

CMSC 473/673

Natural Language Processing

Instructor: Lara J. Martin (she/they)

TA: Duong Ta (he)

Slides modified from Dr. Frank Ferraro & Dr. Jason Eisner

Learning Objectives

Define featurization & other ML terminology

Define some “classification” terminology

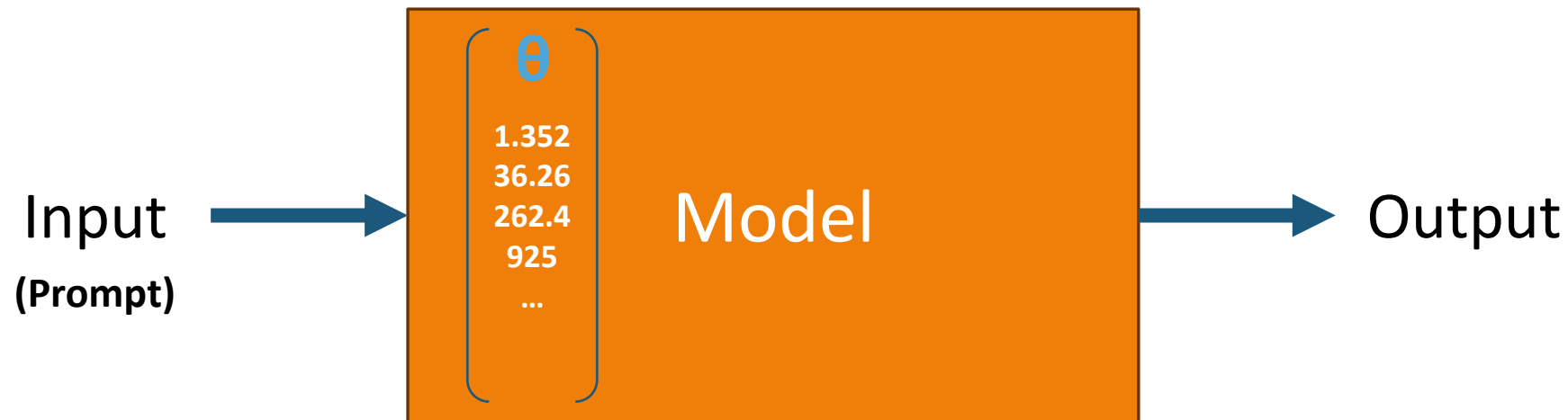
Formalize NLP Tasks at a high-level

- Document classification
- Part of speech tagging
- Syntactic parsing
- Entity id/coreference

Helpful ML Terminology

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

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



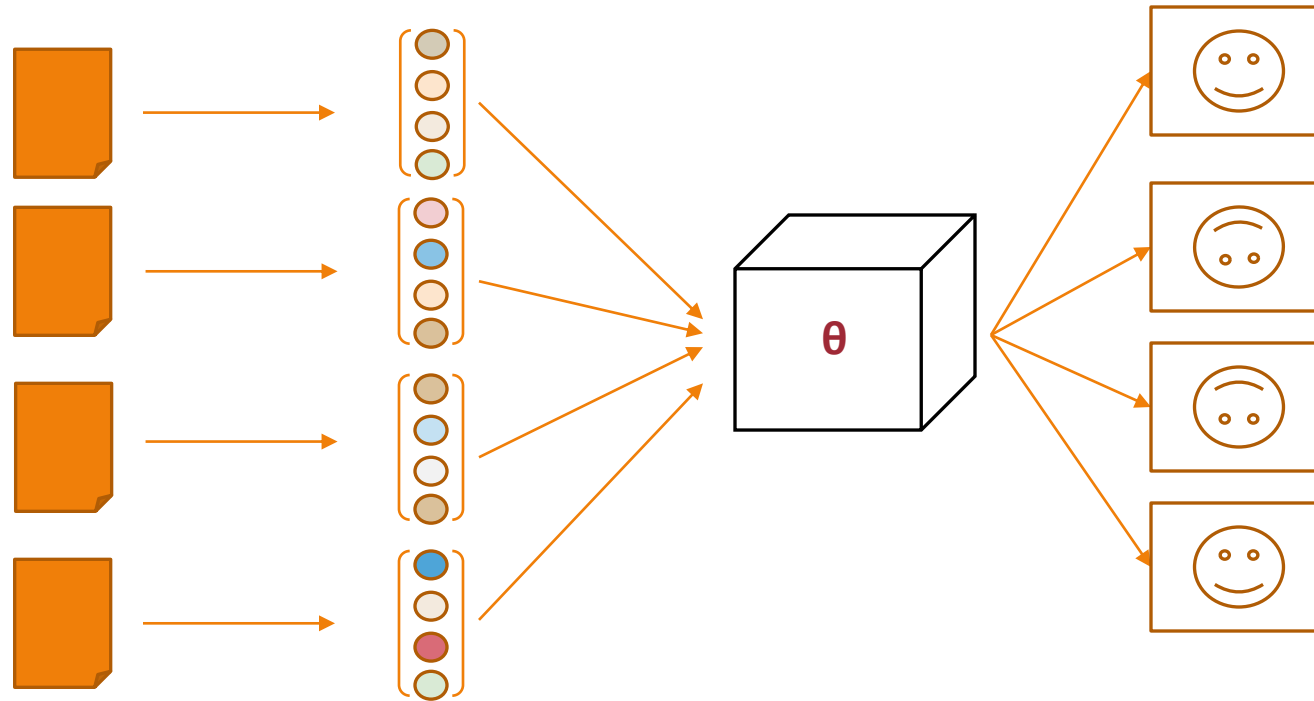
ML/NLP Framework

instances

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

ML model:

- take in featurized input
- output scores/labels
- contains weights θ



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

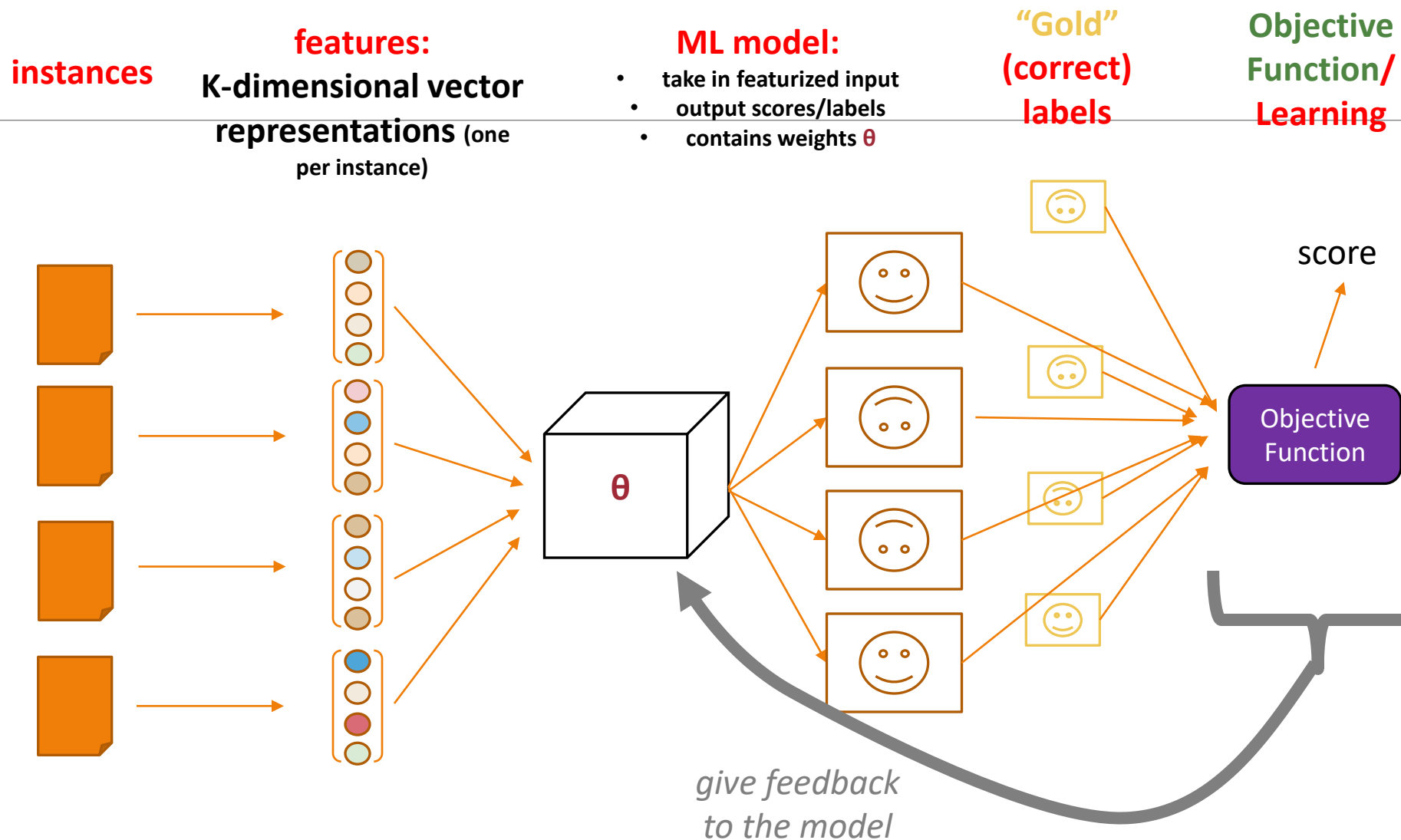
Learning:

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

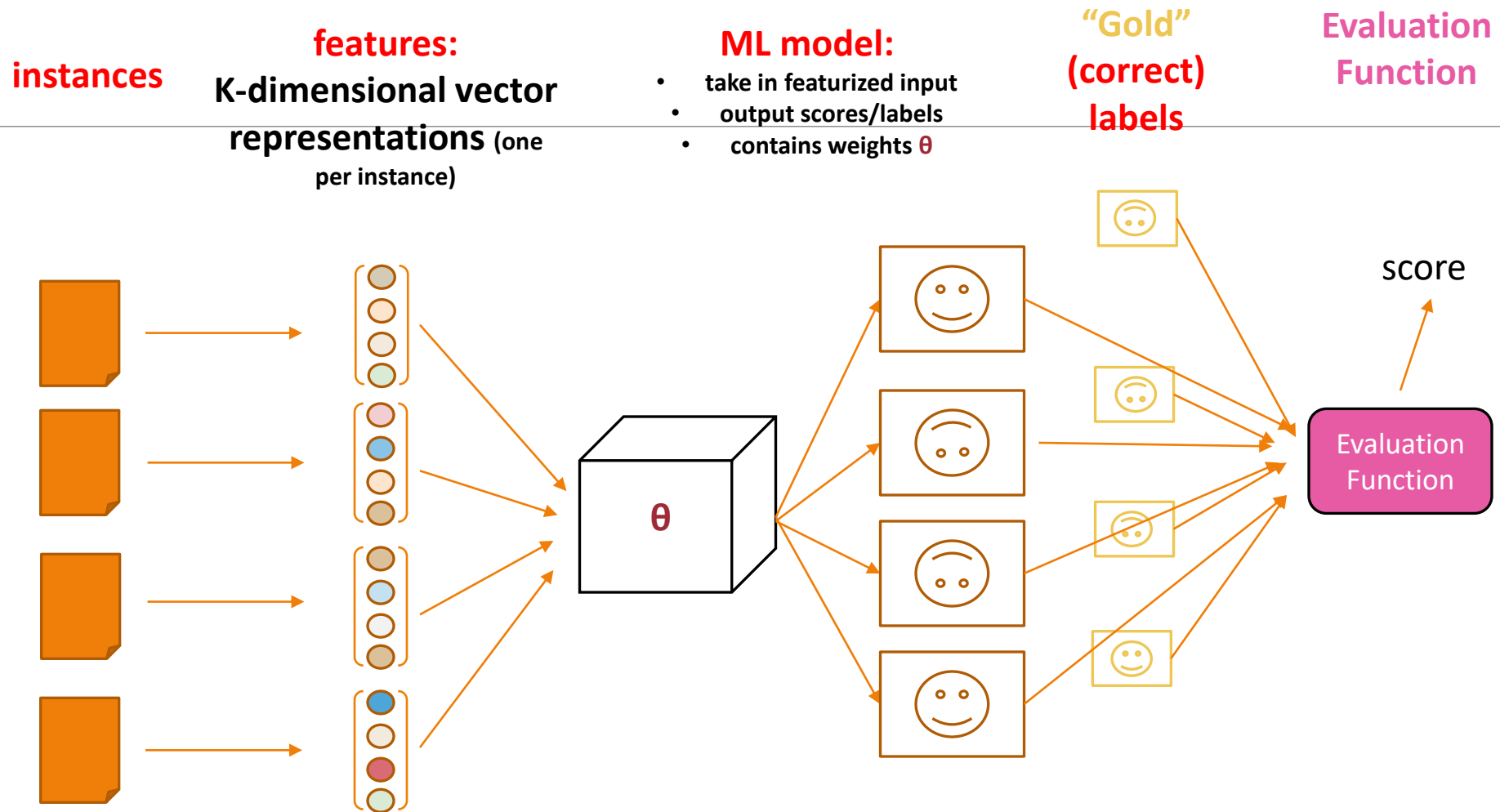
Inference / Prediction / Decoding / Classification:

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

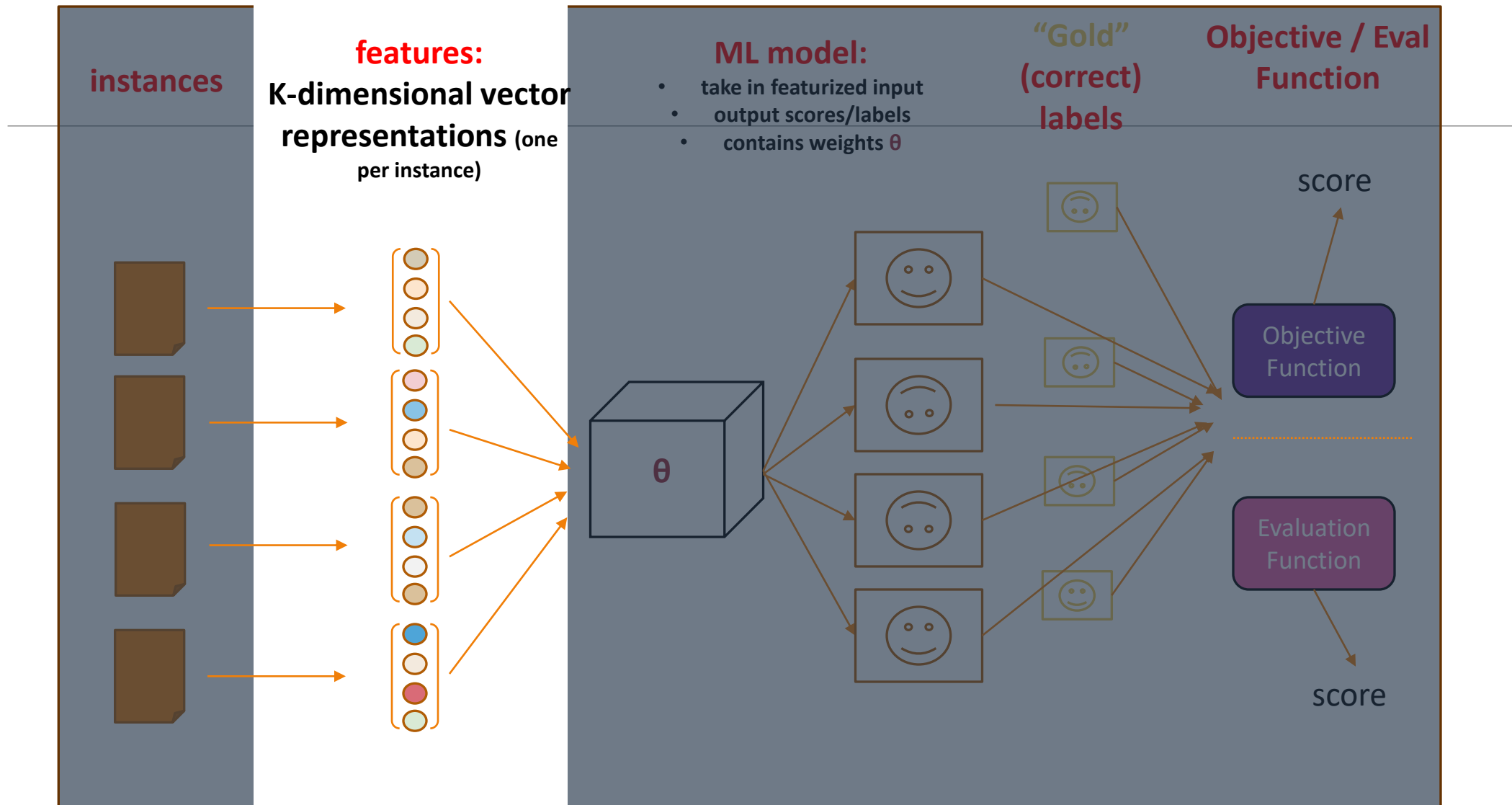
ML/NLP Framework for Learning



ML/NLP Framework for 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

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

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.

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

feature extraction

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

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

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

TECH
NOT TECH

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.

