

# NLP Tasks

---

Instructor: Lara J. Martin (she/they)

TA: Omkar Kulkarni (he)

<https://laramartin.net/NLP-class/>

# Learning Objectives

---

Define featurization & other ML terminology

Define some “classification” terminology

Distinguish between different text classification tasks

Formalize NLP Tasks at a high-level:

- What are the input/output for a particular task?
- What might the features be?
- What types of applications could the task be used for?



Similar to HW 1

Calculate elementary processes on a dataset

# Speaking of HW 1...

---

Due Feb 20

## Homework 1: Being up to the Task

### Learning Objectives

- Searching for basic information about NLP tasks.
- Exploring a dataset.
- Coming up with appropriate tasks for an application & providing your reasoning behind it.
- Determining appropriate inputs and outputs for tasks.
- Creating a system diagram.

### Description

You work for SuperDuperAI (SDAI), a start-up company that makes AI tools that their customers can use. You are their NLP specialist. One of SDAI's customers recently came to the company with a [database of textbooks](#) that they collected. They want SDAI to make them an app that can quiz people when they select a textbook.

The flow of the app will look like this:

- a. The user types in a keyword that they're interested in, and the app finds relevant textbooks.
- b. They select the textbook and chapter they want to use.
- c. The app displays a question relevant to the chapter.
- d. The user answers the question.
- e. The app gives a numerical score for how well the user answered the question.

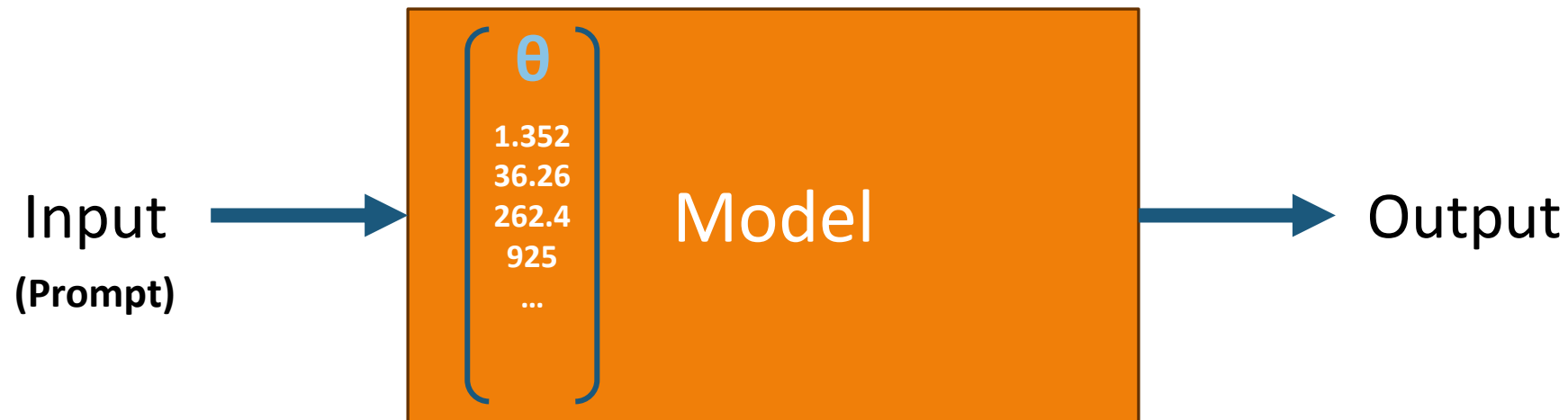
Being the NLP specialist on the team, **you are in charge of figuring out what is needed to create parts a, c, and e.**

# Helpful ML Terminology

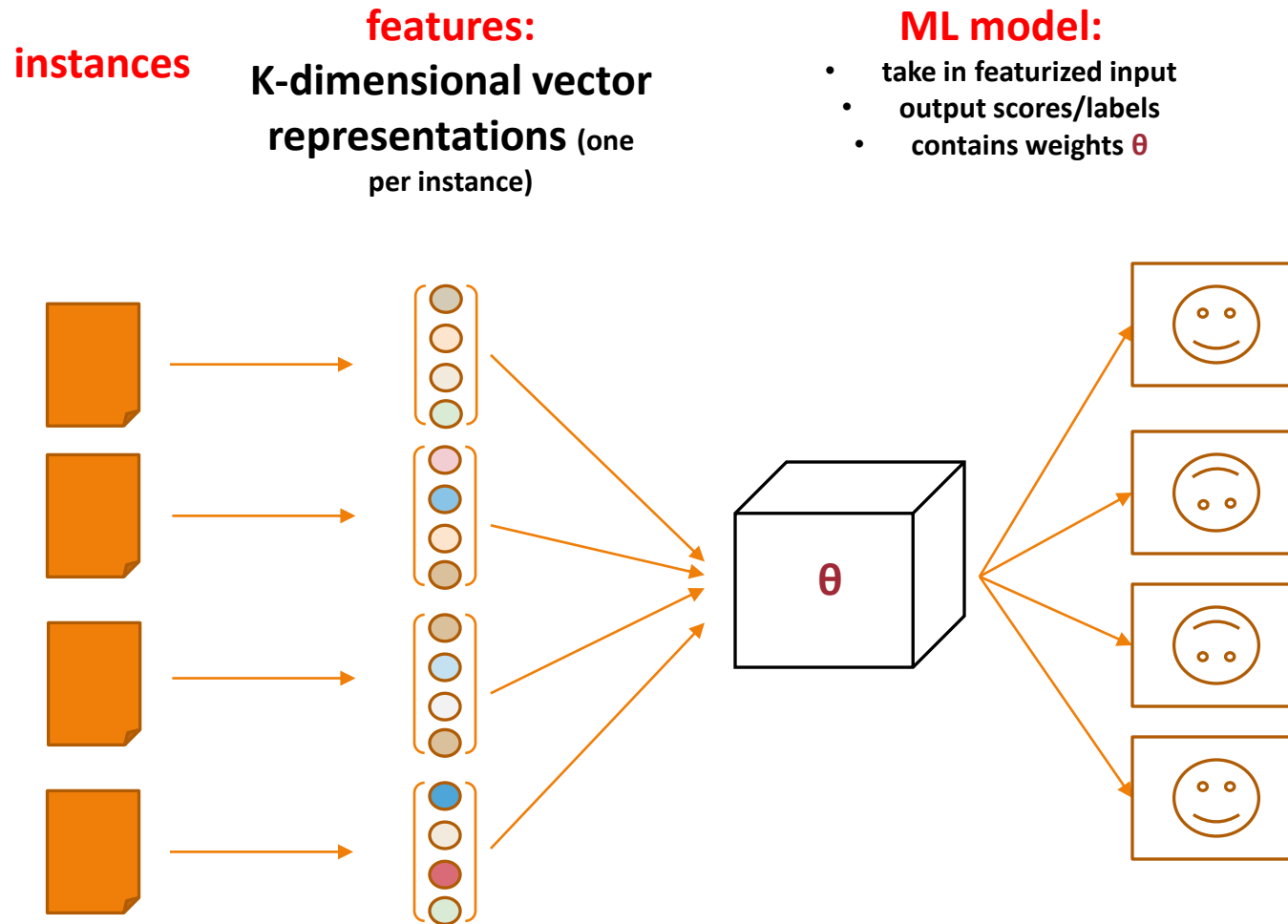
---

**Model:** the (computable) way to go from **features** (input) to labels/scores (output)

**Weights/parameters ( $\theta$ ):** vectors of numbers that control how the model produces labels/scores from inputs. These are learned through **training**.



# ML/NLP Framework



# Helpful ML Terminology

---

**Model:** the (computable) way to go from **features** (input) to labels/scores (output)

**Weights/parameters:** vectors of numbers that control how the model produces labels/scores from inputs. These are learned through **training**.

**Objective function:** an algorithm/calculation, whose variables are the **weights** of the **model**, that we numerically optimize in order to learn appropriate weights based on the labels/scores. The **model's** weights are adjusted.

**Evaluation function:** an algorithm/calculation that scores how “correct” the **model's** predictions are. The **model's** weights are not adjusted.

Note: The evaluation and objective functions are often different!

# (More) Helpful ML Terminology

---

## Training / Learning:

- the process of adjusting the model's weights to learn to make good predictions.

## Inference / Prediction / Decoding / Classification:

- the process of using a model's existing weights to make (hopefully!) good predictions

