# NLP Tasks

Instructor: Lara J. Martin (she/they)

TA: Omkar Kulkarni (he)

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

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

# 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.
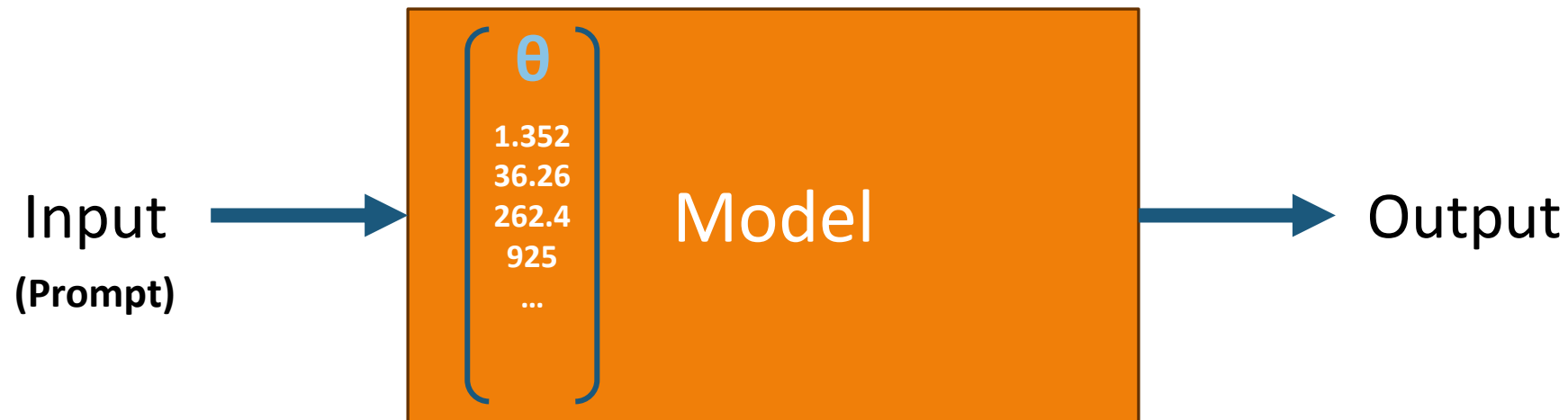
The flow of the app will look like this:
a. The user types in a keyword that they're interested in, and the app finds relevant textbooks.
b. They select the textbook and chapter they want to use.
c. The app displays a question relevant to the chapter.
d. The user answers the question.
e. The app gives 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.**
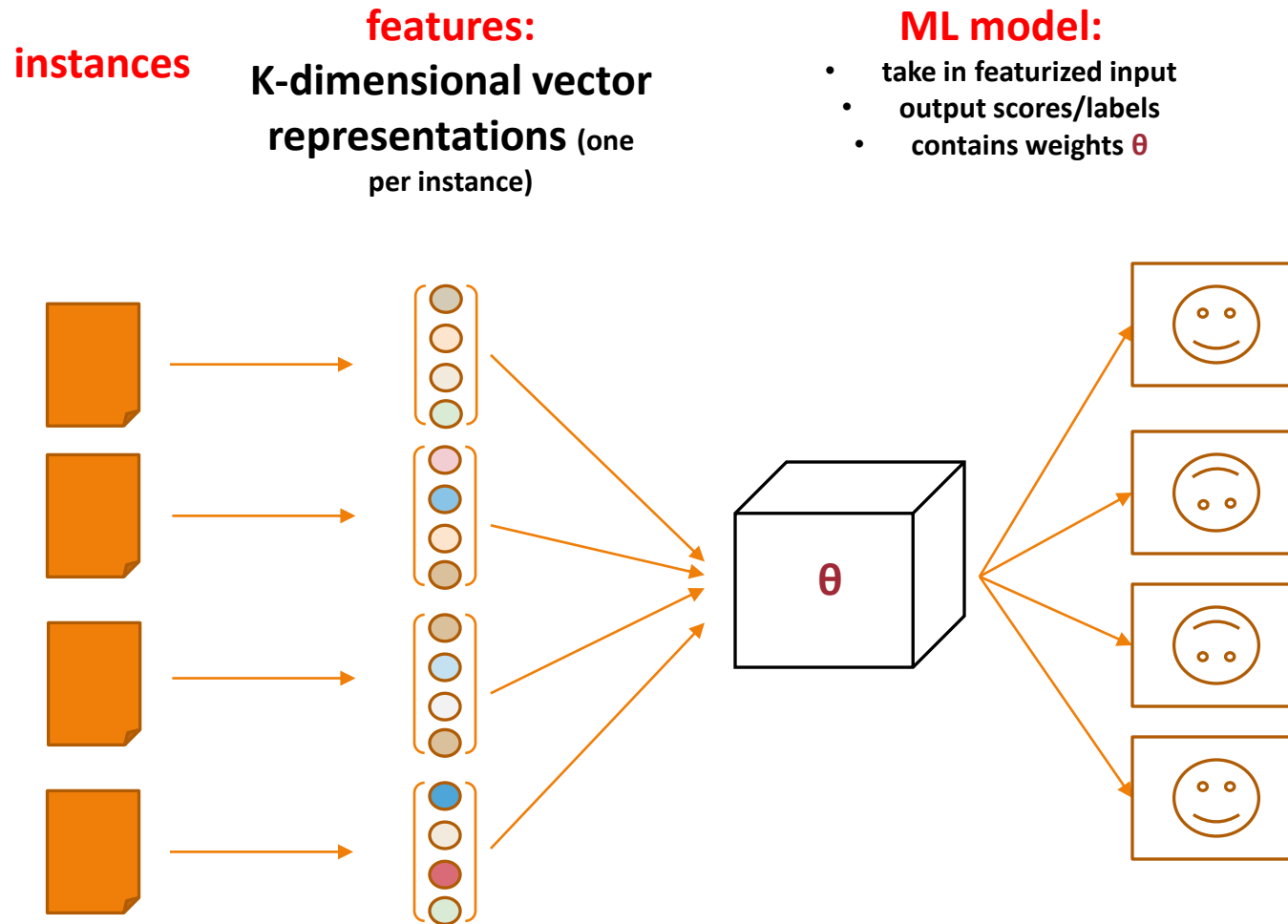
# Helpful ML Terminology