TECH
NOT TECH

$f(x)$: "bag of words"

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

w : weights

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

Three Common Types of Featurization in NLP

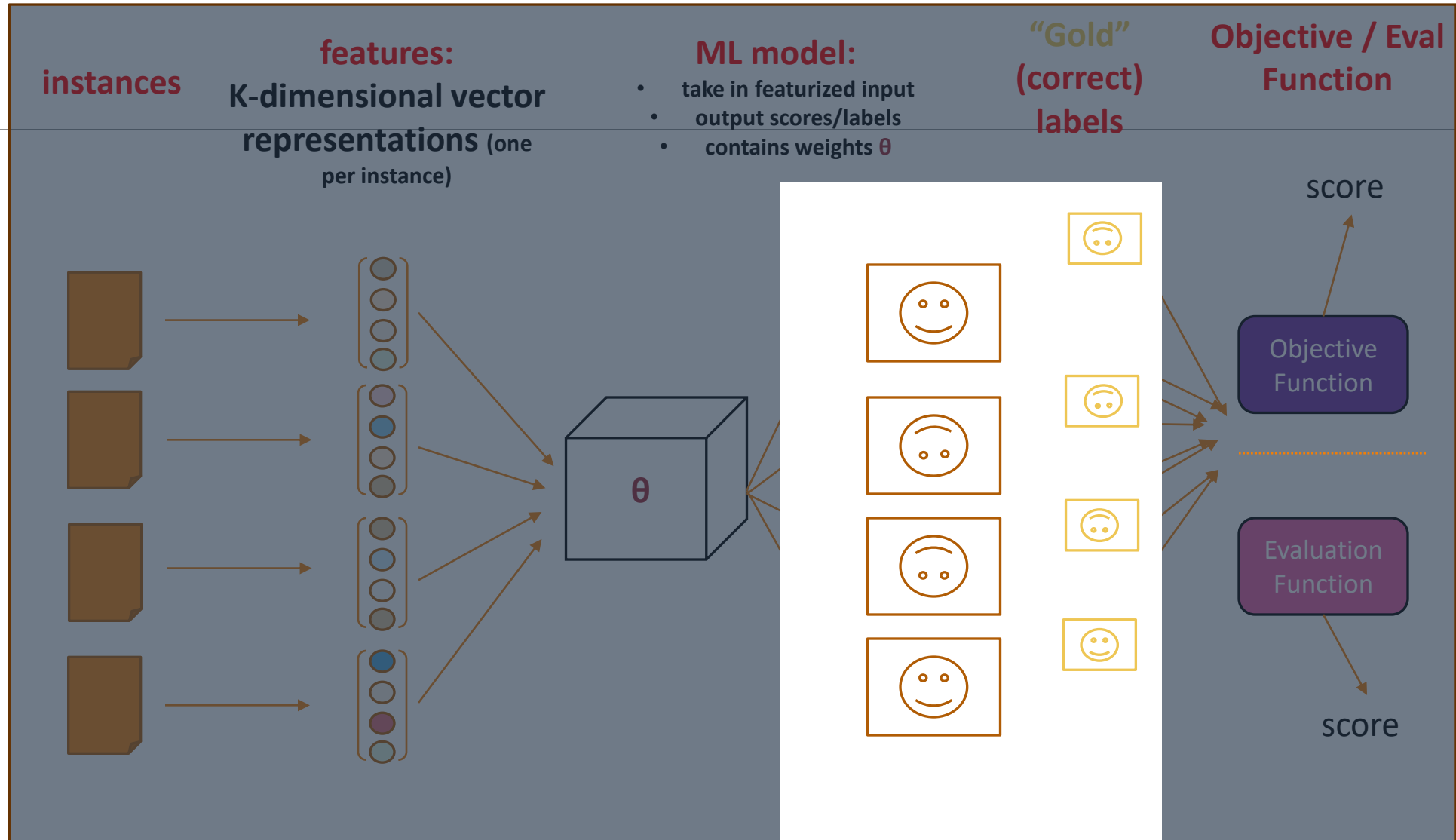
1. Bag-of-words (or bag-of-characters, bag-of-relations)
 - Identify *unique* sufficient atomic sub-parts (e.g., words in a document)
 - Define simple features over these, e.g.,
 - Binary (0 or 1) → indicating presence
 - Natural numbers → indicating number of times in a context
 - Real-valued → various 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

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

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

Text Classification

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Age/gender identification

Language Identification

Sentiment analysis

...

Text Classification

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Age/gender identification

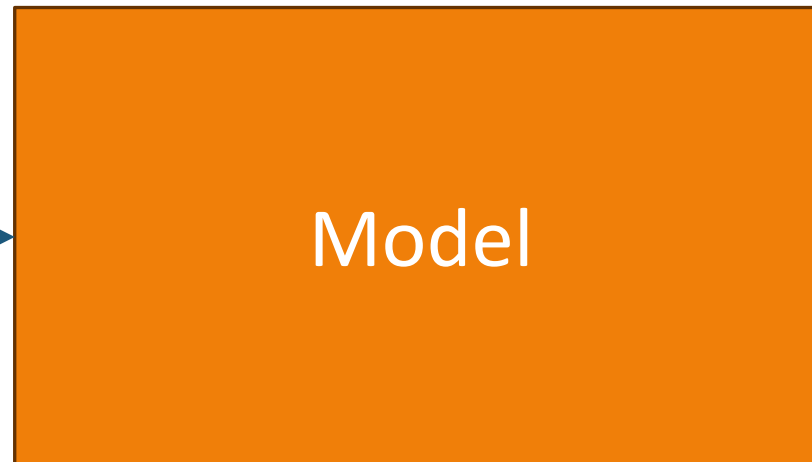
Language Identification

Sentiment analysis

...

a document

a fixed set of classes $C = \{c_1, c_2, \dots, c_j\}$



a predicted class c
from C

Text Classification: Hand-coded Rules?

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Age/gender identification

Language Identification

Sentiment analysis

...