# ML/NLP Framework for Learning

**instances**

**features:**  
**K-dimensional vector  
representations** (one  
per instance)

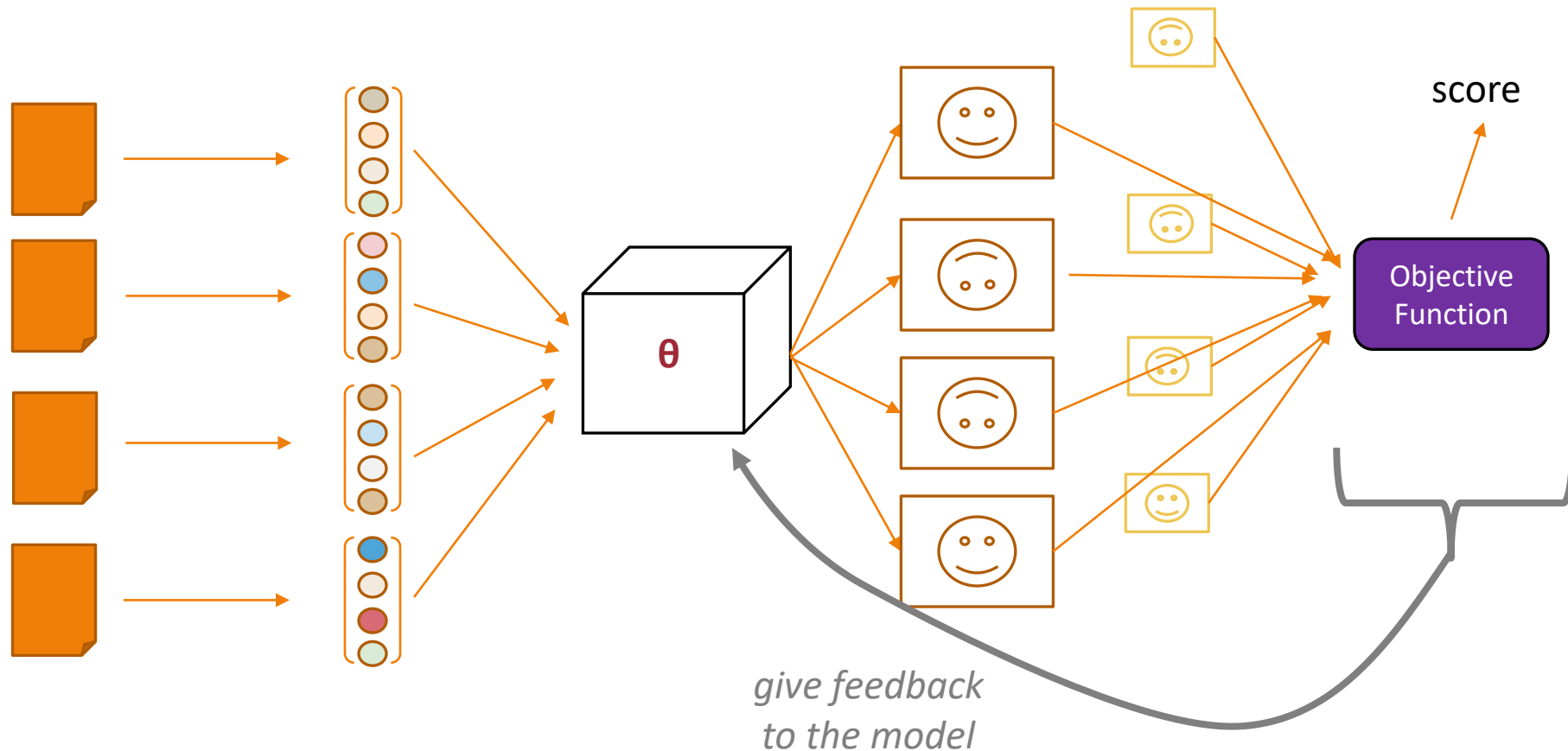
**ML model:**

- take in featurized input
- output scores/labels
- contains weights  $\theta$

**output**

**“Gold”  
(correct)  
labels**

**Objective  
Function/  
Learning**





# ML/NLP Framework for Prediction

**instances**

**features:**  
**K-dimensional vector  
representations** (one  
per instance)

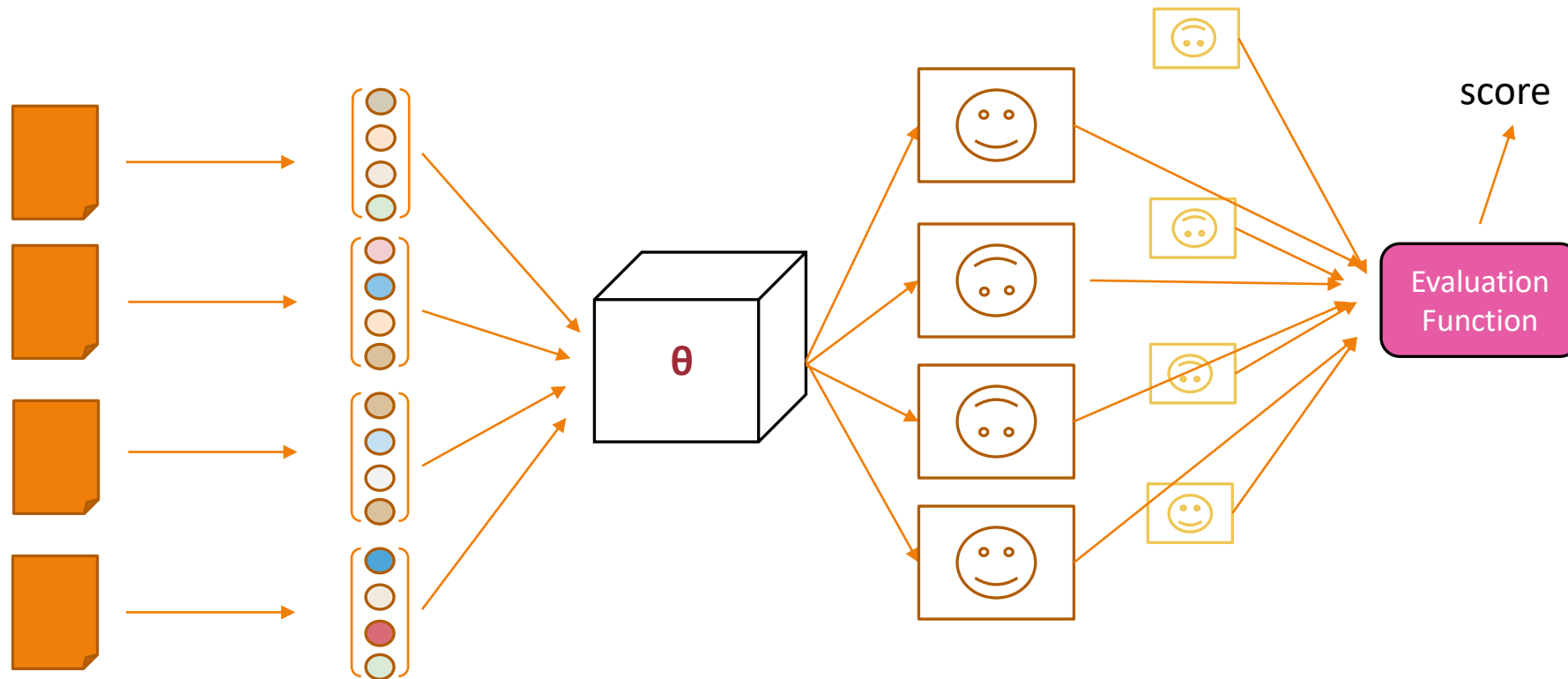
**ML model:**

- take in featurized input
- output scores/labels
- contains weights  $\theta$

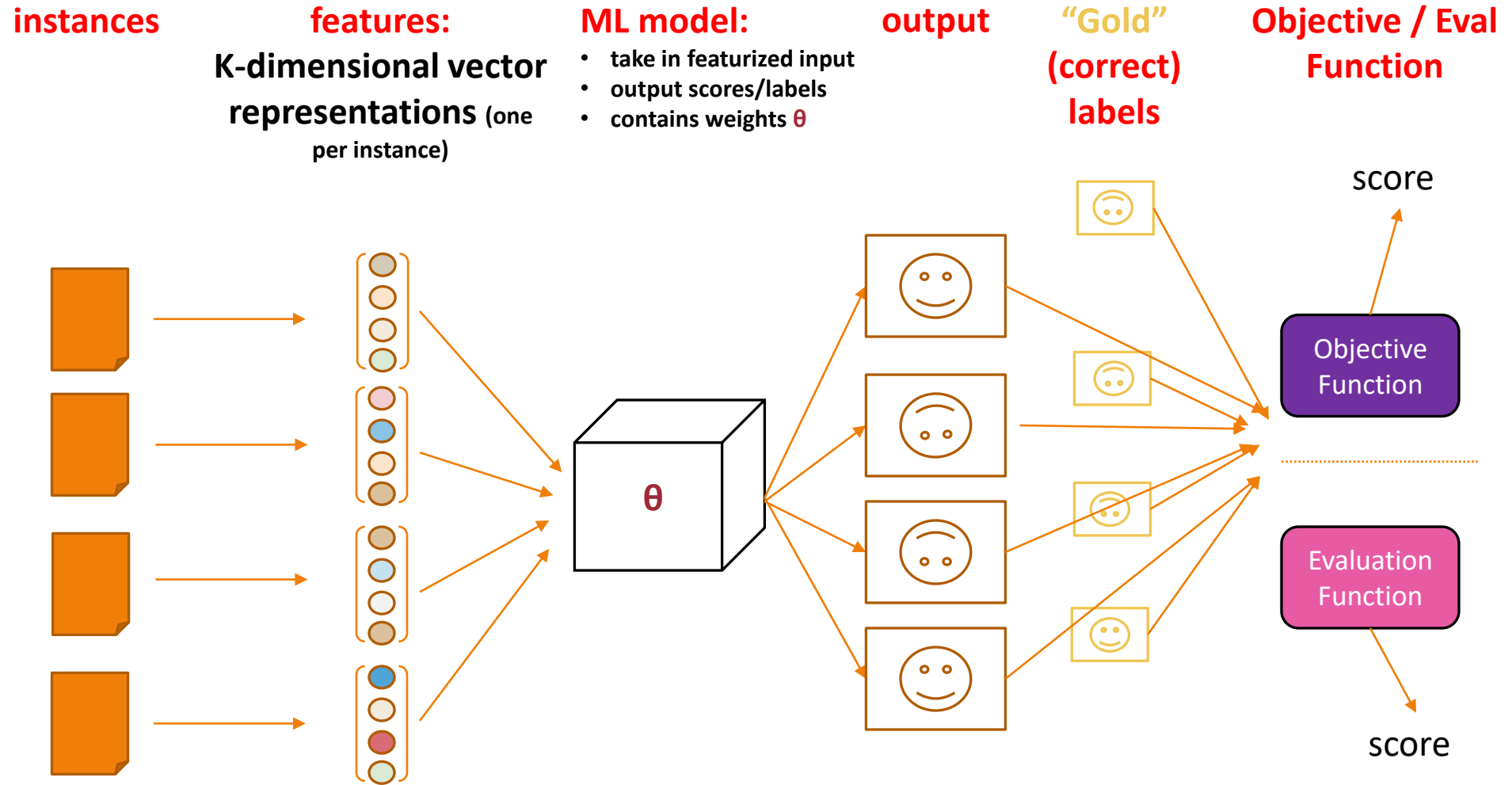
**output**

**“Gold”  
(correct)  
labels**

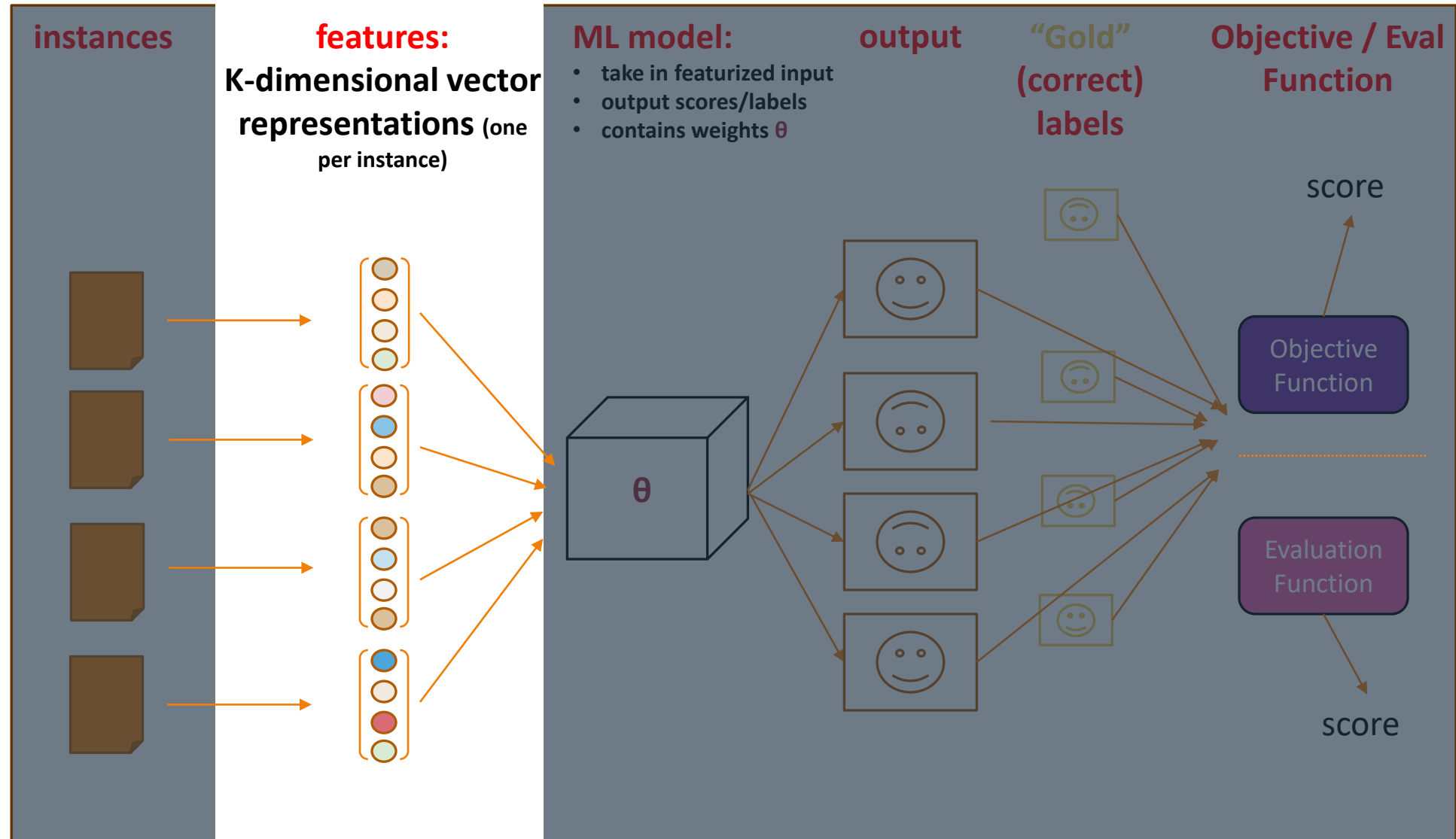
**Evaluation  
Function**



# ML/NLP Framework for Learning & Prediction



# First: Featurization / Encoding / Representation



# ML Term: “Featurization”

---

The procedure of extracting **features** for some input

Often viewed as a K-dimensional vector function  $f$  of the input language  $x$

$$f(x) = (f_1(x), \dots, f_K(x))$$



Each of these is a feature  
(/feature function)

# ML Term: “Featurization”

---

The procedure of extracting **features** for some input

Often viewed as a  $K$ -dimensional vector function  $f$  of the input language  $x$

$$f(x) = (f_1(x), \dots, f_K(x))$$

In supervised settings, it can equivalently be viewed as a  $K$ -dimensional vector function  $f$  of the input language  $x$  and a potential label  $y$

- $f(x, y) = (f_1(x, y), \dots, f_K(x, y))$

Features can be thought of as “soft” rules

- E.g., positive sentiments tweets may be *more likely* to have the word “happy”

# Defining Appropriate Features

---

Feature functions help extract useful features (characteristics) of the data

They turn data into numbers

Features that are not 0 are said to have fired

# Defining Appropriate Features

---

Feature functions help extract useful features (characteristics) of the data

They turn data into numbers

Features that are not 0 are said to have fired

You can define classes of features by templating (we'll come back to this!)

Often binary-valued (0 or 1), but can be real-valued

# Three Common Types of Featurization in NLP

---

1. Bag-of-words (or bag-of-characters, bag-of-relations)
2. Linguistically-inspired features
3. Dense features via embeddings



# Three Common Types of Featurization in NLP

---

1. Bag-of-words (or bag-of-characters, bag-of-relations)



- easy to define / extract
- sometimes still very useful

2. Linguistically-inspired features

3. Dense features via embeddings

# Three Common Types of Featurization in NLP

---

1. Bag-of-words (or bag-of-characters, bag-of-relations)



- easy to define / extract
- sometimes still very useful

2. Linguistically-inspired features



- harder to define
- helpful for interpretation
- depending on task: conceptually helpful
- currently, not freq. used

3. Dense features via embeddings

# Three Common Types of Featurization in NLP

---

1. Bag-of-words (or bag-of-characters, bag-of-relations)



- easy to define / extract
- sometimes still very useful

2. Linguistically-inspired features



- harder to define
- helpful for interpretation
- depending on task: conceptually helpful
- currently, not freq. used

3. Dense features via embeddings



- harder to define
- harder to extract (unless there's a model to run)
- currently: freq. used

# Three Common Types of Featurization in NLP

---

1. Bag-of-words (or bag-of-characters, bag-of-relations)
  - Identify **unique** sufficient atomic sub-parts (e.g., words in a document)
  - Define simple features over these, e.g.,
    - Binary (0 or 1) → indicating presence
    - Natural numbers → indicating number of times in a context
    - Real-valued → various other score (we'll see examples throughout the semester)
2. Linguistically-inspired features
3. Dense features via embeddings

# Example: Document Classification via Bag-of-Words Features

---

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

TECH  
NOT TECH

# Questions to consider...

---

- What are the input/output for this task?
- What might the features be?
- What types of applications could the task be used for?

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

TECH  
NOT TECH

# Questions to consider...

---

- **What are the input/output for this task?**
- What might the features be?
- What types of applications could the task be used for?

## Input

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

## Output

TECH  
NOT TECH

# Questions to consider...

---

- What are the input/output for this task?
- **What might the features be?**
- What types of applications could the task be used for?

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

TECH  
NOT TECH

Let's make a core assumption: the **label** can be predicted from **counts of individual word types**



# Example: Document Classification via Bag-of-Words Features

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

*feature extraction*

TECH

NOT TECH

With  $V$  word types,  
define  $V$  feature  
functions  $f_i(x)$  as

$f_i(x)$  = # of times word  
type  $i$  appears  
in document  $x$

# Example: Document Classification via Bag-of-Words Features

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

TECH

NOT TECH

*feature extraction*

$$f(x) = (f_i(x))_i^V$$

With  $V$  word types,  
define  $V$  feature  
functions  $f_i(x)$  as

$f_i(x)$  = # of times word  
type  $i$  appears  
in document  $x$

# Example: Document Classification via Bag-of-Words Features

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

*feature extraction*

TECH  
NOT TECH

feature $f_i(x)$	value
alerts	1
assist	1
bombing	1
Boston	2
...	
sniffle	0
...	

# Example: Document Classification via Bag-of-Words Features

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

TECH  
NOT TECH

$f(\mathbf{x})$ : "bag of words"

feature $f_i(x)$	value
alerts	1
assist	1
bombing	1
Boston	2
...	
sniffle	0
...	

$\mathbf{w}$ : weights

feature	weight
alerts	.043
assist	-0.25
bombing	0.8
Boston	-0.00001
...	

