?

- ▶ We mapped each example x to a vector of numbers $\phi(x)$.
- ▶ We mapped each answer y to a number.
- ▶ Instead of $P(y|x)$ we now have $P(y|\phi(x))$ $y \in \{-1, +1\}$.
 - ▶ x was a sequence of words.
 - ▶ $\phi(x)$ is a vector of numbers.
- ▶ The field of **Machine Learning** is very good at learning to classify vectors.

Machine Learning

Dataset

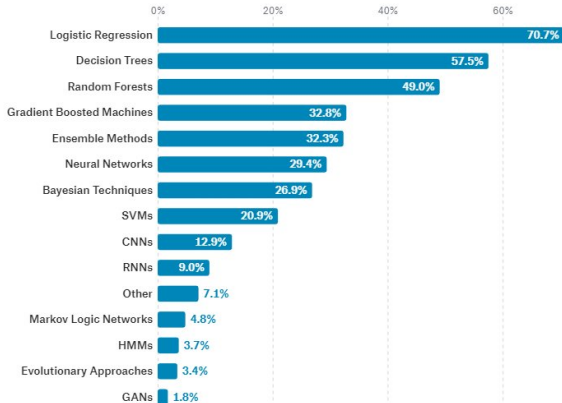
+1 1:1 1045:1 1092:1 ...
-1 2:1 1039:1 1092:1 ...
-1 57:1 1040:1 ...
...

- ▶ We can feed this dataset to a machine-learning algorithm.
- ▶ Many machine learning algorithms exist:
 - ▶ SVM, boosting, decision trees, knn, logistic regression, random forests, perceptron, neural networks, ...
- ▶ The important distinctions:
 - ▶ output type:
 - ▶ Binary, Multiclass, Numeric, Structured
 - ▶ scoring type:
 - ▶ Linear vs. Non Linear

What data science methods are used at work?

Logistic regression is the most commonly reported data science method used at work for all industries *except* **Military and Security** where Neural Networks are used slightly more frequently.

Company Size ▾ Financial ▾ Job Title ▾



789 responses

 [View code in Kaggle Kernels](#)

<https://www.kaggle.com/surveys/2017>

Types of learning problems

Binary

- ▶ Answer is binary: $y \in \{-1, +1\}$
 - ▶ boundary/no-boundary, verb-attach/noun-attach, yes/no

Multiclass

- ▶ Answer is one of k options: $y \in \{1, \dots, k\}$
 - ▶ choose the part-of-speech of a word.
 - ▶ identify the language of a document.

Numeric / Regression

- ▶ Answer is a number: $y \in \mathbb{R}$
 - ▶ Predict the height of a person.

Structured

- ▶ Answer is a complex object:
 - ▶ Predict a sequence of tags for a sentence.

Generic NLP Solution

- ▶ Find an annotated corpus
- ▶ Split it into train and test parts
- ▶ Convert it to a vector representation
 - ▶ Decide on output type
 - ▶ Decide on features
 - ▶ Convert each training example to feature vector
- ▶ Train a machine-learning model on training set
- ▶ Apply machine-learning model to test set
- ▶ Measure accuracy

Linear Models

Machine Learning

- ▶ Many learning algorithms exist.
- ▶ Some of them are based on a linear model:
 - ▶ SVM, boosting, logistic-regression, perceptron
- ▶ Others are non-linear:
 - ▶ Kernel-SVM, decision-trees, knn, random-forests, neural-networks.
- ▶ We will work mostly with linear classifiers in this course.
- ▶ Neural networks are useful and popular. But there's a separate course for them.

Linear Models

- ▶ Some features are indicators for “boundary” and others are good indicators for “no boundary”.
- ▶ We will assign a score (weight) to each feature.
 - ▶ Features in favor of boundary will receive a positive score.
 - ▶ Features in favor of no boundary will receive a negative score.
- ▶ **Use the sum of the scores to classify.**

Linear Models

- ▶ Some features are indicators for “boundary” and others are good indicators for “no boundary”.
- ▶ We will assign a score (weight) to each feature.
 - ▶ Features in favor of boundary will receive a positive score.
 - ▶ Features in favor of no boundary will receive a negative score.
- ▶ **Use the sum of the scores to classify.**