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

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

Input
**(Prompt)**

θ
1.352
36.26
262.4
925
…

Model

Output

# 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

**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 <u>Learning</u>

# ML/NLP Framework for <u>Prediction</u>

**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**

$\theta$

score

Evaluation Function

# ML/NLP Framework for Learning & Prediction

**instances**

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

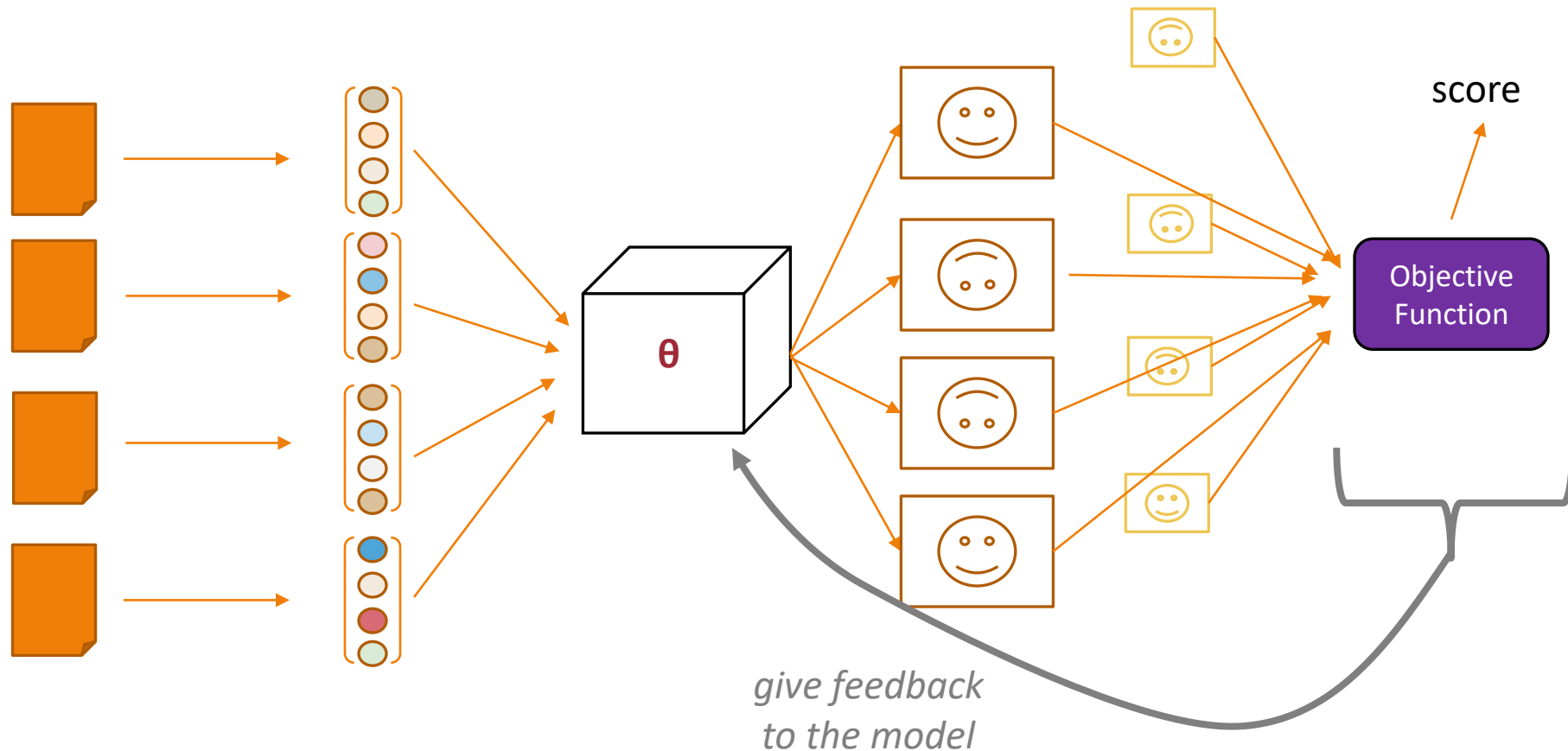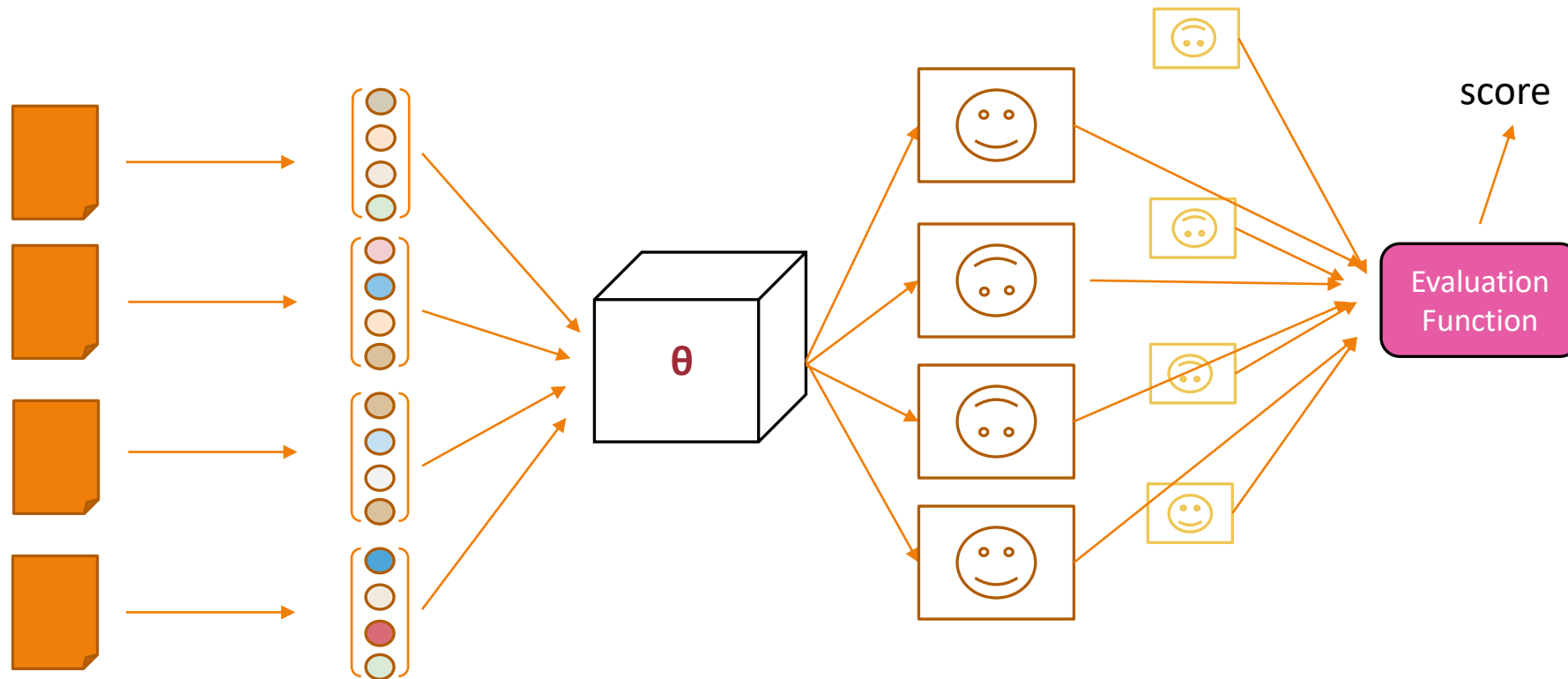**ML model:**
- take in featurized input
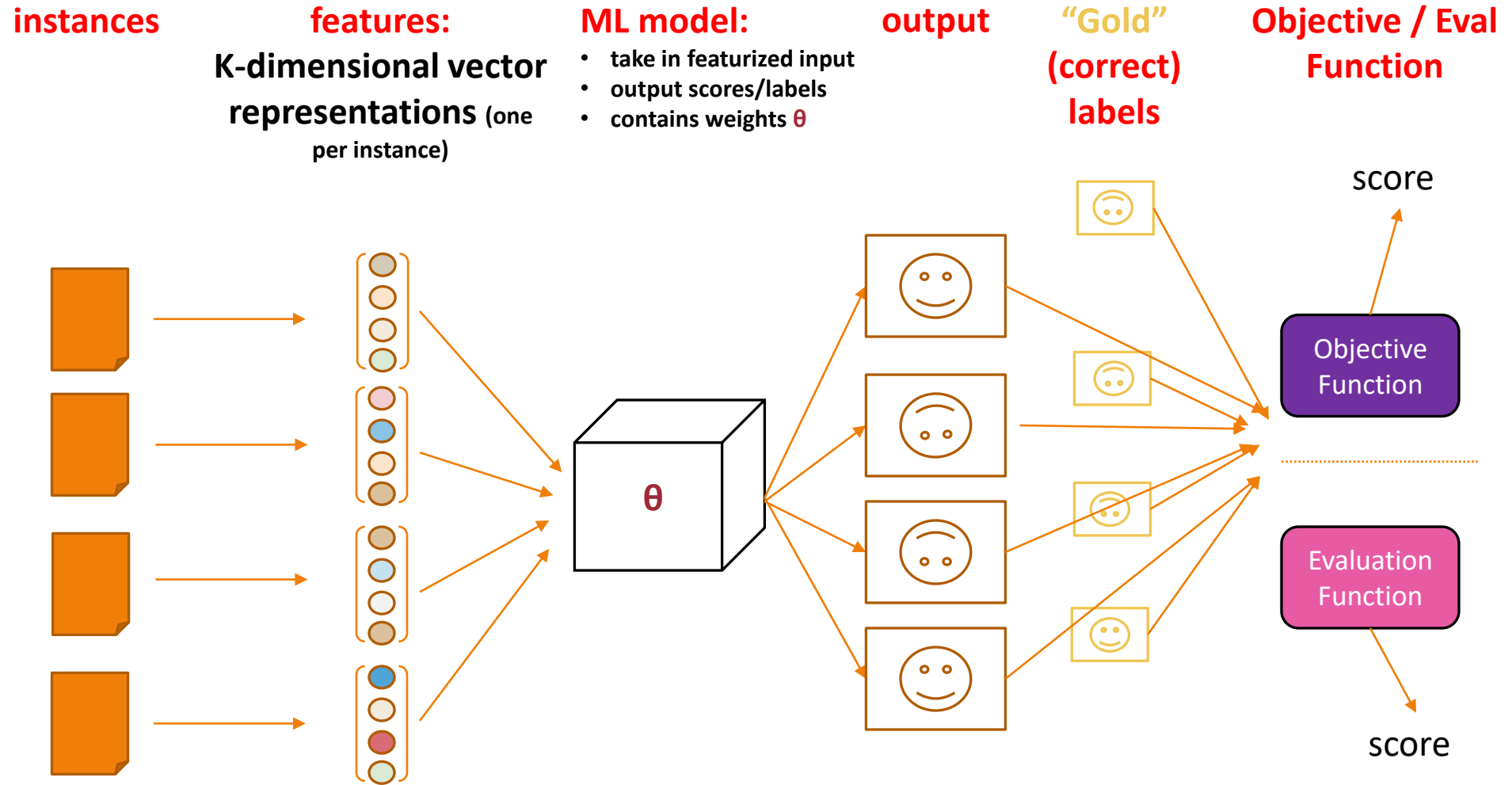- output scores/labels
- contains weights θ
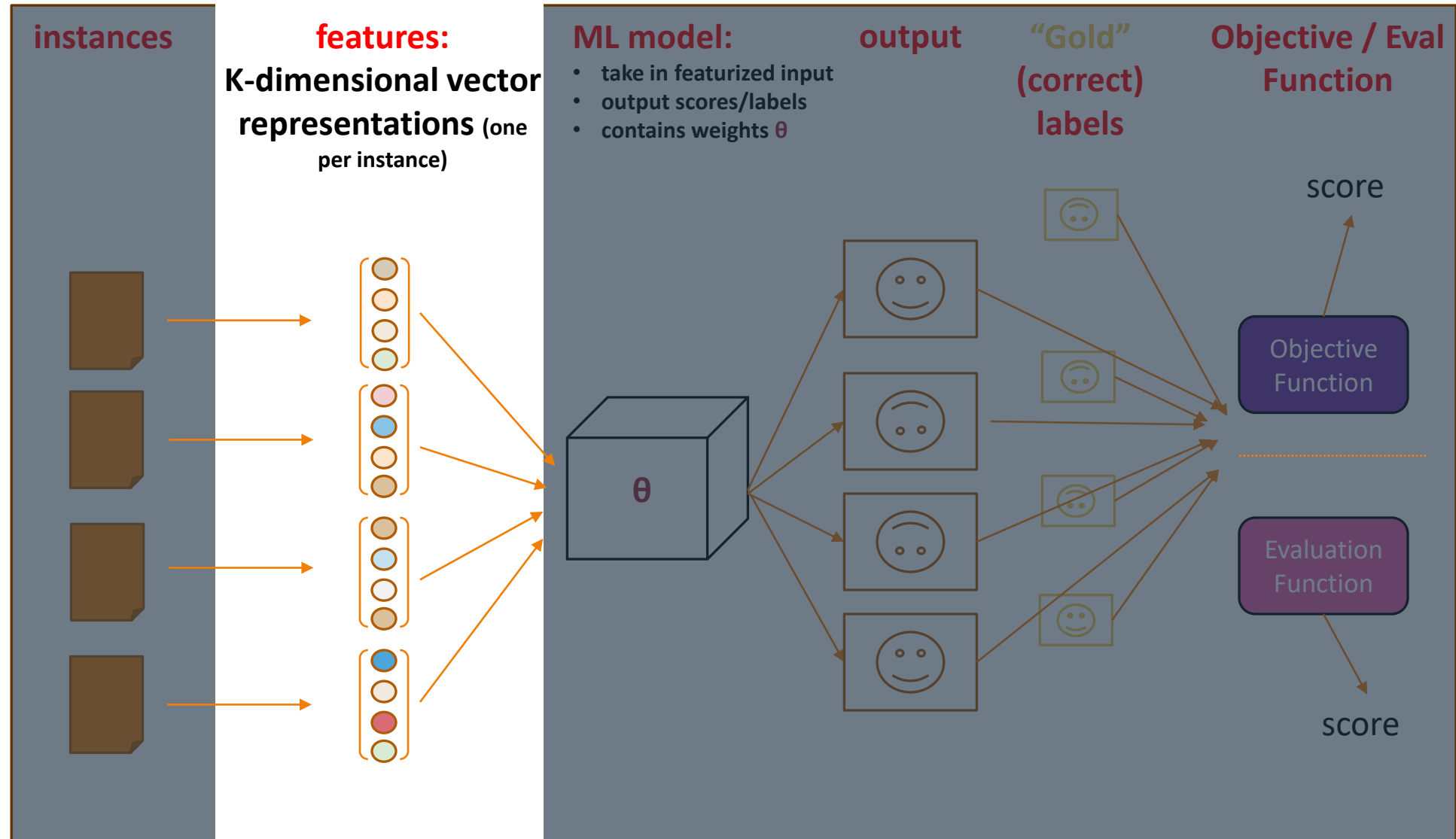
**output**