# Three Common Types of Featurization in NLP

---

1. Bag-of-words (or bag-of-characters, bag-of-relations)
  - Identify **unique** sufficient atomic sub-parts (e.g., words in a document)
  - Define simple features over these, e.g.,
    - Binary (0 or 1) → indicating presence
    - Natural numbers → indicating number of times in a context
    - Real-valued → various other score (we'll see examples throughout the semester)
2. Linguistically-inspired features
  - Define features from words, word spans, or linguistic-based annotations extracted from the document
3. Dense features via embeddings

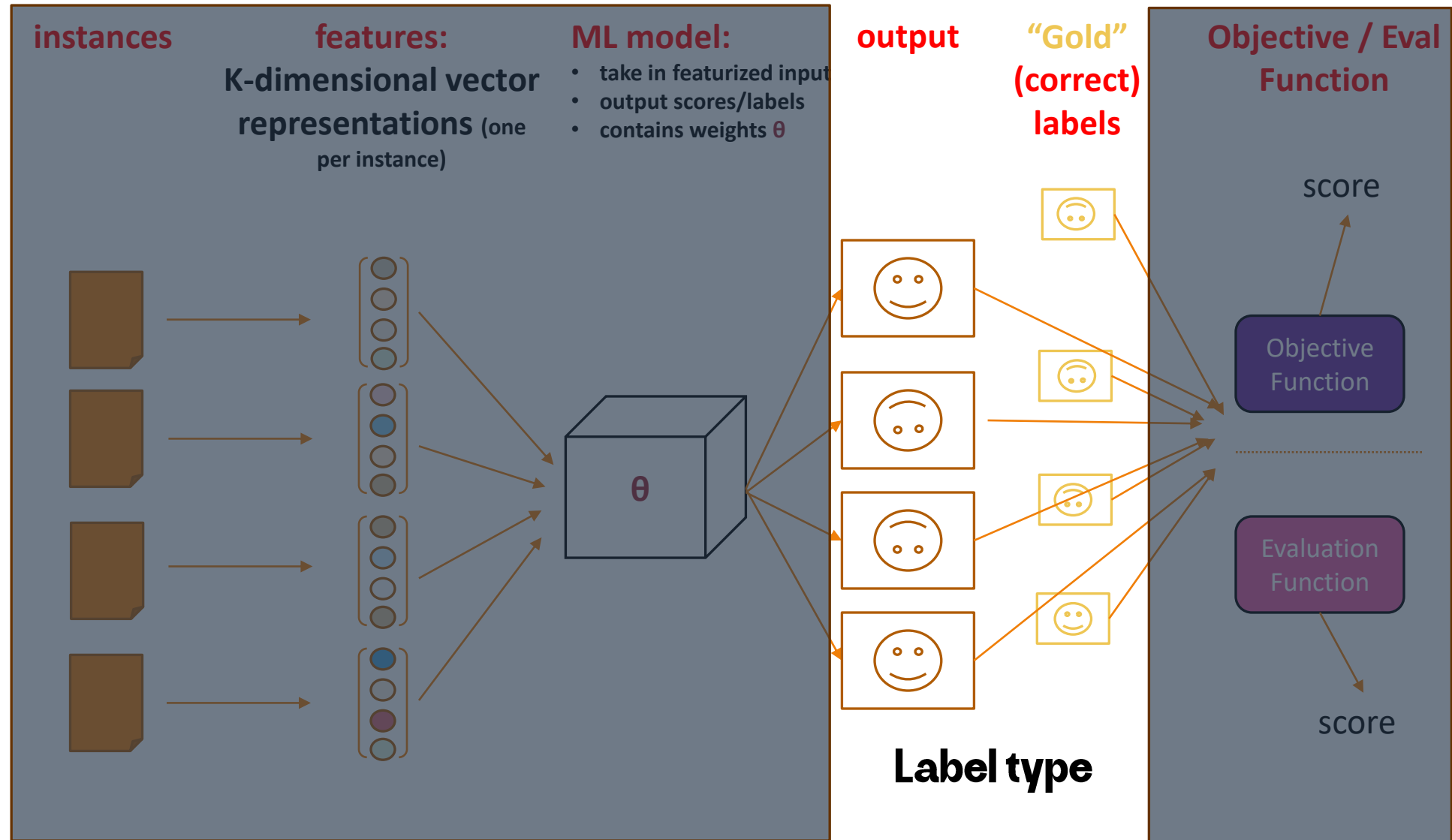
# Three Common Types of Featurization in NLP

---

1. Bag-of-words (or bag-of-characters, bag-of-relations)
  - Identify **unique** sufficient atomic sub-parts (e.g., words in a document)
  - Define simple features over these, e.g.,
    - Binary (0 or 1) → indicating presence
    - Natural numbers → indicating number of times in a context
    - Real-valued → various other score (we'll see examples throughout the semester)
2. Linguistically-inspired features
  - Define features from words, word spans, or linguistic-based annotations extracted from the document
3. Dense features via embeddings
  - Compute/extract a real-valued vector, e.g., from word2vec, ELMO, BERT, ...

Will be  
discussed  
in a future  
lecture

# Second: Classification Terminology



# Classification Types (Terminology)

---

Name	Number of Tasks (Domains) Labels are Associated with	# Label Types	Example
(Binary) Classification			
Multi-class Classification			
Multi-label Classification			
Multi-task Classification			



# Classification Types (Terminology)

---

Name	Number of Tasks (Domains) Labels are Associated with	# Label Types	Example
(Binary) Classification	1	2	Sentiment: Choose one of {positive or negative}
Multi-class Classification			
Multi-label Classification			
Multi-task Classification			

# Classification Types (Terminology)

---

Name	Number of Tasks (Domains) Labels are Associated with	# Label Types	Example
(Binary) Classification	1	2	Sentiment: Choose one of {positive or negative}
Multi-class Classification	1	> 2	Part-of-speech: Choose one of {Noun, Verb, Det, Prep, ...}
Multi-label Classification			
Multi-task Classification			

# Classification Types (Terminology)

---

Name	Number of Tasks (Domains) Labels are Associated with	# Label Types	Example
(Binary) Classification	1	2	Sentiment: Choose one of {positive or negative}
Multi-class Classification	1	> 2	Part-of-speech: Choose one of {Noun, Verb, Det, Prep, ...}
Multi-label Classification	1	> 2	Sentiment: Choose multiple of {positive, angry, sad, excited, ...}
Multi-task Classification			

# Classification Types (Terminology)

Name	Number of Tasks (Domains) Labels are Associated with	# Label Types	Example
(Binary) Classification	1	2	Sentiment: Choose one of {positive or negative}
Multi-class Classification	1	> 2	Part-of-speech: Choose one of {Noun, Verb, Det, Prep, ...}
Multi-label Classification	1	> 2	Sentiment: Choose multiple of {positive, angry, sad, excited, ...}
Multi-task Classification	> 1	Per task: 2 or > 2 (can apply to binary or multi-class)	Task 1: part-of-speech Task 2: named entity tagging ... ----- Task 1: document labeling Task 2: sentiment

# Text Annotation Tasks ("Classification" Tasks)

---

1. Classify the entire document ("text categorization")
2. Classify word tokens individually
3. Classify word tokens in a sequence
4. Identify phrases ("chunking")
5. Syntactic annotation (parsing)
6. Semantic annotation
7. Text generation

*Slide courtesy Jason Eisner, with mild edits*

# Text Annotation Tasks ("Classification" Tasks)

---

1. Classify the entire document ("text categorization")
2. Classify word tokens individually
3. Classify word tokens in a sequence
4. Identify phrases ("chunking")
5. Syntactic annotation (parsing)
6. Semantic annotation

*Slide courtesy Jason Eisner, with mild edits*

# Questions to consider...

---

- What are the input/output for this task?
- What might the features be?
- **What types of applications could the task be used for?**

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

TECH

NOT TECH

# Text Classification

---

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Language Identification

Sentiment analysis

...



# Text Classification

Assigning subject categories, topics, or genres

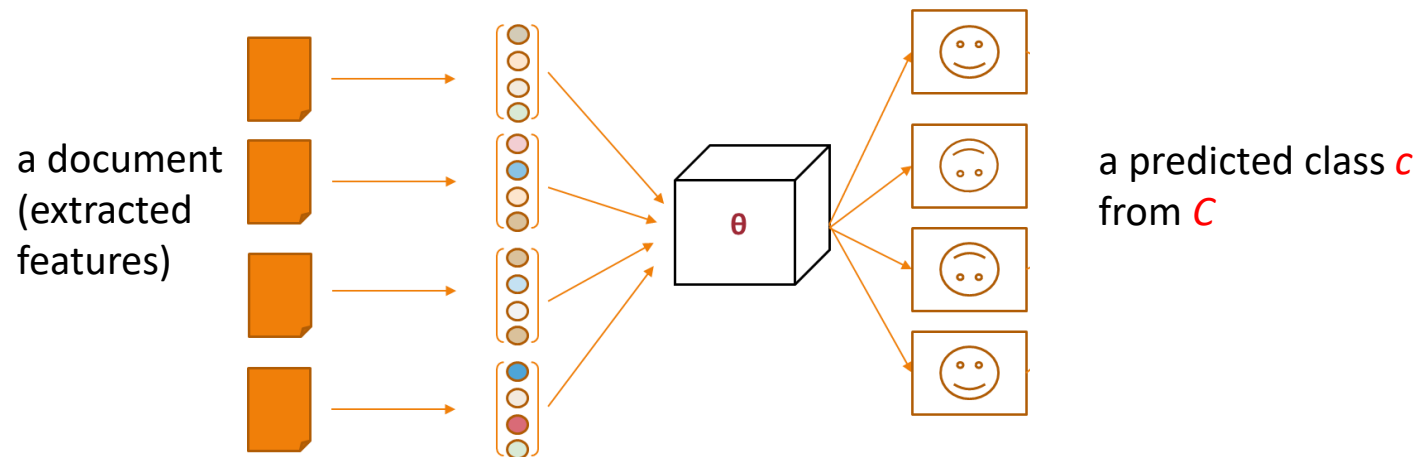
Spam detection

Authorship identification

Language Identification

Sentiment analysis

...



# Text Classification: Hand-coded Rules?

---

Assigning subject categories, topics, or genres

Language Identification

Sentiment analysis

Spam detection

...

Authorship identification

Rules based on combinations of words or other features

spam: black-list-address OR (“dollars” AND “have been selected”)

Accuracy can be high

If rules carefully refined by expert

Building and maintaining these rules is expensive

Can humans faithfully assign uncertainty?

# Text Classification: Supervised Machine Learning

---

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Language Identification

Sentiment analysis

...

a fixed set of classes

$C = \{c_1, c_2, \dots, c_J\}$

a training set of  $m$  hand-labeled documents  $D$  with corresponding labels  $(d_1, y_1), \dots, (d_m, y_m), y \in C$



“Training Process”



a learned classifier  $\gamma$  that maps documents to classes

# Questions to consider...

---

- What are the input/output for this task?
- What might the features be?
- What types of applications could the task be used for?

## Input

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after in 2013, when were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

## Output

TECH

An alternate view of this is...

# Text Classification: Supervised Machine Learning - Training

Assigning subject categories, topics, or genres

Spam detection

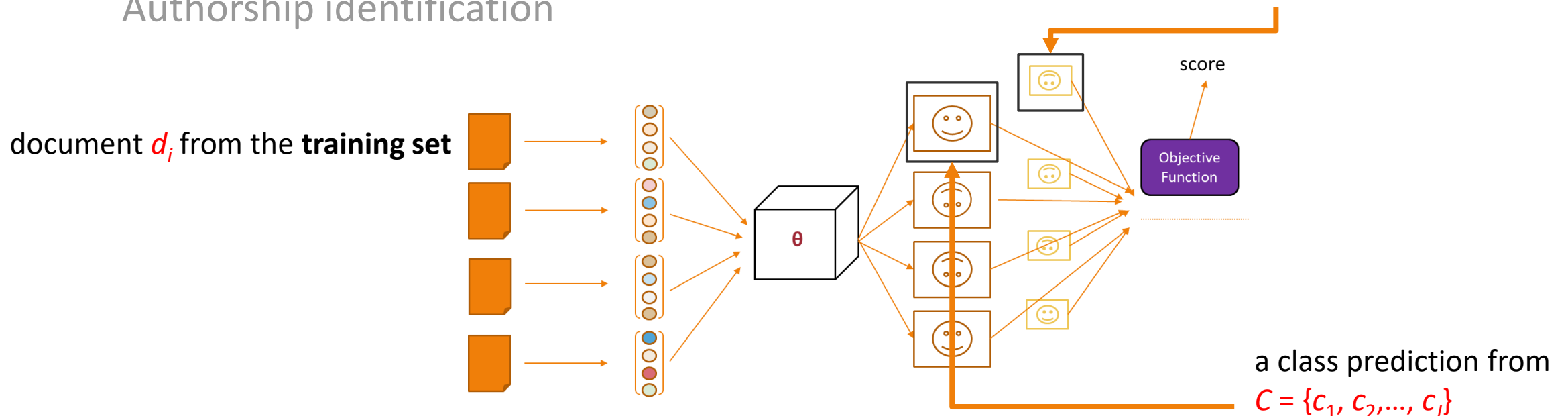
Authorship identification

Language Identification

Sentiment analysis

...

$y_i$  corresponding to the gold label for  $d_i$



# Text Classification: Supervised Machine Learning - Testing

Assigning subject categories, topics, or genres

Language Identification

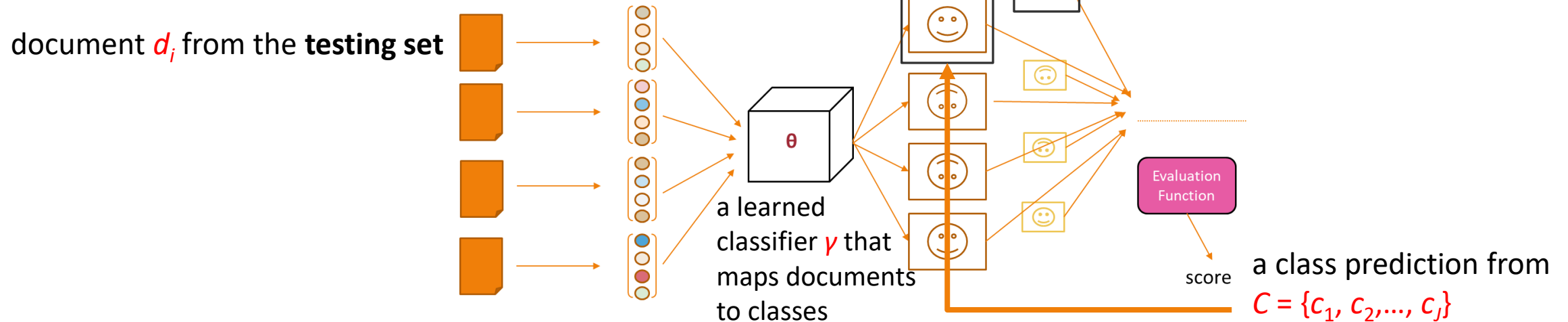
Sentiment analysis

Spam detection

Authorship identification

...