Binary Linear Classifier

$$score(x) = \sum_{i \in \{1, \dots, m\}} w_i \times \phi_i(x) = \mathbf{w} \cdot \phi(x)$$

$$\hat{y} = predict(x) = \begin{cases} 1 & score(x) \geq 0 \\ -1 & score(x) < 0 \end{cases}$$

Setting the Weights

Feature Examples

$$\phi_1(L.R) = \begin{cases} 1 & \text{if } L = \text{'subeditor'}$$

$$\phi_2(L.R) = \begin{cases} 1 & \text{if } L = \text{'O'}$$

$$\phi_{57}(L.R) = \begin{cases} 1 & \text{if } R = \text{'George'}$$

$$\phi_{1039}(L.R) = \begin{cases} 1 & \text{if } \text{len}(L) = 1$$

$$\phi_{1040}(L.R) = \begin{cases} 1 & \text{if } \text{len}(L) = 2$$

$$\phi_{1045}(L.R) = \begin{cases} 1 & \text{if } \text{len}(L) > 4$$

$$\phi_{1092}(L.R) = \begin{cases} 1 & \text{if } \text{isCap}(R)$$

$$\phi_{1093}(L.R) = \begin{cases} 1 & \text{if } \text{isNumber}(L)$$

$$\phi_{2045}(L.R) = \begin{cases} 1 & \text{if } L = \text{'US' and } \text{isCap}(R)$$

$$\phi_{3066}(L.R) = \begin{cases} 1 & \text{if } \text{len}(L) = 1 \text{ and } \text{isCap}(R)$$

$$\phi_{5032}(L.R) = \begin{cases} 1 & \text{if } L = \text{'Ca' and } R = \text{'The'}$$

$$\phi_{5033}(L.R) = \begin{cases} 1 & \text{if } L = \text{'Ca' and } R = \text{'snow'}$$

Setting the Weights

Binary Linear Classifier

$$score(x) = \sum_{i \in \{1, \dots, m\}} w_i \times \phi_i(x) = \mathbf{w} \cdot \phi(x)$$

$$\hat{y} = predict(x) = \begin{cases} 1 & score(x) \geq 0 \\ -1 & score(x) < 0 \end{cases}$$

Machine Learning

Provide algorithms for **learning** the values of w from examples.

- ▶ A supervised binary learning problem:
 - ▶ **Input:** a set of (x, y) examples, features extractor ϕ .
 - ▶ **Output:** weights \mathbf{w} such that $\mathbf{w} \cdot \phi(x)$ classifies many of the y 's correctly.

Setting weights: the perceptron algorithm

- 1: $\mathbf{w} \leftarrow \mathbf{0}$
- 2: **for** x, y in training examples **do**
- 3: $\hat{y} \leftarrow \text{sign}(\mathbf{w} \cdot \phi(x))$
- 4: **if** $\hat{y} < 0$ and $y \geq 0$ **then**
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \phi(x)$
- 6: **if** $\hat{y} \geq 0$ and $y < 0$ **then**
- 7: $\mathbf{w} \leftarrow \mathbf{w} - \phi(x)$
- 8: **return** \mathbf{w}

Setting the Weights

Multiclass Linear Classifier

$$\text{score}(x, y) = \sum_{i \in 1, \dots, m} w_i^{(y)} \times \phi_i(x, y) = \mathbf{w} \cdot \phi(x, y)$$

$$\hat{y} = \text{predict}(x) = \arg \max_y \text{score}(x, y)$$

Setting the Weights

Multiclass Linear Classifier

$$\text{score}(x, y) = \sum_{i \in 1, \dots, m} w_i^{(y)} \times \phi_i(x, y) = \mathbf{w} \cdot \phi(x, y)$$

$$\hat{y} = \text{predict}(x) = \arg \max_y \text{score}(x, y)$$

Alternative view

$$\text{score}(x, y) = \sum_{i \in 1, \dots, m} w_i \times \phi_i(x) = \mathbf{w}^{(y)} \cdot \phi(x)$$

$$\hat{y} = \text{predict}(x) = \arg \max_y \text{score}(x, y)$$

Here, each class receives its own weight vector.