**"Gold" (correct) labels**

**Objective / Eval Function**



θ

score

Objective Function

Evaluation Function

score

# 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*

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

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

$$f(x) = \left( f_i(x) \right)_i^V$$

# 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*

| 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(**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)
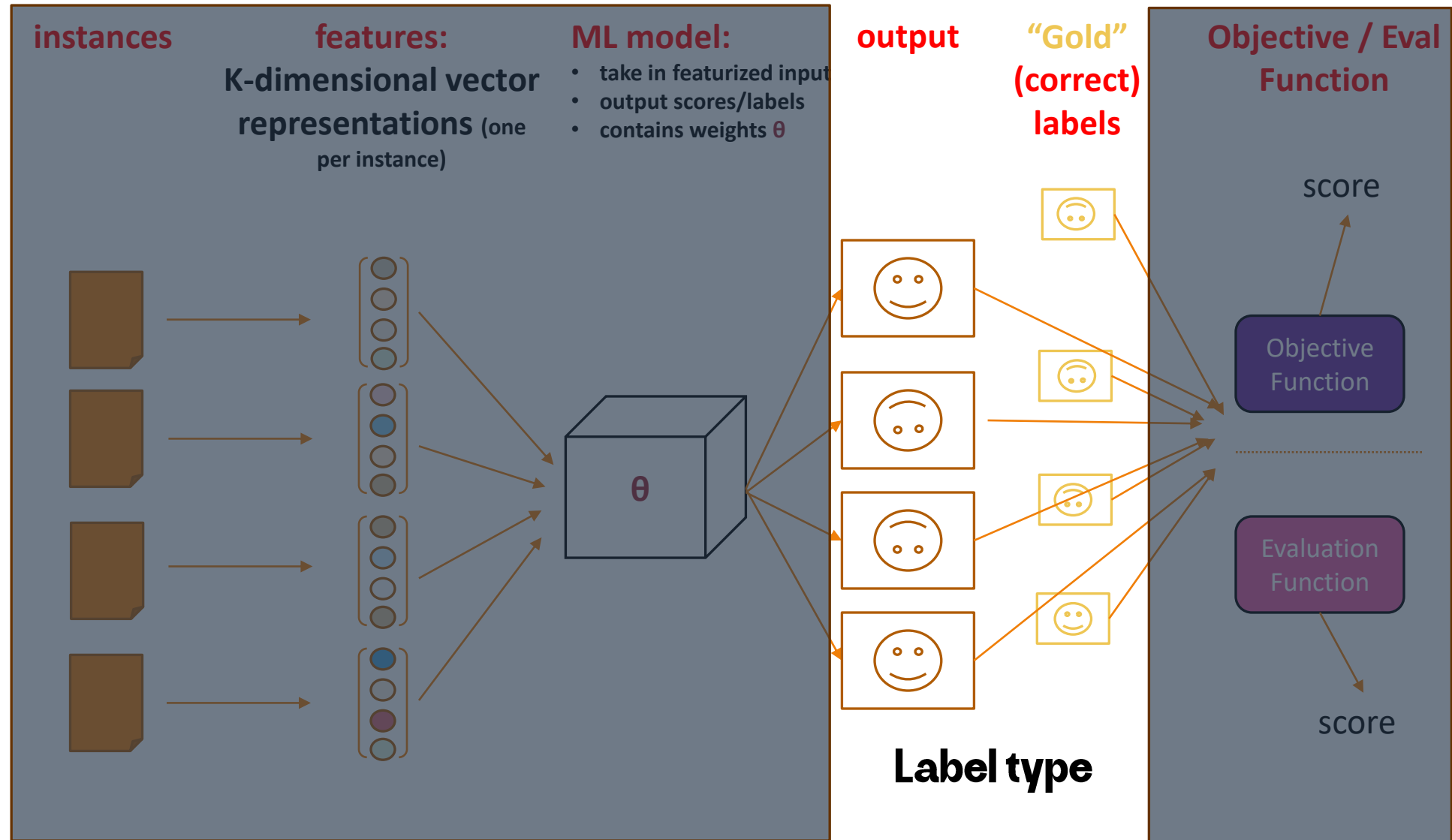
2. Linguistically-inspired features