$y_i$  corresponding to the gold label for  $d_i$



# Text Classification: Supervised Machine Learning – Model examples

Assigning subject categories, topics, or genres

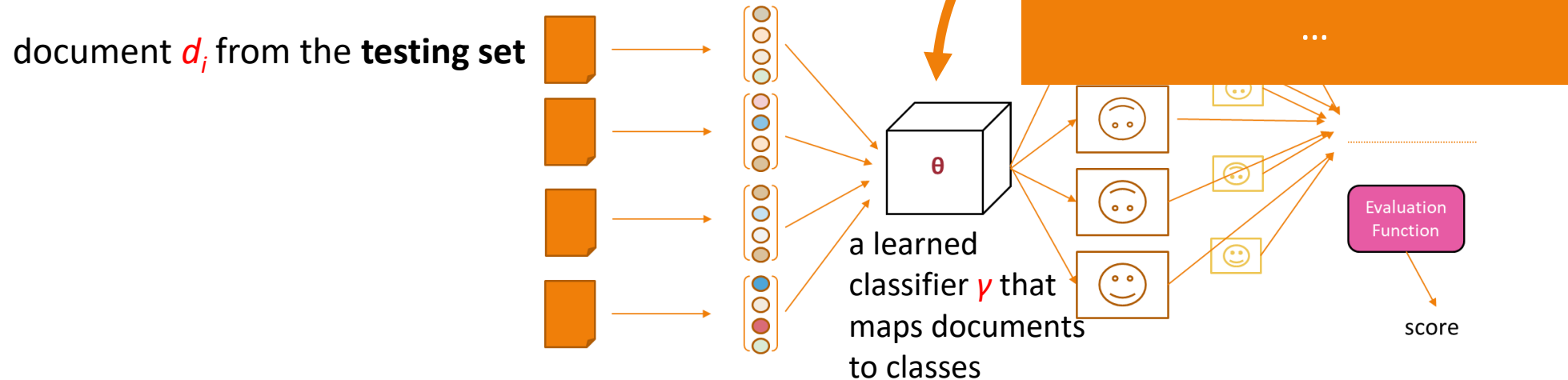
Spam detection

Authorship identification

Language Identification

Sentiment analysis

...



# Knowledge Check: Handling Types and Tokens

- 10 minutes to do it in class
- You can complete it after class
- Then submit it to Blackboard
- I'll release my answer 2/13 (please finish before then)

CMSC 473/673 NLP @ UMBC

About Schedule Homework ▾ Knowledge Checks ▾

2/5 - Handling Types and Tokens

## CMSC 473/673 Natural Language Processing at UMBC Spring 2026

**Jump to class policies:** [\[Late Day\]](#) [\[Academic Integrity\]](#) [\[Generative AI\]](#) [\[GitHub Use\]](#) [\[Collaboration\]](#)

### Course Description

Natural language processing (NLP) is the field of working with language to automatically perform a variety of tasks, instead of or in collaboration with people. NLP can focus on the Generation (NLG) and/or Understanding (NLU) of natural language. Recently, large language models (LLMs) like ChatGPT have gotten the attention of the general public, but they have also greatly changed the landscape of modern NLP research. This course will show you both old & new techniques that are still used today and will give you a basic understanding of why & how we do NLP.

### Learning Objectives

By the end of the course, you will be able to...



# Text Annotation Tasks ("Classification" Tasks)

---

1. Classify the entire document ("text categorization")
2. Classify word tokens individually
3. Classify word tokens in a sequence
4. Identify phrases ("chunking")
5. Syntactic annotation (parsing)
6. Semantic annotation
7. Text generation

*Slide courtesy Jason Eisner, with mild edits*

## Word Sense Disambiguation (WSD)

### Problem:

The company said the *plant* is still operating ...

⇒ (A) Manufacturing plant    or

⇒ (B) Living plant

**Training Data:**    Build a special classifier just for tokens of “plant”

Sense	Context
<b>(1) Manufacturing</b>	... union responses to <i>plant</i> closures . ...
” ”	... computer disk drive <i>plant</i> located in ...
” ”	company manufacturing <i>plant</i> is in Orlando ...
<b>(2) Living</b>	... animal rather than <i>plant</i> tissues can be ...
” ”	... to strain microscopic <i>plant</i> life from the ...
” ”	and Golgi apparatus of <i>plant</i> and animal cells

### Test Data:

Sense	Context
???	... vinyl chloride monomer <i>plant</i> , which is ...
???	... molecules found in <i>plant</i> tissue from the ...

slide courtesy of D. Yarowsky (modified)

## WSD for Machine Translation (English → Spanish)

### Problem:

... He wrote the last **sentence** two years later ...

⇒ *sentencia* (legal sentence) or

⇒ *frase* (grammatical sentence)

**Training Data:** Build a special classifier just for tokens of “sentence”

Translation	Context
<b>(1) sentencia</b>	... for a maximum <i>sentence</i> for a young offender ...
” ”	... of the minimum <i>sentence</i> of seven years in jail ...
” ”	... were under the <i>sentence</i> of death at that time ...
<b>(2) frase</b>	... read the second <i>sentence</i> because it is just as ...
” ”	... The next <i>sentence</i> is a very important ...
” ”	... It is the second <i>sentence</i> which I think is at ...

### Test Data:

Translation	Context
???	... cannot criticize a <i>sentence</i> handed down by ...
???	... listen to this <i>sentence</i> uttered by a former ...

slide courtesy of D. Yarowsky (modified)

## Accent Restoration in Spanish & French

### Problem:

**Input:** ... deja travaille cote a cote ...



**Output:** ... déjà travaillé côte à côte ...

### Examples:

... appeler l'autre **cote** de l'atlantique ...

⇒ *côté* (meaning side) or

⇒ *côte* (meaning coast)

... une famille des **pecheurs** ...

⇒ *pêcheurs* (meaning fishermen) or

⇒ *pécheurs* (meaning sinners)

## Accent Restoration in Spanish & French

### Training Data:

Pattern	Context
(1) côté	... du laisser de <i>cote</i> faute de temps ...
” ”	... appeler l’ autre <i>cote</i> de l’ atlantique ...
” ”	... passe de notre <i>cote</i> de la frontiere ...
(2) côte	... vivre sur notre <i>cote</i> ouest toujours ...
” ”	... creer sur la <i>cote</i> du labrador des ...
” ”	travaillaient cote a <i>cote</i> , ils avaient ...

### Test Data:

Pattern	Context
???	... passe de notre <i>cote</i> de la frontiere ...
???	... creer sur la <i>cote</i> du labrador des ...

slide courtesy of D. Yarowsky (modified)

## Text-to-Speech Synthesis

### Problem:

... slightly elevated *lead* levels ...

⇒ *lɛd* (as in *lead mine*)    or

⇒ *li:d* (as in *lead role*)

### Training Data:

Pronunciation	Context
<b>(1) lɛd</b>	... it monitors the <i>lead</i> levels in drinking ...
” ”	... conference on <i>lead</i> poisoning in ...
” ”	... strontium and <i>lead</i> isotope zonation ...
<b>(2) li:d</b>	... maintained their <i>lead</i> Thursday over ...
” ”	... to Boston and <i>lead</i> singer for Purple ...
” ”	... Bush a 17-point <i>lead</i> in Texas , only 3 ...

### Test Data:

Pronunciation	Context
???	... median blood <i>lead</i> concentration was ..
???	... his double-digit <i>lead</i> nationwide . The ...

slide courtesy of D. Yarowsky (modified)

## Spelling Correction

### Problem:

... and he fired presidential **aid/aide** Dick Morris after ...

⇒ *aid* or

⇒ *aide*

### Training Data:

Spelling	Context
<b>(1) aid</b>	... and cut the foreign <i>aid/aide</i> budget in fiscal 1996 ...
” ”	... they offered federal <i>aid/aide</i> for flood-ravaged states ...
<b>(2) aide</b>	... fired presidential <i>aid/aide</i> Dick Morris after ...
” ”	... and said the chief <i>aid/aide</i> to Sen. Baker, Mr. John ...

### Test Data:

Spelling	Context
???	... said the longtime <i>aid/aide</i> to the Mayor of St. ...
???	... will squander the <i>aid/aide</i> it receives from the ...

slide courtesy of D. Yarowsky (modified)

# What features? Example: “word to [the] left [of correction]”

Word to left	Frequency as <b>Aid</b>	Frequency as <b>Aide</b>
foreign	718	1
federal	297	0
western	146	0
provide	88	0
covert	26	0
oppose	13	0
future	9	0
similar	6	0
presidential	0	63
chief	0	40
longtime	0	26
aids-infected	0	2
sleepy	0	1
disaffected	0	1
indispensable	2	1
practical	2	0
squander	1	0

Spelling correction using an n-gram language model ( $n \geq 2$ ) would use words to left and right to help predict the true word.

Similarly, an HMM would predict a word's class using classes to left and right.

But we'd like to throw in all kinds of other features, too ...

*slide courtesy of D. Yarowsky (modified)*



# An assortment of possible cues ...

	Position	Collocation	led	li:d
<b>N-grams</b>  (word, lemma, part-of-speech)	+1 L	lead <i>level/N</i>	219	0
	-1 W	<i>narrow</i> lead	0	70
	+1 W	lead <i>in</i>	207	898
	-1 W,+1 W	<i>of</i> lead <i>in</i>	162	0
	-1 W,+1 W	<i>the</i> lead <i>in</i>	0	301
	+1 P,+2 P	lead , < <i>NOUN</i> >	234	7
<b>Wide-context collocations</b>	$\pm k$ W	<i>zinc</i> (in $\pm k$ words)	235	0
	$\pm k$ W	<i>copper</i> (in $\pm k$ words)	130	0
<b>Verb-object relationships</b>	-V L	<i>follow/V</i> + lead	0	527
	-V L	<i>take/V</i> + lead	1	665

generates a whole bunch of potential cues – use data to find out which ones work best

Word to left	Frequency as Aid	Frequency as Aide
foreign	718	1
federal	297	0
western	146	0
provide	88	0

slide courtesy of D. Yarowsky (modified)

# An assortment of possible cues ...

	Position	Collocation	l <sub>ed</sub>	li:d
<b>N-grams</b>  (word, lemma, part-of-speech)	+1 L	lead <i>level/N</i>	219	0
	-1 W	<i>narrow</i> lead	0	70
	+1 W	lead <i>in</i>	207	898
	-1 W,+1 W	<i>of</i> lead <i>in</i>	162	0
	-1 W,+1 W	<i>the</i> lead <i>in</i>	0	301
	+1 P,+2 P	lead , < <i>NOUN</i> >	234	7
<b>Wide-context collocations</b>	±k W	<i>zinc</i> (in ±k words)	235	0
	±k W	<i>copper</i> (in ±k words)	130	0
<b>Verb-object relationships</b>	-V L	<i>follow/V</i> + lead	0	527
	-V L	<i>take/V</i> + lead	1	665

This feature is relatively weak, but weak features are still useful, especially since very few features will fire in a given context.

merged ranking  
of all cues  
of all these types

11.40	<i>follow/V</i> + lead	⇒ li:d
11.20	<i>zinc</i> (in ±k words)	⇒ l <sub>ed</sub>
11.10	lead <i>level/N</i>	⇒ l <sub>ed</sub>
10.66	<i>of</i> lead <i>in</i>	⇒ l <sub>ed</sub>
10.59	<i>the</i> lead <i>in</i>	⇒ li:d
10.51	lead <i>role</i>	⇒ li:d

slide courtesy of D. Yarowsky (modified)

# Final decision list for *lead* (abbreviated)

What are the input/output?  
What are the features?  
What types of applications?

List of all features,  
ranked by their weight.

(These weights are for a simple  
“decision list” model where the single  
highest-weighted feature that fires  
gets to make the decision all by itself.

However, a log-linear model, which  
adds up the weights of all features  
that fire, would be roughly similar.)

LogL	Evidence	Pronunciation
11.40	<i>follow/V</i> + lead	⇒ li:d
11.20	<i>zinc</i> (in $\pm k$ words)	⇒ lɛd
11.10	lead <i>level/N</i>	⇒ lɛd
10.66	<i>of</i> lead <i>in</i>	⇒ lɛd
10.59	<i>the</i> lead <i>in</i>	⇒ li:d
10.51	lead <i>role</i>	⇒ li:d
10.35	<i>copper</i> (in $\pm k$ words)	⇒ lɛd
10.28	lead <i>time</i>	⇒ li:d
10.24	lead <i>levels</i>	⇒ lɛd
10.16	lead <i>poisoning</i>	⇒ lɛd
8.55	<i>big</i> lead	⇒ li:d
8.49	<i>narrow</i> lead	⇒ li:d
7.76	<i>take/V</i> + lead	⇒ li:d
5.99	lead , <i>NOUN</i>	⇒ lɛd
1.15	lead <i>in</i>	⇒ li:d
	○ ○ ○	

slide courtesy of D. Yarowsky (modified)

# Text-to-Speech Synthesis

## Problem:

... slightly elevated *lead* levels ...

⇒ *lɛd* (as in *lead mine*)    or

⇒ *li:d* (as in *lead role*)

## Training Data:

Pronunciation	Context
<b>(1) lɛd</b>	... it monitors the <i>lead</i> levels in drinking ...
” ”	... conference on <i>lead</i> poisoning in ...
” ”	... strontium and <i>lead</i> isotope zonation ...
<b>(2) li:d</b>	... maintained their <i>lead</i> Thursday over ...
” ”	... to Boston and <i>lead</i> singer for Purple ...
” ”	... Bush a 17-point <i>lead</i> in Texas , only 3 ...

## Test Data:

Pronunciation	Context
???	... median blood <i>lead</i> concentration was ..
???	... his double-digit <i>lead</i> nationwide . The ...