Setting weights: the multiclass perceptron algorithm

```
1:  $\mathbf{w} \leftarrow 0$ 
2: for  $x, y$  in training examples do
3:    $\hat{y} \leftarrow \arg \max_{y'} (\mathbf{w} \cdot \phi(x, y'))$ 
4:   if  $\hat{y} \neq y$  then
5:      $\mathbf{w} \leftarrow \mathbf{w} + \phi(x, y)$ 
6:      $\mathbf{w} \leftarrow \mathbf{w} - \phi(x, \hat{y})$ 
7: return  $\mathbf{w}$ 
```

Setting weights: the multiclass perceptron algorithm

- 1: $\mathbf{w} \leftarrow 0$
- 2: **for** x, y in training examples **do**
- 3: $\hat{y} \leftarrow \arg \max_{y'} (\mathbf{w} \cdot \phi(x, y'))$
- 4: **if** $\hat{y} \neq y$ **then**
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \phi(x, y) - \phi(x, \hat{y})$
- 6: **return** \mathbf{w}

Perceptron \rightarrow Averaged Perceptron

- 1: $\mathbf{w} \leftarrow 0$
- 2: **for** x, y in training examples **do**
- 3: $\hat{y} \leftarrow \arg \max_{y'} (\mathbf{w} \cdot \phi(x, y'))$
- 4: **if** $\hat{y} \neq y$ **then**
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \phi(x, y) - \phi(x, \hat{y})$
- 6: **return** \mathbf{w}

Averaged Perceptron

- ▶ The perceptron algorithm is not a very good learner
- ▶ But can be easily improved

Perceptron \rightarrow Averaged Perceptron

- 1: $\mathbf{w} \leftarrow 0$
- 2: **for** x, y in training examples **do**
- 3: $\hat{y} \leftarrow \arg \max_{y'} (\mathbf{w} \cdot \phi(x, y'))$
- 4: **if** $\hat{y} \neq y$ **then**
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \phi(x, y) - \phi(x, \hat{y})$
- 6: **return** \mathbf{w}

Averaged Perceptron

- ▶ The perceptron algorithm is not a very good learner
- ▶ But can be easily improved
- ▶ Instead of returning \mathbf{w} , return $avg(\mathbf{w})$
 - ▶ $avg(\mathbf{w})$ is the average of all versions of \mathbf{w} seen in training

Perceptron \rightarrow Averaged Perceptron

- 1: $\mathbf{w} \leftarrow 0$
- 2: **for** x, y in training examples **do**
- 3: $\hat{y} \leftarrow \arg \max_{y'} (\mathbf{w} \cdot \phi(x, y'))$
- 4: **if** $\hat{y} \neq y$ **then**
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \phi(x, y) - \phi(x, \hat{y})$
- 6: **return** \mathbf{w}

Averaged Perceptron

- ▶ The perceptron algorithm is not a very good learner
- ▶ But can be easily improved
- ▶ Instead of returning \mathbf{w} , return $avg(\mathbf{w})$
 - ▶ $avg(\mathbf{w})$ is the average of all versions of \mathbf{w} seen in training
- ▶ Require a bit more book-keeping for calculating the average at the end
- ▶ **The result is a very competitive algorithm**

PP-attachment revisited

We calculated

$P(V|v = \text{saw}, n1 = \text{mouse}, p = \text{with}, n2 = \text{telescope})$

Problems

- ▶ Was not trivial to come up with a formula.
- ▶ Hard to add more sources of information.

New solution

- ▶ Encode as a binary or multiclass classification.
- ▶ Decide on features.
- ▶ Apply learning algorithm.

PP-attachment

Multiclass classification problem

- ▶ Previously, we had a binary classification problem:
 - ⇒ $x = (v, n1, p, n2)$
 - ⇒ $y \in \{V, N\}$
- ▶ Let's do a multiclass problem:
 - ⇒ $y \in \{V, N, \text{Other}\}$

PP-attachment

Types of features

- ▶ Single items
 - ▶ Identity of v
 - ▶ Identity of p
 - ▶ Identity of n_1
 - ▶ Identity of n_2
- ▶ Pairs
 - ▶ Identity of (v,p)
 - ▶ Identity of (n_1,p)
 - ▶ Identity of (p,n_1)
- ▶ Triplets
 - ▶ Identity of (v,n_1,p)
 - ▶ Identity of (v,p,n_2)
 - ▶ Identity of (n_1,p,n_2)
- ▶ Quadruple
 - ▶ Identity of (v,n_1,p,n_2)