   - Define features from words, word spans, or linguistic-based annotations extracted from the document

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



instances

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

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

output

"Gold" (correct) labels

Objective / Eval Function

score

Objective Function

Evaluation Function

score

**Label type**

# 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

# 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
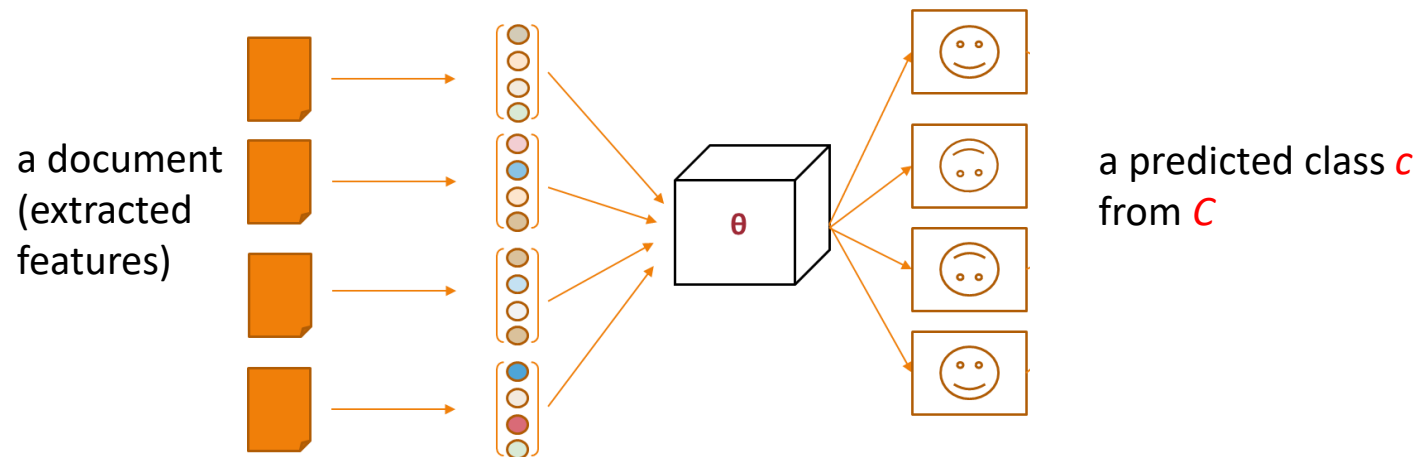
Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Language Identification