*slide courtesy of D. Yarowsky (modified)*

# Token Classification

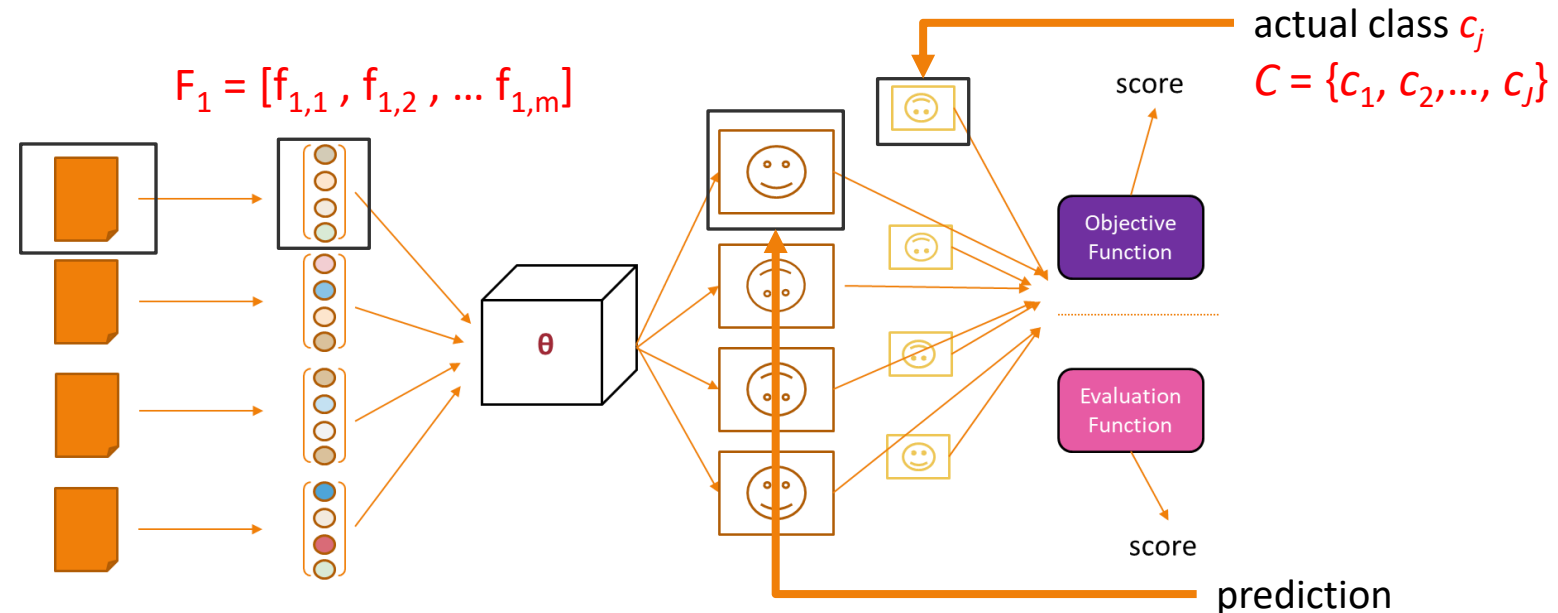
Word pronunciation

Word sense disambiguation (WSD)  
within or across languages

Accent restoration

...

features  $F_1$  extracted from  
word  $w_1$  and its surrounding  
words (context)



# Text Annotation Tasks ("Classification" Tasks)

---

1. Classify the entire document ("text categorization")
2. Classify word tokens individually
3. Classify word tokens in a sequence (i.e., order matters)
4. Identify phrases ("chunking")
5. Syntactic annotation (parsing)
6. Semantic annotation
7. Text generation

*Slide courtesy Jason Eisner, with mild edits*

# Example: Part of Speech Tagging

---

We could treat tagging as a token classification problem

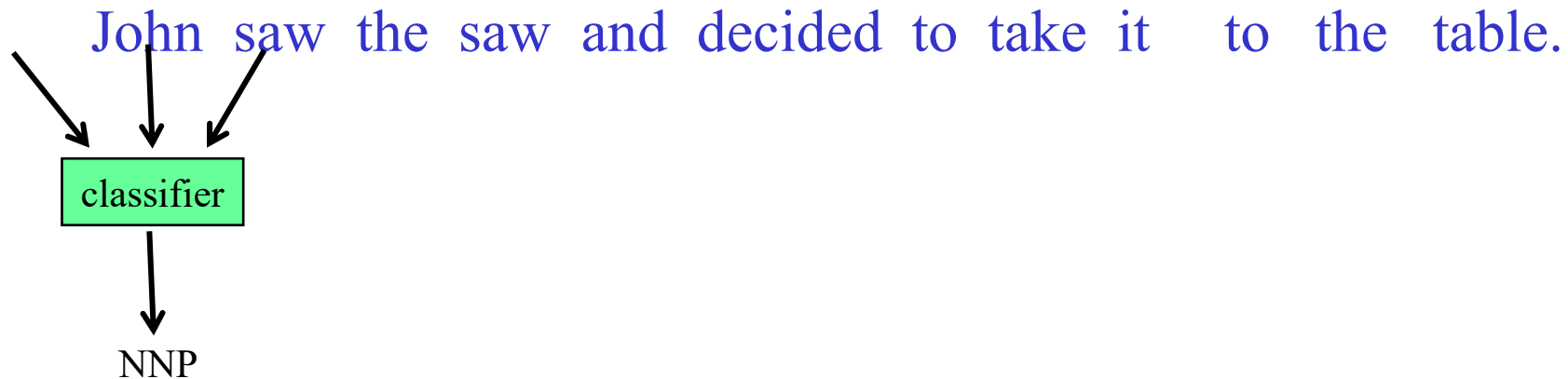
- Tag each word independently given features of context
- And features of the word's spelling (suffixes, capitalization)

*Slide courtesy Jason Eisner, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).



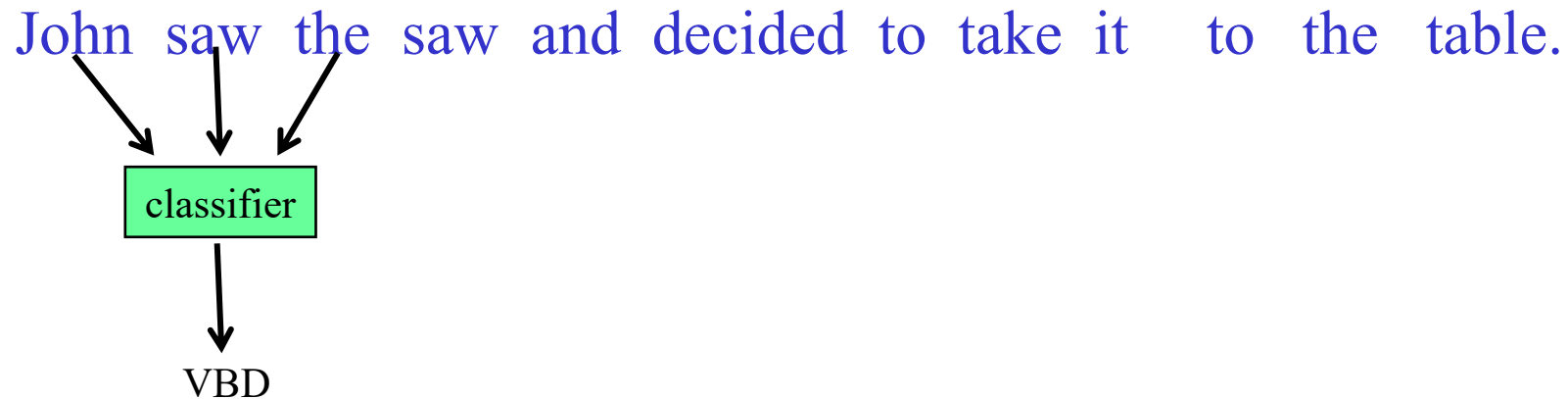
*Slide courtesy Ray Mooney, with mild edits*



# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

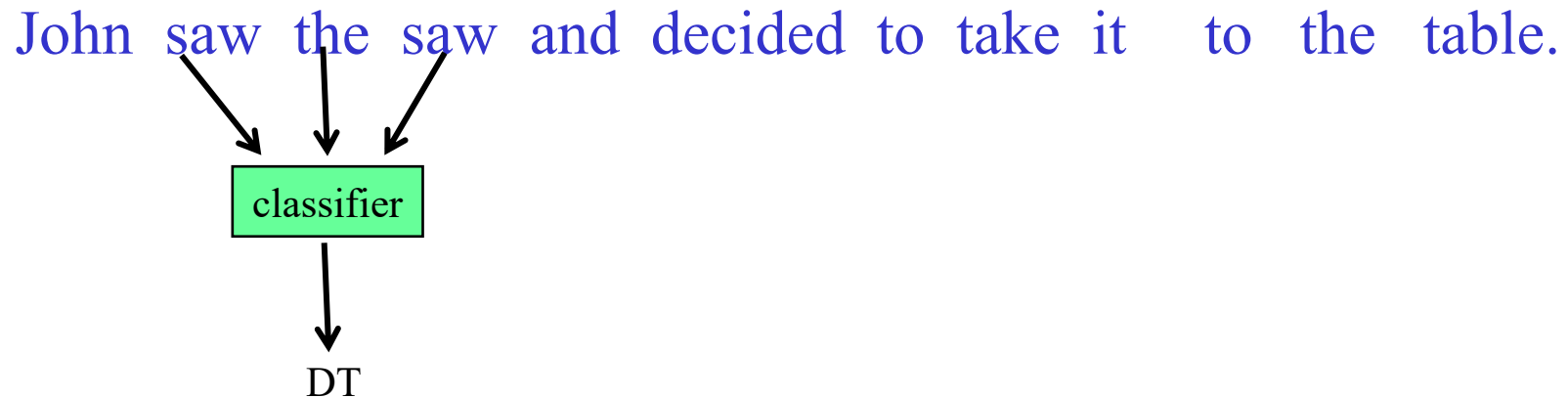


*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

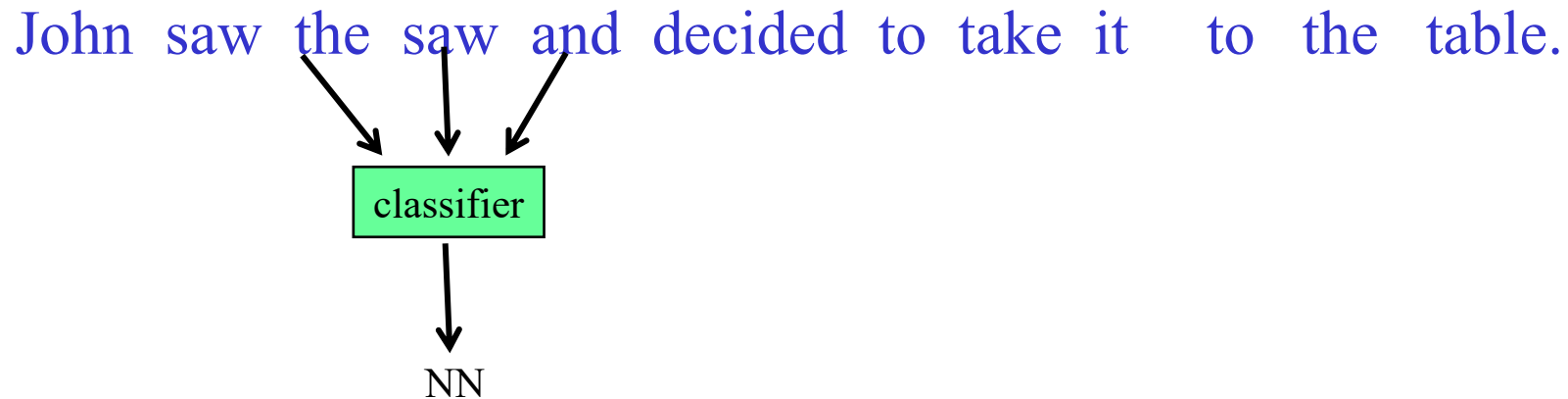


*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

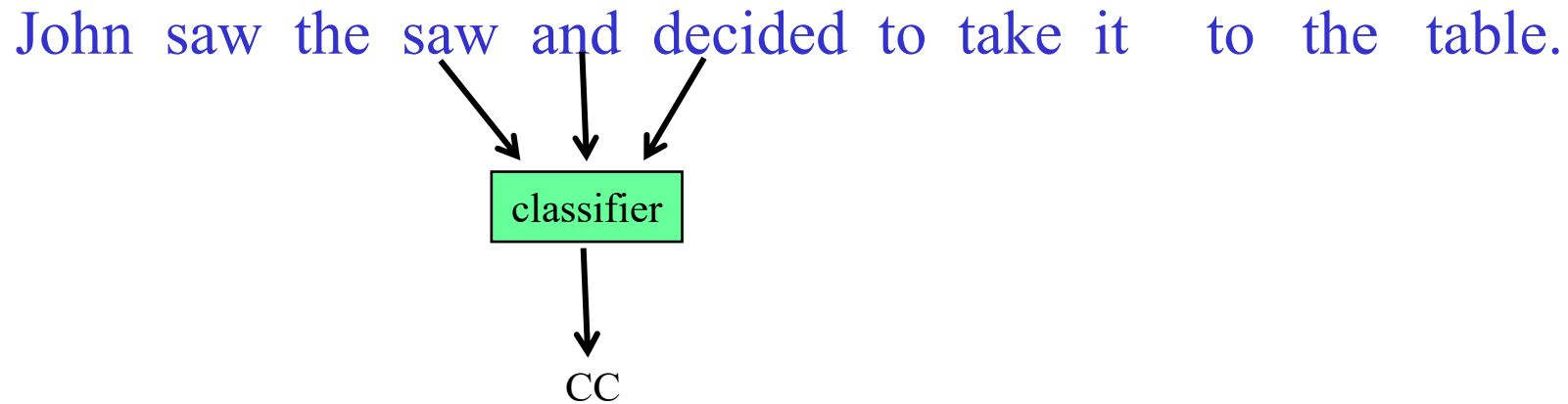


*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

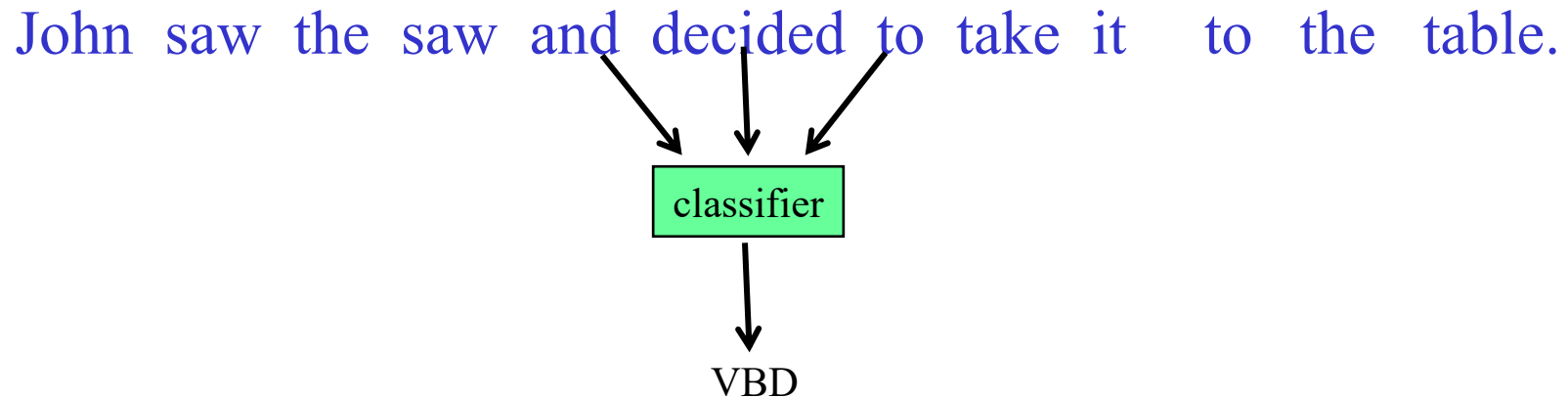


*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).



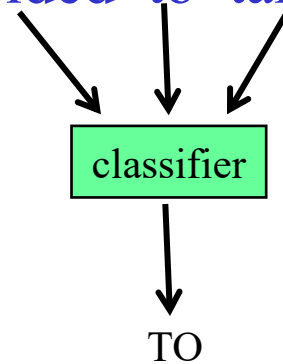
*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.