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

Age/gender identification

Language Identification

Sentiment analysis

...

a document d

a fixed set of classes $C = \{c_1, c_2, \dots, c_j\}$

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



a learned classifier γ that maps documents to classes

Text Classification: Supervised Machine Learning

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Age/gender identification

Language Identification

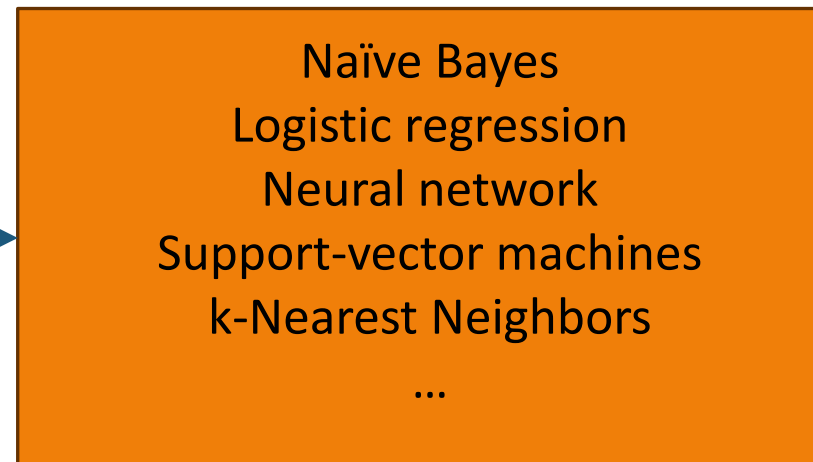
Sentiment analysis

...

a document d

a fixed set of classes $C = \{c_1, c_2, \dots, c_j\}$

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



a learned classifier γ that maps documents to classes

Text Annotation 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 |
|--------------------------|--|
| (1) Manufacturing | ... union responses to <i>plant</i> closures |
| ” ” | ... computer disk drive <i>plant</i> located in ... |
| ” ” | company manufacturing <i>plant</i> is in Orlando ... |
| (2) Living | ... 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 |
|----------------------|---|
| (1) sentencia | ... 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 ... |
| (2) frase | ... 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 ... |

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 |
|-----------------|--|
| (1) lɛd | ... it monitors the <i>lead</i> levels in drinking ... |
| ” ” | ... conference on <i>lead</i> poisoning in ... |
| ” ” | ... strontium and <i>lead</i> isotope zonation ... |
| (2) li:d | ... 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 |
|-----------------|---|
| (1) aid | ... and cut the foreign <i>aid/aide</i> budget in fiscal 1996 ... |
| ” ” | ... they offered federal <i>aid/aide</i> for flood-ravaged states ... |
| (2) aide | ... 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 Aid | Frequency as Aide |
|---------------|-------------------------|--------------------------|
| 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 |
|---|-----------|----------------------------------|-----|------|
| N-grams (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 |
| | +1P,+2P | lead , < <i>NOUN</i> > | 234 | 7 |
| Wide-context collocations | $\pm k$ W | <i>zinc</i> (in $\pm k$ words) | 235 | 0 |
| | $\pm k$ W | <i>copper</i> (in $\pm k$ words) | 130 | 0 |
| Verb-object relationships | -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

| | Frequency as Aid | Frequency as Aide |
|---------------------|----------------------------|-----------------------------|
| Word to left | | |
| 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 | led | li:d |
|---|-----------|----------------------------------|-----|------|
| N-grams (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 |
| Wide-context collocations | $\pm k$ W | <i>zinc</i> (in $\pm k$ words) | 235 | 0 |
| | $\pm k$ W | <i>copper</i> (in $\pm k$ words) | 130 | 0 |
| Verb-object relationships | -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 $\pm k$ words) | ⇒ led |
| 11.10 | lead <i>level/N</i> | ⇒ led |
| 10.66 | <i>of</i> lead <i>in</i> | ⇒ led |
| 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)

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

slide courtesy of D. Yarowsky (modified)

Text Annotation 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

Part of Speech Tagging

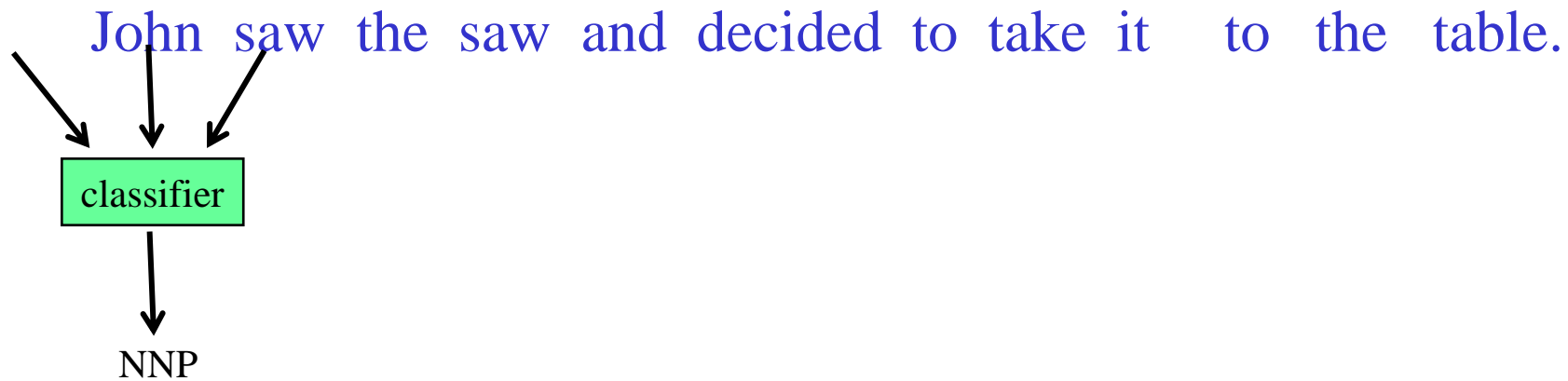
We could treat tagging as a token classification problem