Sentiment analysis

...



a document (extracted features)

$\theta$

a predicted class $c$ from $C$

# Text Classification: Hand-coded Rules?

Assigning subject categories, topics, or genres

Spam detection

Authorship 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

Language Identification

Sentiment analysis

…

a fixed set of classes
$C = \{c_1, c_2, ..., c_J\}$

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

"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**

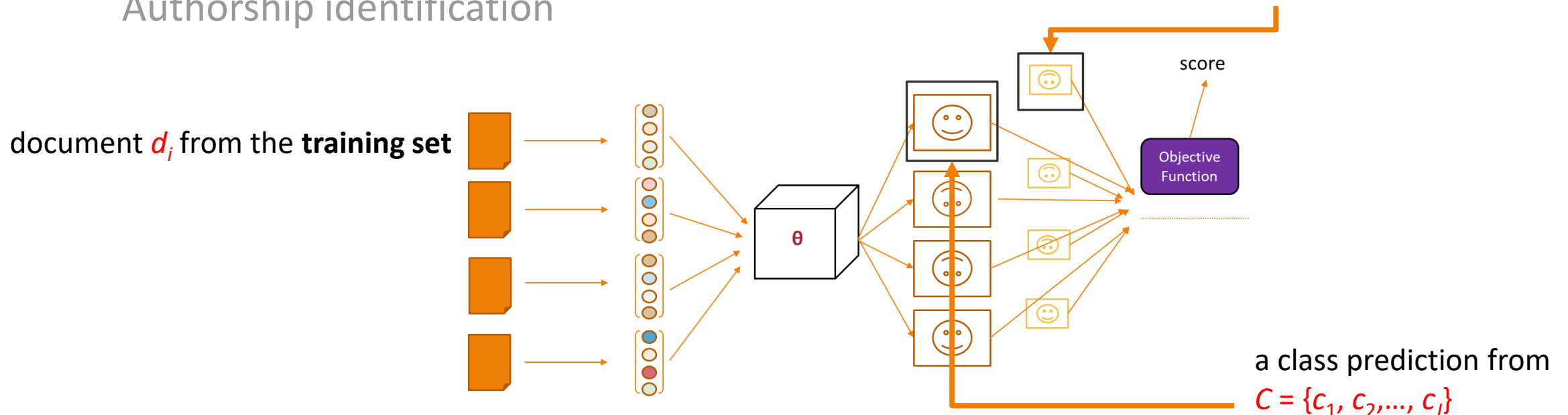Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Language Identification

Sentiment analysis

...

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

score

document $d_i$ from the **training set**

θ

Objective Function

a class prediction from

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

# Text Classification: Supervised Machine Learning - **Testing**

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Language Identification

Sentiment analysis

…



document $d_i$ from the **testing set**

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

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

Evaluation Function

score

a class prediction from $C = \{c_1, c_2, ..., c_J\}$

# Text Classification: Supervised Machine Learning – **Model examples**

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Language Identification

Sentiment analysis

...

Naïve Bayes
Logistic regression
Neural network
Support-vector machines
k-Nearest Neighbors

...

document $d_i$ from the **testing set**

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

θ

Evaluation Function

score

# 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

## 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 *plant* closures . ... |
| " " | ... computer disk drive *plant* located in ... |
| " " | company manufacturing *plant* is in Orlando ... |
| **(2) Living** | ... animal rather than *plant* tissues can be ... |
| " " | ... to strain microscopic *plant* life from the ... |
| " " | and Golgi apparatus of *plant* and animal cells |

**Test Data:**

| Sense | Context |
|---|---|
| ??? | ... vinyl chloride monomer *plant* , which is ... |
| ??? | ... molecules found in *plant* tissue from the ... |

*slide courtesy of D. Yarowsky (modified)*

# p(class | token in context)

**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 *sentence* for a young offender ... |
| "   " | ... of the minimum *sentence* of seven years in jail ... |
| "   " | ... were under the *sentence* of death at that time ... |
| **(2) frase** | ... read the second *sentence* because it is just as ... |
| "   " | ... The next *sentence* is a very important ... |
| "   " | ... It is the second *sentence* which I think is at ... |

**Test Data:**

| Translation | Context |
|---|---|
| ??? | ... cannot criticize a *sentence* handed down by ... |
| ??? | ... listen to this *sentence* uttered by a former ... |