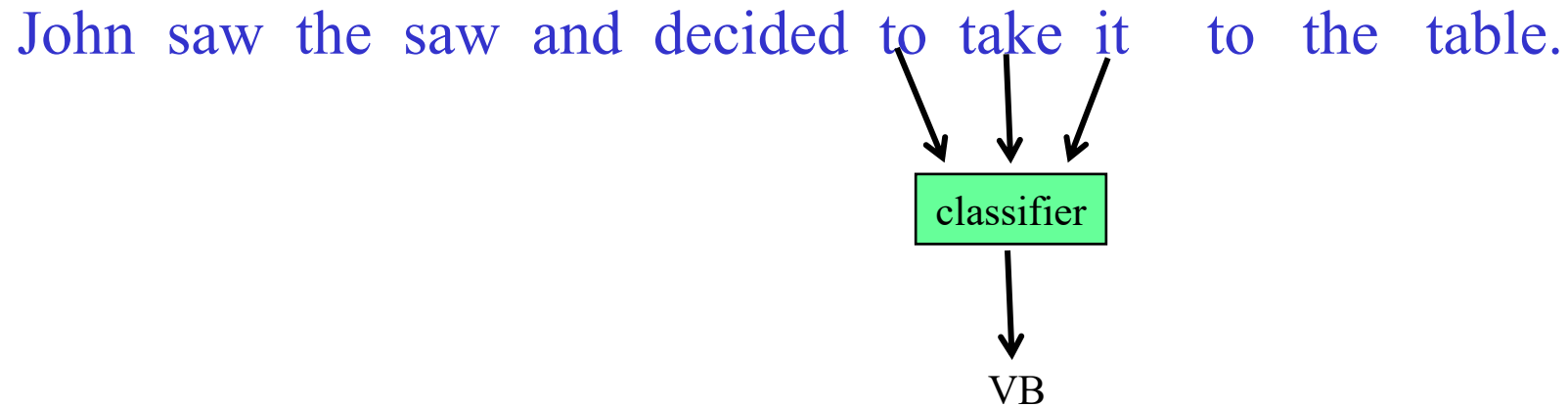


*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).



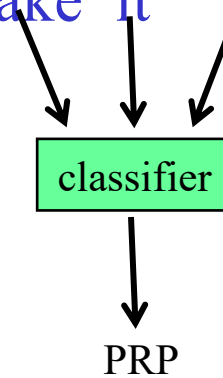
*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.



*Slide courtesy Ray Mooney, with mild edits*

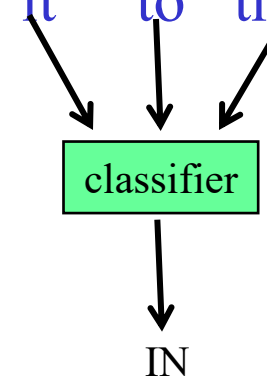


# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.

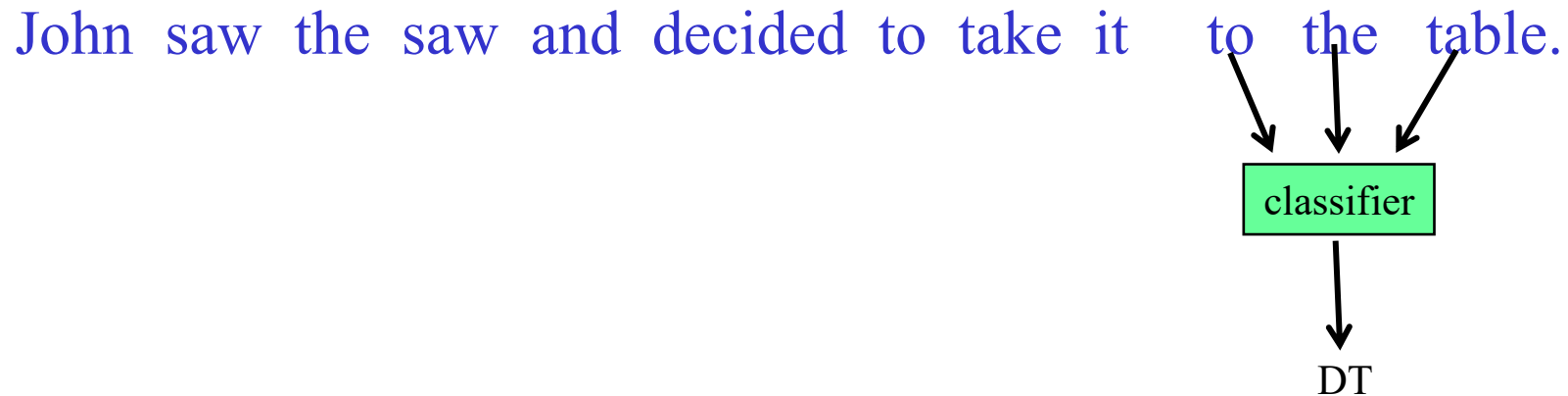


*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

---

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

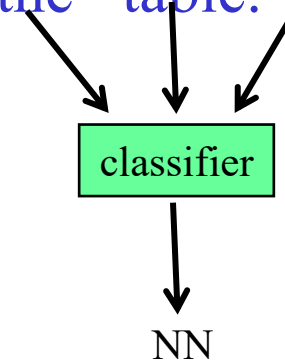


*Slide courtesy Ray Mooney, with mild edits*

# Sequence Labeling as Classification

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.



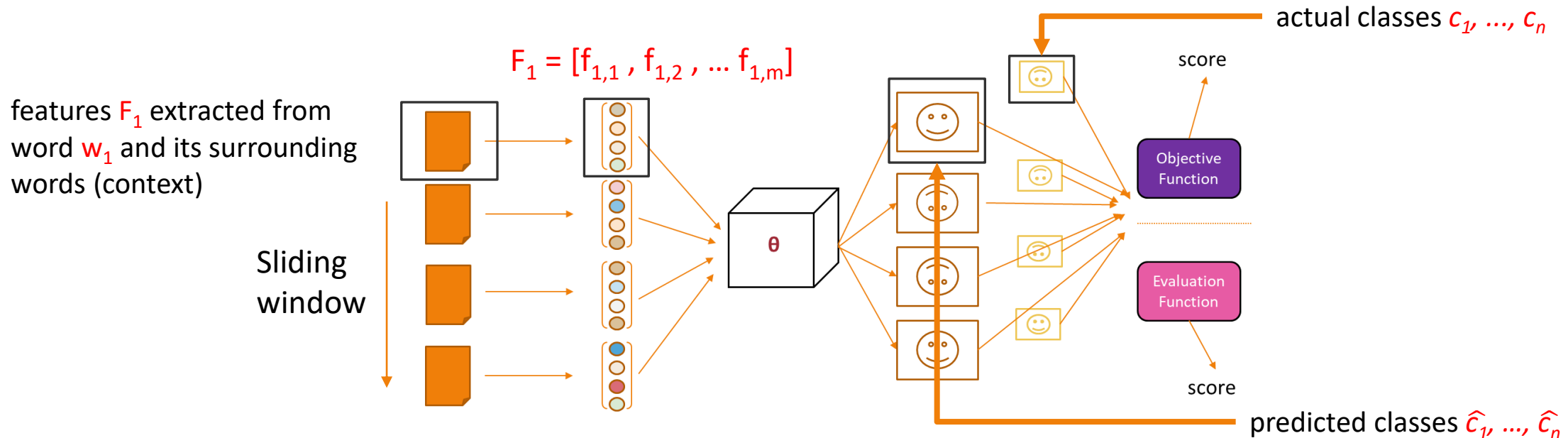
What are the input/output?  
What are the features?  
What types of applications?

*Slide courtesy Ray Mooney, with mild edits*

# Token Classification in a Sequence

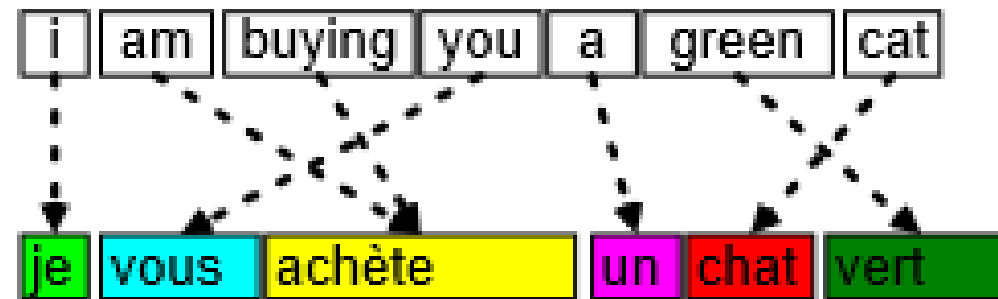
Part of speech tagging

Word alignment



# Machine Translation: Word Alignment

---



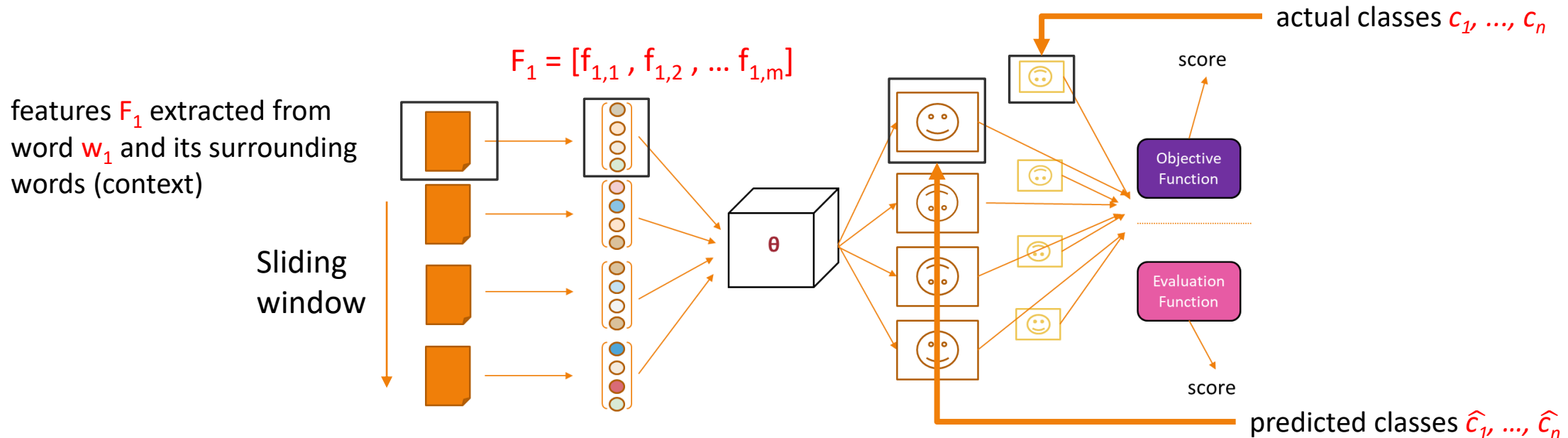
What kinds of features might we want to consider here?

# Token Classification in a Sequence

Part of speech tagging

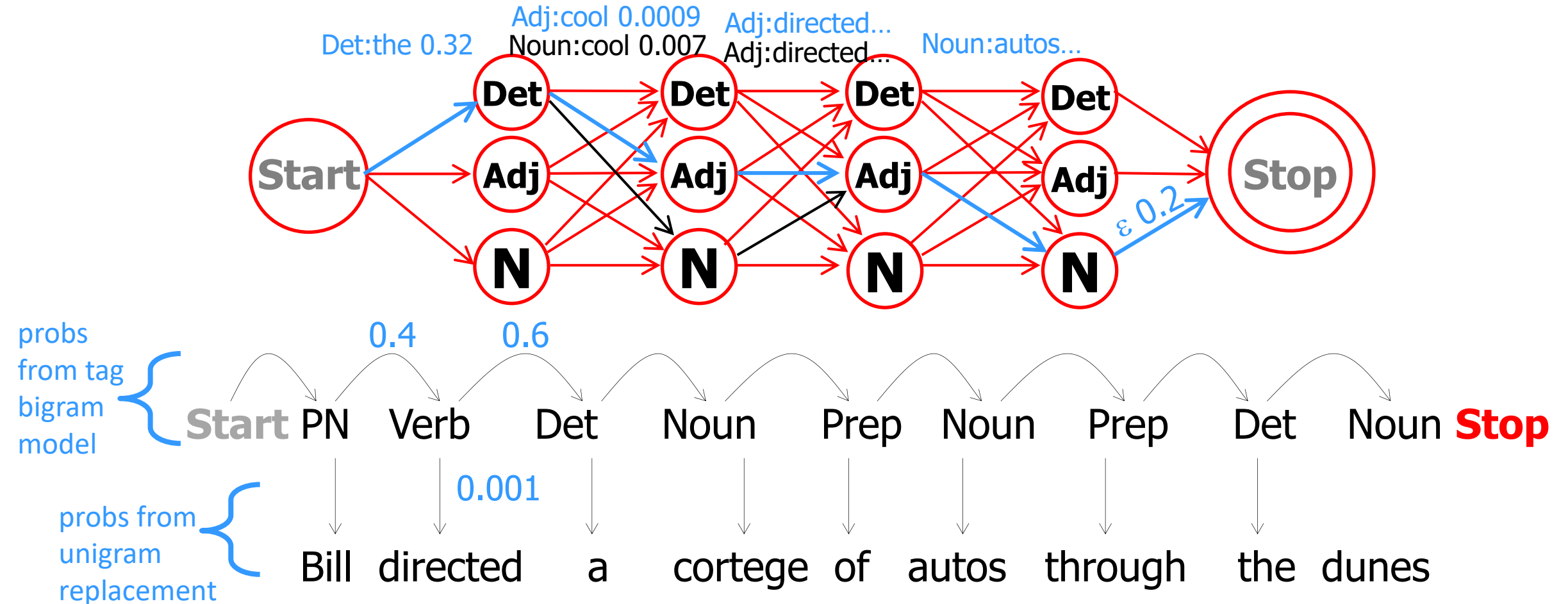
Word alignment

**Other examples?**



# Part of Speech Tagging

Or we could use an HMM:



Slide courtesy Jason Eisner, with mild edits

# Part of Speech Tagging

---

We could treat tagging as a token classification problem

- Tag each word independently given features of context
- And features of the word's spelling (suffixes, capitalization)

Or we could use an HMM:

- The point of the HMM is basically that the tag of one word might depend on the tags of adjacent words.