- 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

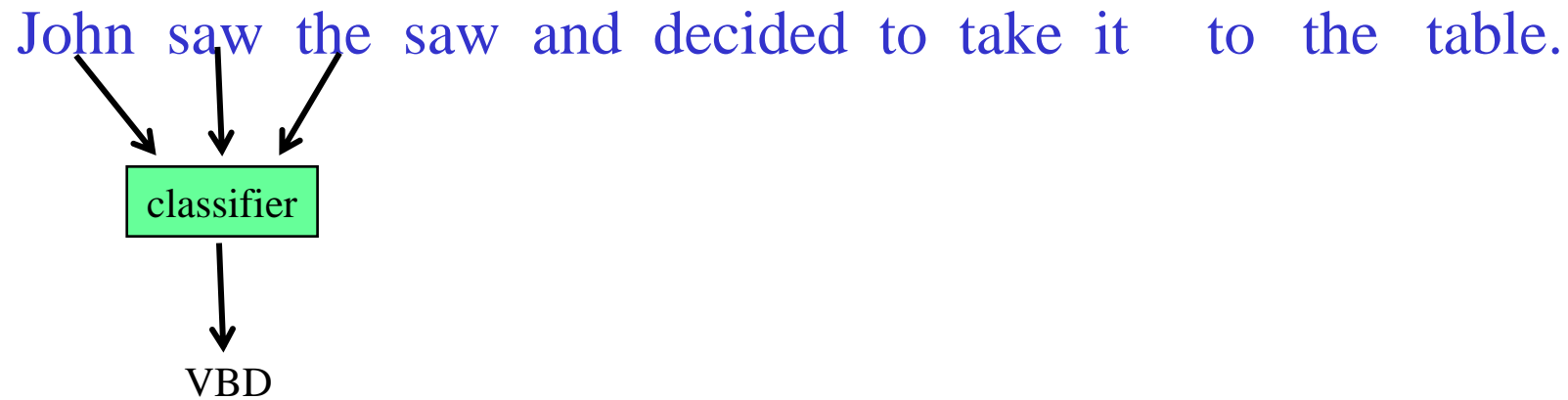
Classify each token independently but use as input features, information about the surrounding tokens (sliding window).



Slide courtesy Ray Mooney, with mild edits

Sequence Labeling as Classification

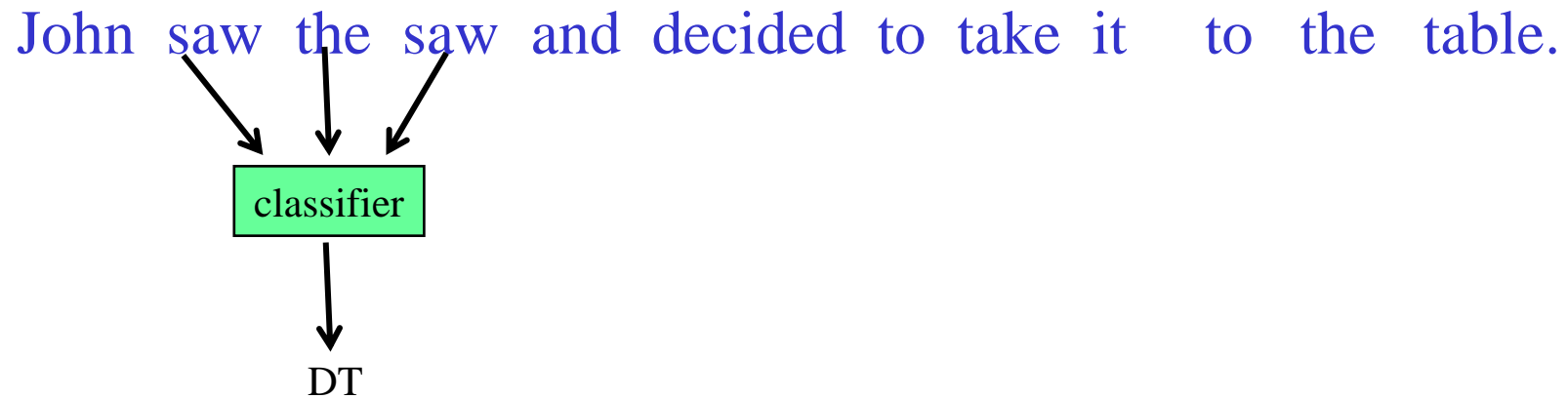
Classify each token independently but use as input features, information about the surrounding tokens (sliding window).



Slide courtesy Ray Mooney, with mild edits

Sequence Labeling as Classification

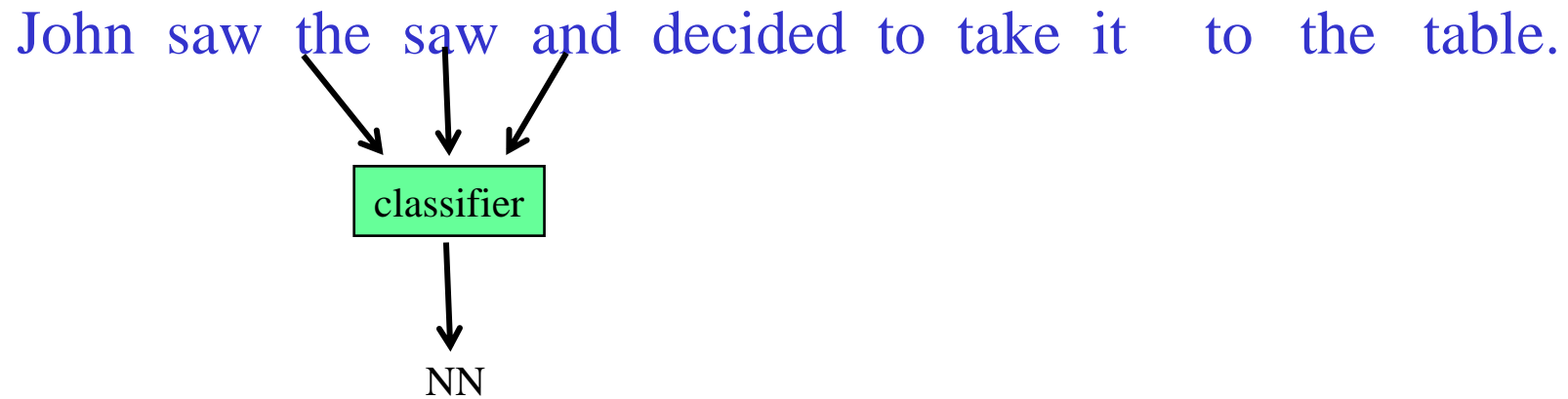
Classify each token independently but use as input features, information about the surrounding tokens (sliding window).