PP-attachment

Types of features

- ▶ Single items
 - ▶ **Identity of v**
 - ▶ Identity of p
 - ▶ Identity of n1
 - ▶ Identity of n2
- ▶ Pairs
 - ▶ Identity of (v,p)
 - ▶ Identity of (n1,p)
 - ▶ Identity of (p,n1)
- ▶ Triplets
 - ▶ Identity of (v,n1,p)
 - ▶ Identity of (v,p,n2)
 - ▶ Identity of (n1,p,n2)
- ▶ Quadruple
 - ▶ Identity of (v,n1,p,n2)

$$\begin{aligned}\phi_{1039} &= \begin{cases} 1 & \text{if } v='ate' \text{ and } y='N' \\ 0 & \textit{otherwise} \end{cases} \\ \phi_{1040} &= \begin{cases} 1 & \text{if } v='ate' \text{ and } y='V' \\ 0 & \textit{otherwise} \end{cases} \\ \phi_{1041} &= \begin{cases} 1 & \text{if } v='ate' \text{ and } y='O' \\ 0 & \textit{otherwise} \end{cases} \\ \phi_{1042} &= \begin{cases} 1 & \text{if } v='saw' \text{ and } y='N' \\ 0 & \textit{otherwise} \end{cases} \\ \dots & \end{aligned}$$

PP-attachment

Types of features

- ▶ Single items
 - ▶ Identity of v
 - ▶ Identity of p
 - ▶ Identity of $n1$
 - ▶ Identity of $n2$
- ▶ Pairs
 - ▶ **Identity of (v,p)**
 - ▶ Identity of $(n1,p)$
 - ▶ Identity of $(p,n1)$
- ▶ Triplets
 - ▶ Identity of $(v,n1,p)$
 - ▶ Identity of $(v,p,n2)$
 - ▶ Identity of $(n1,p,n2)$
- ▶ Quadruple
 - ▶ Identity of $(v,n1,p,n2)$

$$\phi_{2349} = \begin{cases} 1 & \text{if } v='ate' \text{ } p='with' \text{ and } y='N' \\ 0 & \textit{otherwise} \end{cases}$$

$$\phi_{2350} = \begin{cases} 1 & \text{if } v='ate' \text{ } p='with' \text{ and } y='V' \\ 0 & \textit{otherwise} \end{cases}$$

$$\phi_{2351} = \begin{cases} 1 & \text{if } v='ate' \text{ } p='with' \text{ and } y='O' \\ 0 & \textit{otherwise} \end{cases}$$

...

PP-attachment

Types of features

▶ Single items

- ▶ Identity of v
- ▶ Identity of p
- ▶ Identity of $n1$
- ▶ Identity of $n2$

▶ Pairs

- ▶ Identity of (v,p)
- ▶ Identity of $(n1,p)$
- ▶ Identity of $(p,n1)$

▶ Triplets

- ▶ Identity of $(v,n1,p)$
- ▶ Identity of $(v,p,n2)$
- ▶ Identity of $(n1,p,n2)$

▶ Quadruple

- ▶ Identity of $(v,n1,p,n2)$

$$\phi_{2349} = \begin{cases} 1 & \text{if } v='ate' \text{ } p='with' \text{ and } y='N' \\ 0 & \textit{otherwise} \end{cases}$$

$$\phi_{2350} = \begin{cases} 1 & \text{if } v='ate' \text{ } p='with' \text{ and } y='V' \\ 0 & \textit{otherwise} \end{cases}$$

$$\phi_{2351} = \begin{cases} 1 & \text{if } v='ate' \text{ } p='with' \text{ and } y='O' \\ 0 & \textit{otherwise} \end{cases}$$

...

PP-attachment

More features?

- ▶ Corpus Level
 - ▶ Did we see the (v,p) pair in a 5-word window in a big corpus?
 - ▶ Did we see the $(n1,p)$ pair in a 5-word window in a big corpus?
 - ▶ Did we see the $(n1,p,n2)$ triplet in a 5-word window in a big corpus?
 - ▶ We can also use counts, or binned counts.
- ▶ Distance
 - ▶ Distance (in words) between v and p
 - ▶ Distance (in words) between $n1$ and p
 - ▶ ...

PP-attachment

- ▶ Compute features for each example.
- ▶ Feed into a machine learning algorithm (for example SVM or perceptron).
- ▶ Get a weight vector.
- ▶ Classify new examples.