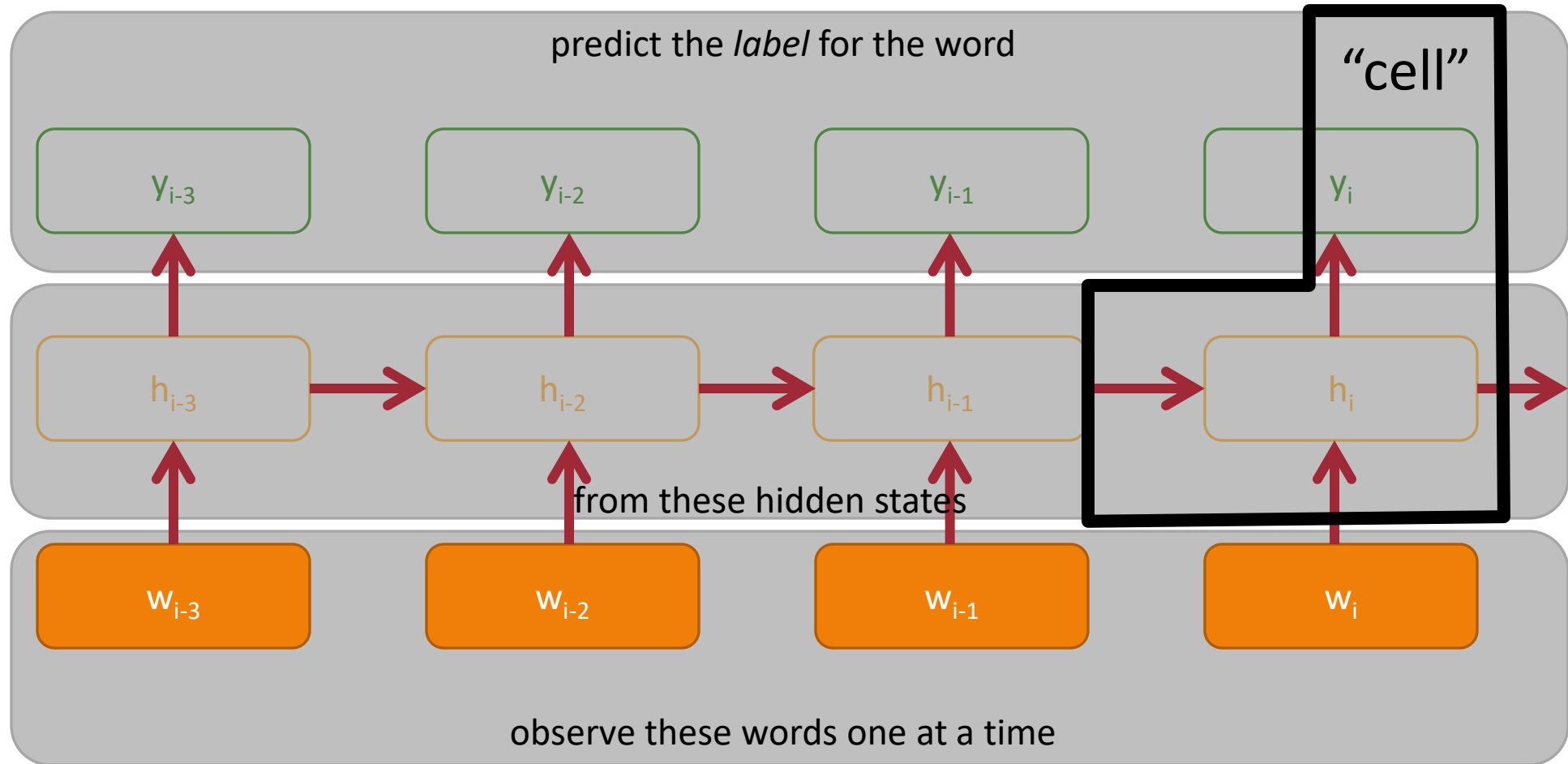
Combine these two ideas??

- We'd like rich features (e.g., in a **log-linear model**), but we'd also like our feature functions to depend on adjacent tags.
- So, the problem is to predict **all** tags together.

*Slide courtesy Jason Eisner, with mild edits*



# Can We Use Neural, Recurrent Methods for PoS Tagging?



# Text Annotation Tasks ("Classification" Tasks)

---

1. Classify the entire document ("text categorization")
2. Classify word tokens individually
3. Classify word tokens in a sequence
4. Identify phrases ("chunking")
5. Syntactic annotation (parsing)
6. Semantic annotation
7. Text generation

*Slide courtesy Jason Eisner, with mild edits*

# Example: Finding Named Entities

---

Named entity recognition (NER)

Identify proper names in texts, and classification into a set of predefined categories of interest

- Person names
- Organizations (companies, government organisations, committees, etc.)
- Locations (cities, countries, rivers, etc.)
- Date and time expressions
- Measures (percent, money, weight, etc.),
- email addresses, web addresses, street addresses, etc.
- Domain-specific: names of drugs, medical conditions,
- names of ships, bibliographic references etc.

# NE Types

Type	Tag	Sample Categories
People	PER	Individuals, fictional characters, small groups
Organization	ORG	Companies, agencies, political parties, religious groups, sports teams
Location	LOC	Physical extents, mountains, lakes, seas
Geo-Political Entity	GPE	Countries, states, provinces, counties
Facility	FAC	Bridges, buildings, airports
Vehicles	VEH	Planes, trains, and automobiles

Type	Example
People	<i>Turing</i> is often considered to be the father of modern computer science.
Organization	The <i>IPCC</i> said it is likely that future tropical cyclones will become more intense.
Location	The <i>Mt. Sanitas</i> loop hike begins at the base of <i>Sunshine Canyon</i> .
Geo-Political Entity	<i>Palo Alto</i> is looking at raising the fees for parking in the University Avenue district.
Facility	Drivers were advised to consider either the <i>Tappan Zee Bridge</i> or the <i>Lincoln Tunnel</i> .
Vehicles	The updated <i>Mini Cooper</i> retains its charm and agility.

Slide courtesy Jim Martin

# Named Entity Recognition

**CHICAGO** (AP) — Citing high fuel prices, **United Airlines** said Friday it has increased fares by \$6 per round trip on flights to some cities also served by lower-cost carriers. **American Airlines**, a unit **AMR**, immediately matched the move, spokesman **Tim Wagner** said. **United**, a unit of **UAL**, said the increase took effect Thursday night and applies to most routes where it competes against discount carriers, such as **Chicago** to **Dallas** and **Atlanta** and **Denver** to **San Francisco**, **Los Angeles** and **New York**.

What are the input/output?  
What are the features?  
What types of applications?

*Slide courtesy Jim Martin*

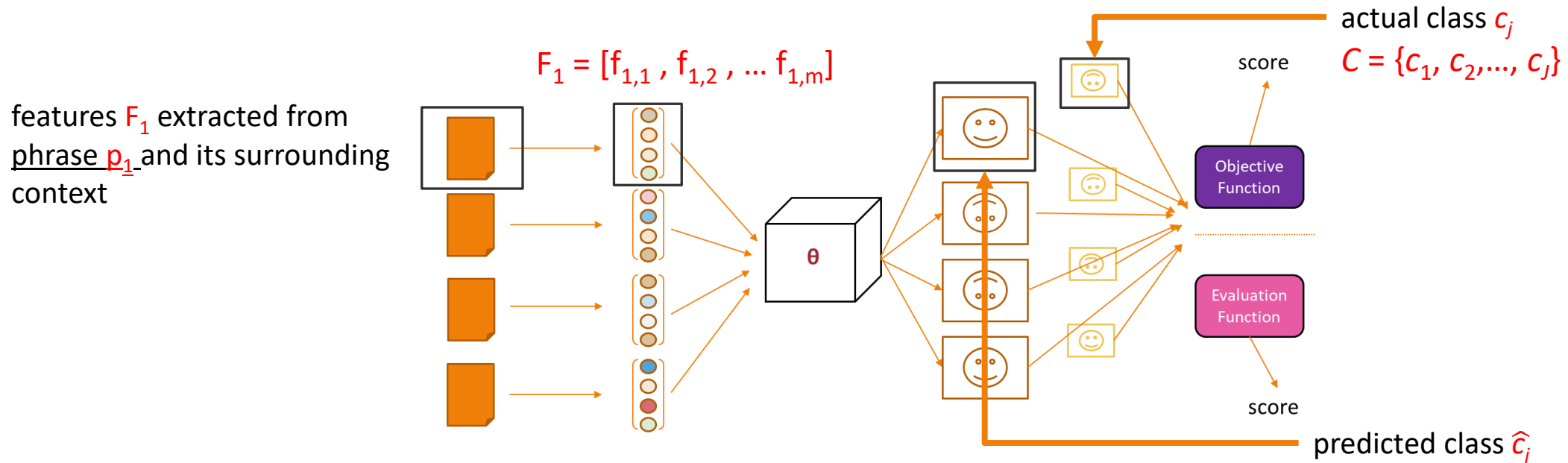
# Chunking

Named entity recognition

Information extraction

Identifying idioms

...



# Example: Information Extraction

As a task:

Filling slots in a database from sub-segments of text.

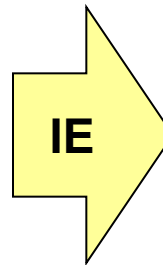
October 14, 2002, 4:00 a.m. PT

For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."

Richard Stallman, founder of the Free Software Foundation, countered saying...



<u>NAME</u>	<u>TITLE</u>	<u>ORGANIZATION</u>
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
Richard Stallman	founder	Free Soft..

Slide from Chris Brew, adapted from slide by William Cohen

# Example *applications* for IE

---

Classified ads

Restaurant reviews

Bibliographic citations

Appointment emails

Legal opinions

Papers describing clinical medical studies

Task vs  
application?

*Slide courtesy Jason Eisner, with mild edits*



# Chunking

Named entity recognition

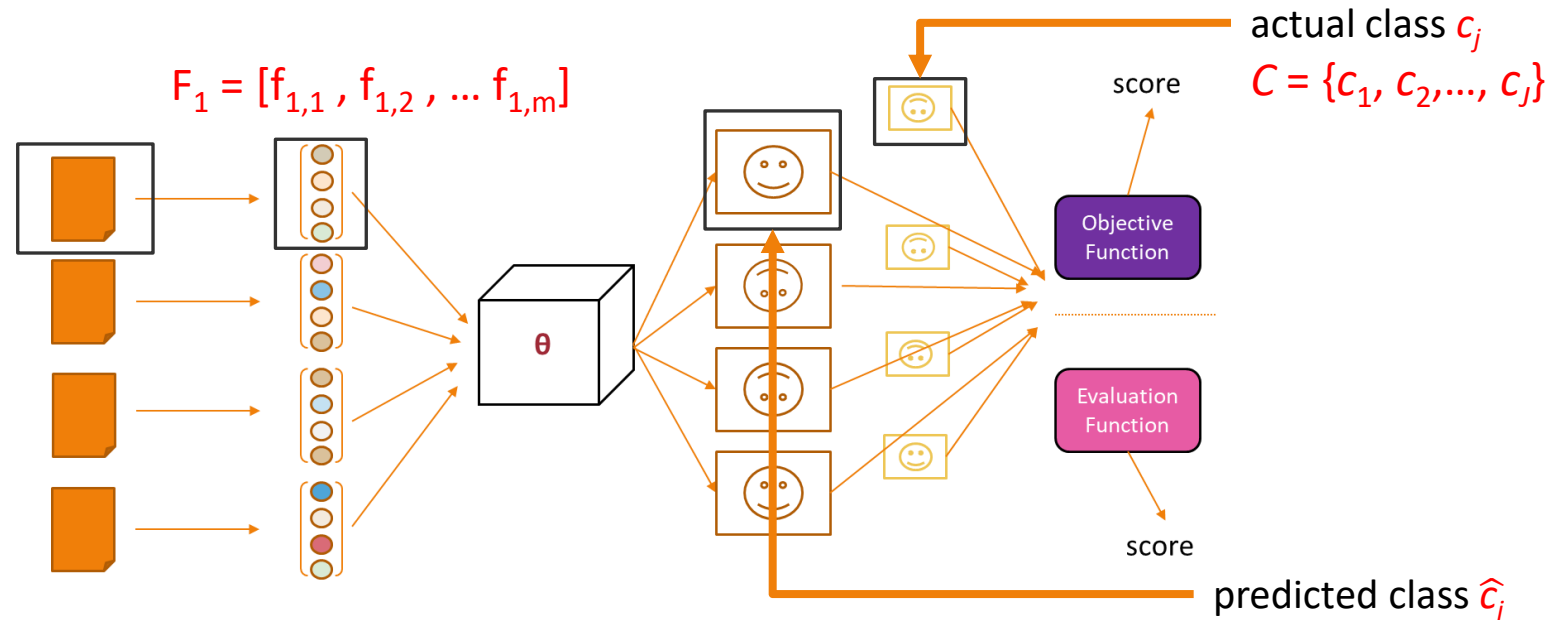
Information extraction

Identifying idioms

...

## Other examples?

features  $F_1$  extracted from  
phrase  $p_1$  and its surrounding  
context



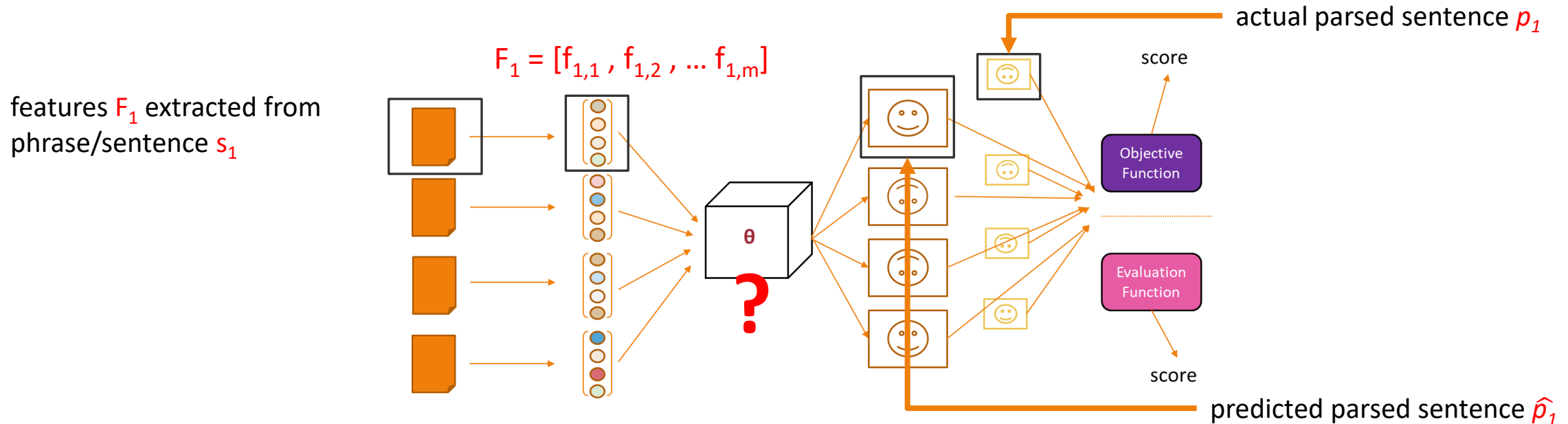
# Text Annotation Tasks ("Classification" Tasks)