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).

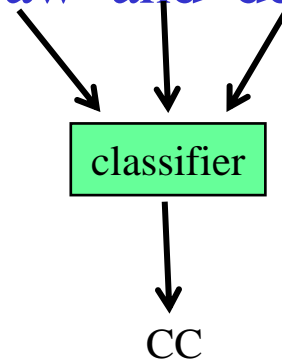


Slide courtesy Ray Mooney, with mild edits

Sequence Labeling as Classification

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

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

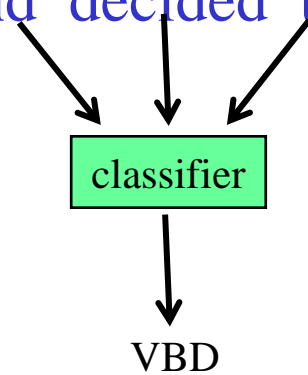


Slide courtesy Ray Mooney, with mild edits

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.



A diagram illustrating the classification process. A green rectangular box labeled "classifier" is positioned below the word "and" in the sentence above. Three black arrows point from the words "saw", "and", and "decided" to the top of the classifier box. A single black arrow points from the bottom of the classifier box to the label "VBD" below it.

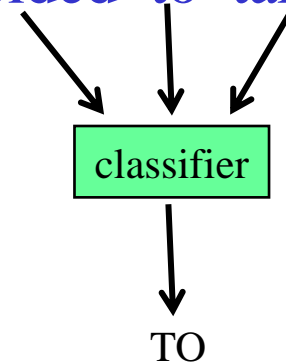
VBD

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.

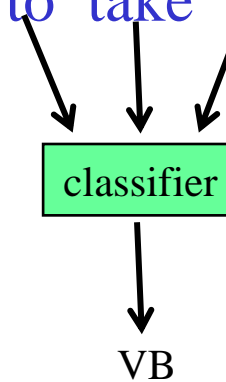


Slide courtesy Ray Mooney, with mild edits

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.

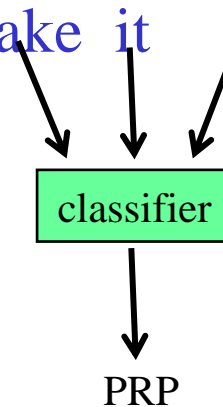


Slide courtesy Ray Mooney, with mild edits

Sequence Labeling as Classification

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

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

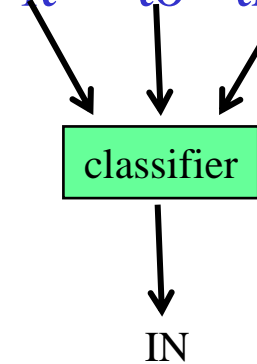


Slide courtesy Ray Mooney, with mild edits

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.