---

1. Classify the entire document ("text categorization")
2. Classify word tokens individually
3. Classify word tokens in a sequence
4. Identify phrases ("chunking")
5. Syntactic annotation (syntax parsing)
6. Semantic annotation
7. Text generation

*Slide courtesy Jason Eisner, with mild edits*


# Syntax Parsing



# Context Free Grammar

---

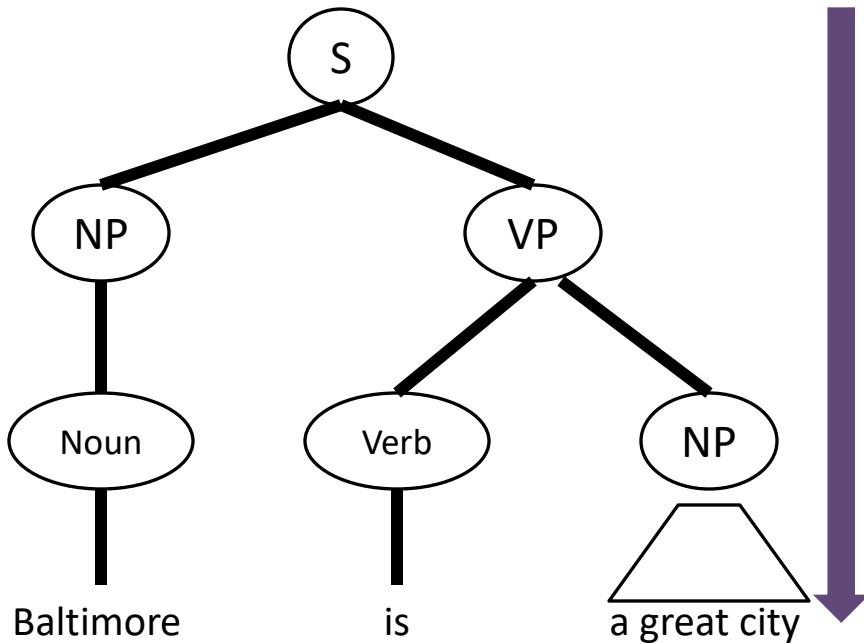
$S \rightarrow NP VP$        $PP \rightarrow P NP$   
 $NP \rightarrow Det Noun$     $AdjP \rightarrow Adj Noun$   
 $NP \rightarrow Noun$          $VP \rightarrow V NP$   
 $NP \rightarrow Det AdjP$     $Noun \rightarrow Baltimore$   
 $NP \rightarrow NP PP$         ...



Set of rewrite rules, comprised of terminals and non-terminals

# Generate from a Context Free Grammar

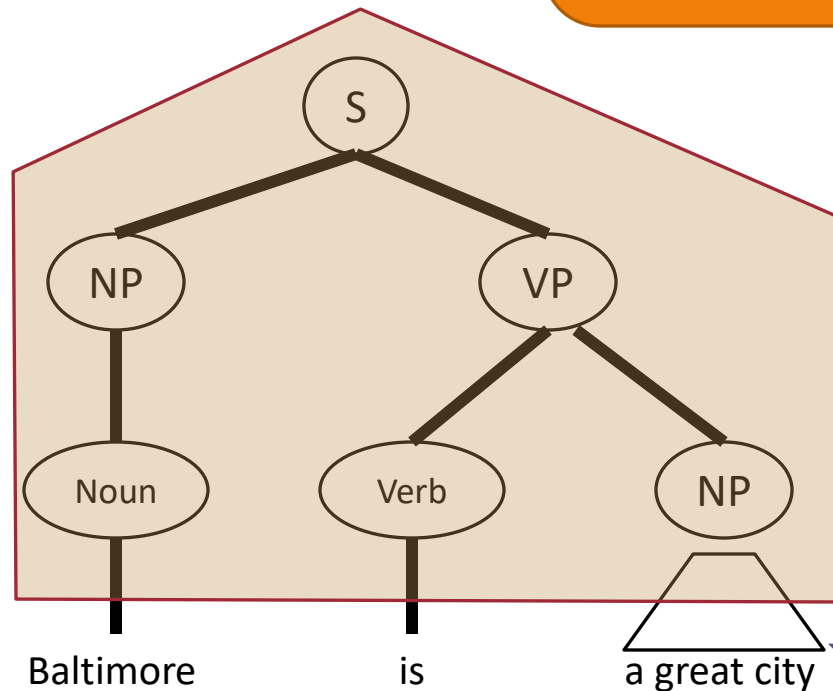
$S \rightarrow NP VP$        $PP \rightarrow P NP$   
 $NP \rightarrow Det Noun$     $AdjP \rightarrow Adj Noun$   
 $NP \rightarrow Noun$          $VP \rightarrow V NP$   
 $NP \rightarrow Det AdjP$     $Noun \rightarrow Baltimore$   
 $NP \rightarrow NP PP$         ...



Baltimore is a great city

# Assign Structure (**Parse**) with a Context Free Grammar

$S \rightarrow NP VP$        $PP \rightarrow P NP$   
 $NP \rightarrow Det Noun$     $AdjP \rightarrow Adj Noun$   
 $NP \rightarrow Noun$          $VP \rightarrow V NP$   
 $NP \rightarrow Det AdjP$     $Noun \rightarrow Baltimore$   
 $NP \rightarrow NP PP$         ...



Baltimore is a great city

$[_S [_{NP} [_{Noun} \text{Baltimore}]] [_{VP} [_{Verb} \text{is}] [_{NP} \text{a great city}]]]$

*bracket notation*

(S (NP (Noun Baltimore))  
(VP (V is)  
(NP a great city)))

*S-expression*

# Why is it useful?

---



<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

The old man the boat .





<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

The old man the boat .



<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

The rat the cat the dog chased killed ate the malt.



<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

The rat *that* the cat the dog chased killed ate the malt.



<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

The rat *that* the cat *that* the dog chased killed ate the malt.





<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

The rat *that* the cat *that* the dog chased killed ate the malt.



<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

The rat *that* the cat *that* the dog chased killed ate the malt.



<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

The rat *that* the cat *that* the dog chased killed ate the malt.



<https://www.housebeautiful.com/uk/garden/g4558287s/garden-path-ideas/>

# Garden Path Sentences

---

[The rat [the cat [the dog chased] killed] ate the malt].

Language can have recursive patterns

**Syntactic parsing** can help identify those



# Text Annotation Tasks ("Classification" Tasks)

---

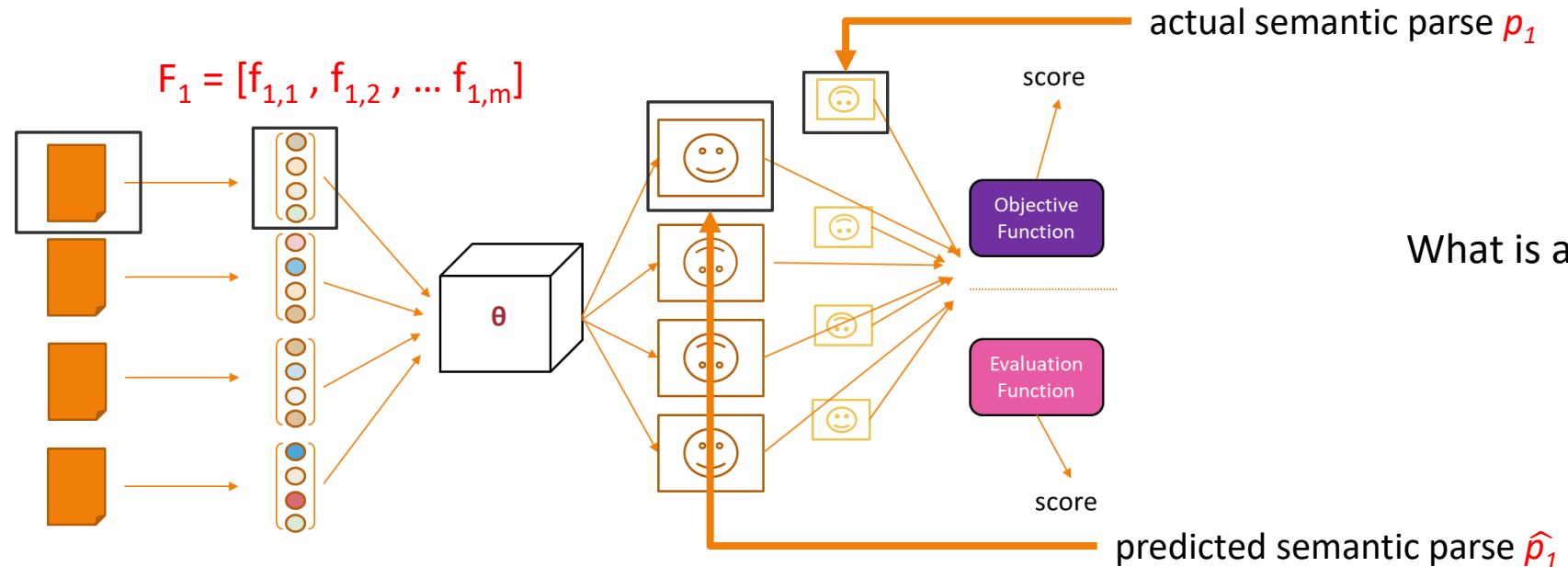
1. Classify the entire document ("text categorization")
2. Classify word tokens individually
3. Classify word tokens in a sequence
4. Identify phrases ("chunking")
5. Syntactic annotation (syntax parsing)
6. Semantic annotation
7. Text generation

*Slide courtesy Jason Eisner, with mild edits*

# Semantic Parsing

## Semantic role labeling (SRL)

features  $F_1$  extracted from phrase/sentence  $s_1$  and its surrounding context



What is a semantic parse?

# Semantic Role Labeling (SRL)

---

For each predicate (e.g., verb)

1. find its arguments (e.g., NPs)
2. determine their **semantic roles**

John drove Mary from Austin to Dallas in his Toyota Prius.

The hammer broke the window.

- **agent**: Actor of an action
- **patient**: Entity affected by the action
- **source**: Origin of the affected entity
- **destination**: Destination of the affected entity
- **instrument**: Tool used in performing action.
- **beneficiary**: Entity for whom action is performed

*Slide thanks to Ray Mooney (modified)*

# Other Current Semantic Annotation Tasks (similar to SRL)

---

PropBank – coarse-grained roles of verbs

NomBank – similar, but for nouns

FrameNet – fine-grained roles of any word

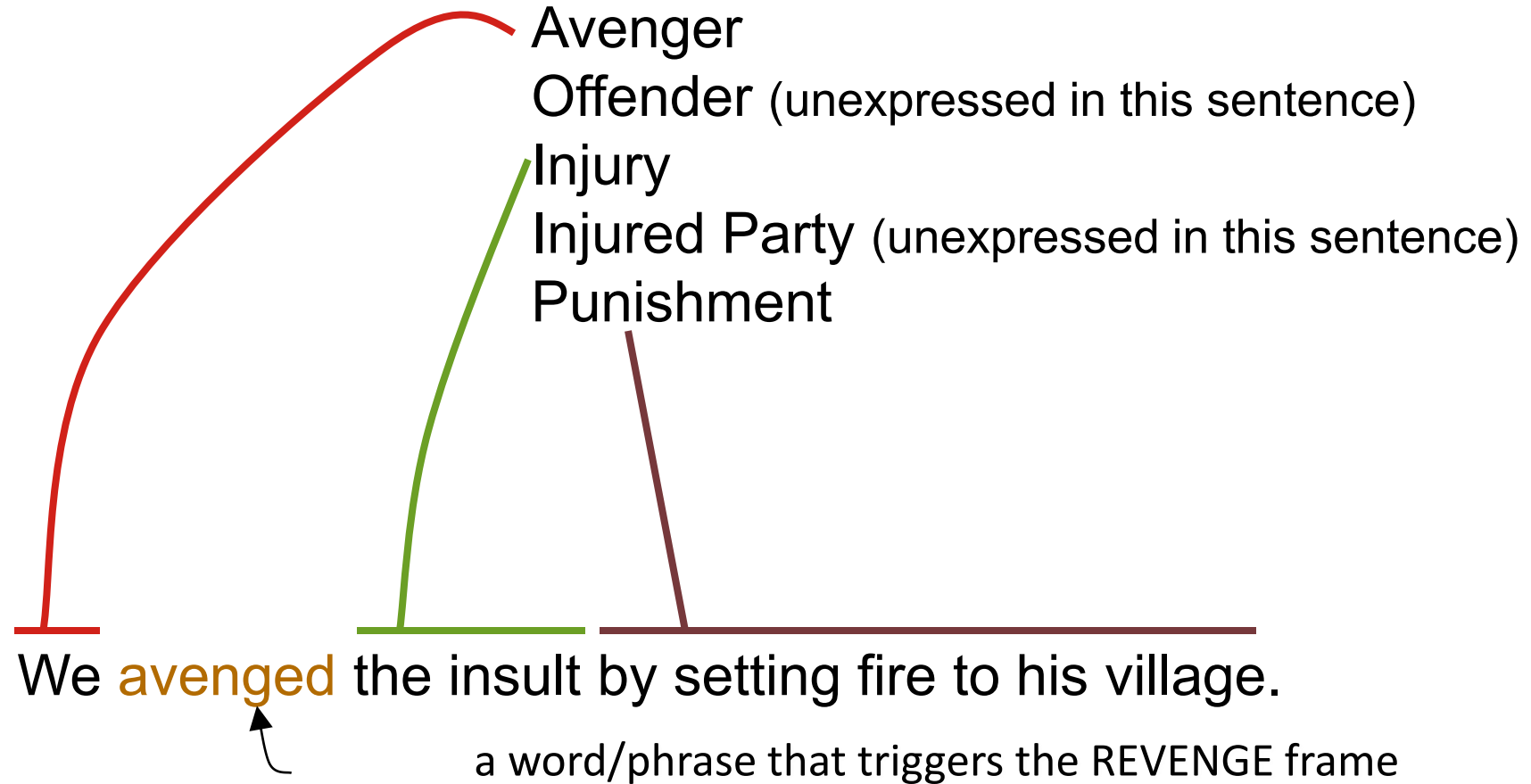
TimeBank – temporal expressions

*Slide courtesy Jason Eisner, with mild edits*

What type of applications might this have?

# FrameNet Example

REVENGE FRAME



Slide thanks to CJ Fillmore (modified)

# Text Annotation Tasks ("Classification" Tasks)

---

1. Classify the entire document ("text categorization")
2. Classify word tokens individually
3. Classify word tokens in a sequence
4. Identify phrases ("chunking")
5. Syntactic annotation (syntax parsing)
6. Semantic annotation
7. Text generation

*Slide courtesy Jason Eisner, with mild edits*

# Text Generation

Question answering (QA)

Speech recognition (ASR)

Machine translation (MT)

Summarization

Generating text from a structured representation