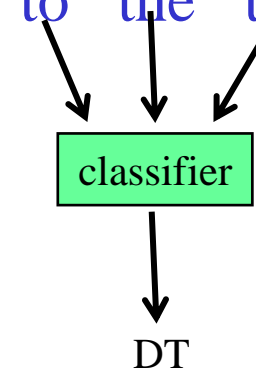


Slide courtesy Ray Mooney, with mild edits

Sequence Labeling as Classification

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

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

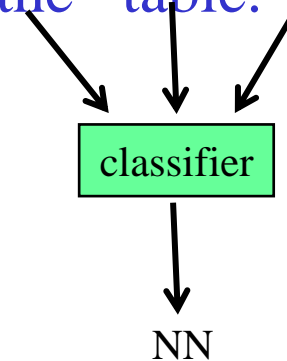


Slide courtesy Ray Mooney, with mild edits

Sequence Labeling as Classification

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

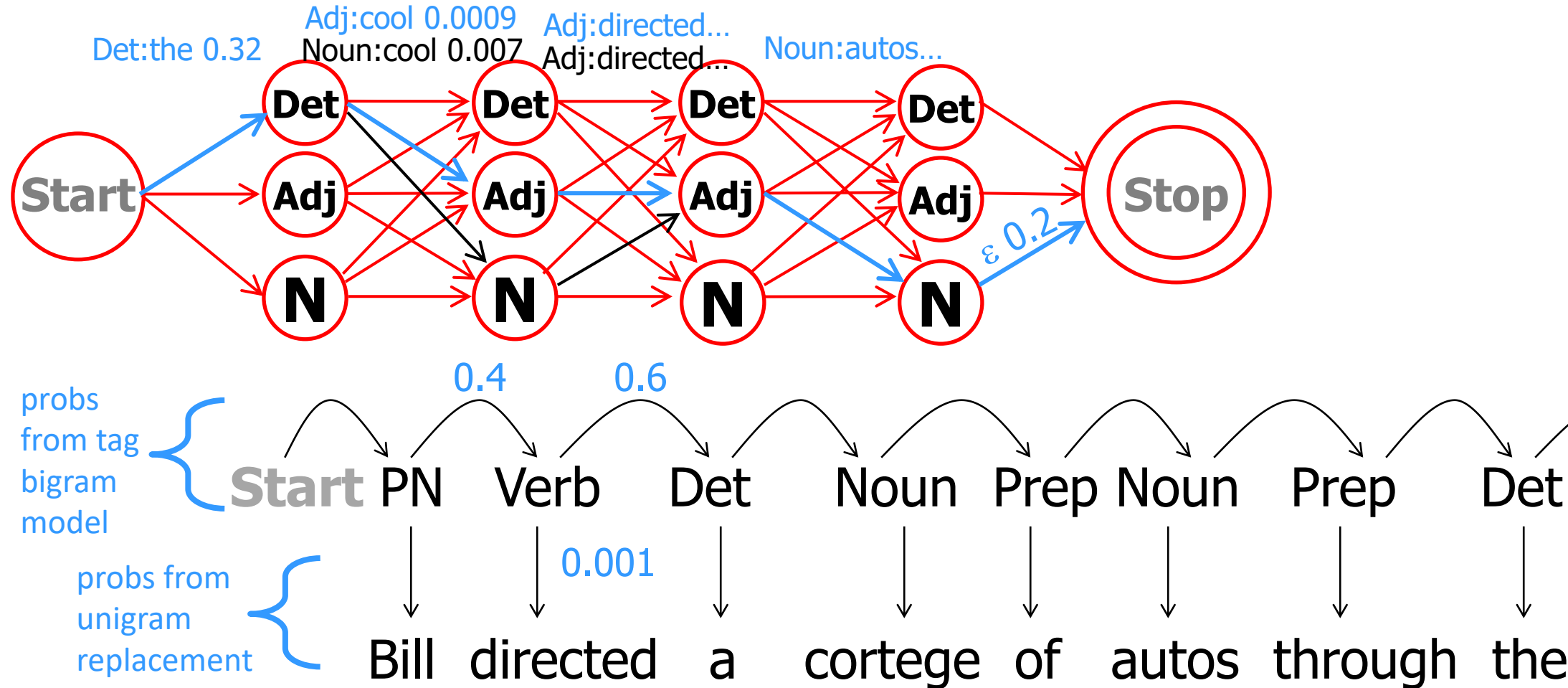
John saw the saw and decided to take it to the table.



Slide courtesy Ray Mooney, with mild edits

Part of Speech Tagging

Or we could use an HMM:



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)

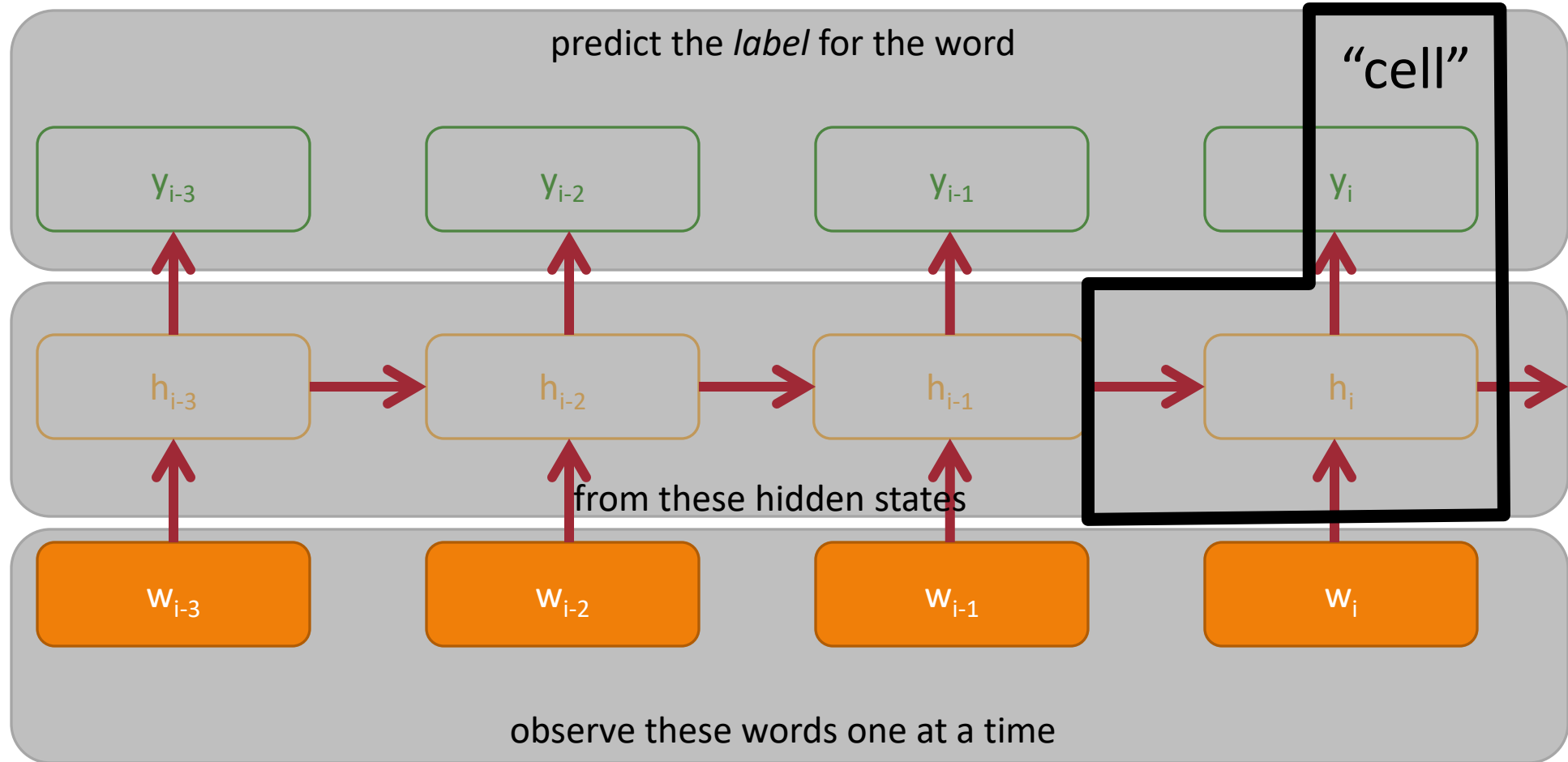
Or we could use an HMM:

- 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.

Can We Use Neural, Recurrent Methods?



Text Annotation 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**.