# p(class | token in context)

## 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)

# p(class | token in context)

## Accent Restoration in Spanish & French

**Training Data:**

| Pattern | Context |
|---|---|
| **(1) côté** | ... du laisser de *cote* faute de temps ... |
| " " | ... appeler l' autre *cote* de l' atlantique ... |
| " " | ... passe de notre *cote* de la frontiere ... |
| **(2) côte** | ... vivre sur notre *cote* ouest toujours ... |
| " " | ... creer sur la *cote* du labrador des ... |
| " " | travaillaient cote a *cote* , ils avaient ... |

**Test Data:**

| Pattern | Context |
|---|---|
| ??? | ... passe de notre *cote* de la frontiere ... |
| ??? | ... creer sur la *cote* du labrador des ... |

*slide courtesy of D. Yarowsky (modified)*

# Text-to-Speech Synthesis

**Problem:**

... slightly elevated *lead* levels ...
- ⇒ l$\epsilon$d (as in *lead mine*)    or
- ⇒ li:d (as in *lead role*)

**Training Data:**

| Pronunciation | Context |
|---|---|
| **(1) l$\epsilon$d** | ... it monitors the *lead* levels in drinking ... |
| " " | ... conference on *lead* poisoning in ... |
| " " | ... strontium and *lead* isotope zonation ... |
| **(2) li:d** | ... maintained their *lead* Thursday over ... |
| " " | ... to Boston and *lead* singer for Purple ... |
| " " | ... Bush a 17-point *lead* in Texas , only 3 ... |

**Test Data:**

| Pronunciation | Context |
|---|---|
| ??? | ... median blood *lead* concentration was .. |
| ??? | ... his double-digit *lead* nationwide . The ... |

*slide courtesy of D. Yarowsky (modified)*

# p(class | token in context)

## Spelling Correction

**Problem:**

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

$\Rightarrow$ *aid*   or

$\Rightarrow$ *aide*

**Training Data:**

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

**Test Data:**

| Spelling | Context |
|---|---|
| ??? | ... said the longtime *aid/aide* to the Mayor of St. ... |
| ??? | ... will squander the *aid/aide* 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 ≥ 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 | lɛd | li:d |
|---|---|---|---|---|
| **N-grams** | +1 L | lead *level/N* | 219 | 0 |
|  | -1 W | *narrow* lead | 0 | 70 |
| (word, | +1 W | lead *in* | 207 | 898 |
| lemma, | -1W,+1W | *of* lead *in* | 162 | 0 |
| part-of-speech) | -1W,+1W | *the* lead *in* | 0 | 301 |
|  | +1P,+2P | lead , *<NOUN>* | 234 | 7 |
| **Wide-context** | ±k W | *zinc* (in ±k words) | 235 | 0 |
| **collocations** | ±k W | *copper* (in ±k words) | 130 | 0 |
| **Verb-object** | -V L | *follow/V* + lead | 0 | 527 |
| **relationships** | -V L | *take/V* + 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 |

# An assortment of possible cues ...

| | Position | Collocation | lɛd | li:d |
|---|---|---|---|---|
| **N-grams** | +1 L | lead *level/N* | 219 | 0 |
| | -1 W | *narrow* lead | 0 | 70 |
| (word, | +1 W | lead *in* | 207 | 898 |
| lemma, | -1W,+1W | *of* lead *in* | 162 | 0 |
| part-of-speech) | -1W,+1W | *the* lead *in* | 0 | 301 |
| | +1P,+2P | lead , *<NOUN>* | 234 | 7 |
| **Wide-context** | ±k W | *zinc* (in ±k words) | 235 | 0 |
| **collocations** | ±k W | *copper* (in ±k words) | 130 | 0 |
| **Verb-object** | -V L | *follow/V* + lead | 0 | 527 |
| **relationships** | -V L | *take/V* + 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 | *follow/V* + lead | ⇒ li:d |
|---|---|---|
| 11.20 | *zinc* (in ±k words) | ⇒ lɛd |
| 11.10 | lead *level/N* | ⇒ lɛd |
| 10.66 | *of* lead *in* | ⇒ lɛd |
| 10.59 | *the* lead *in* | ⇒ li:d |
| 10.51 | lead *role* | ⇒ li:d |

# 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 | *follow/V* + lead | ⇒ li:d |
| 11.20 | *zinc* (in ±$k$ words) | ⇒ lɛd |
| 11.10 | lead *level/N* | ⇒ lɛd |
| 10.66 | *of* lead *in* | ⇒ lɛd |
| 10.59 | *the* lead *in* | ⇒ li:d |
| 10.51 | lead *role* | ⇒ li:d |
| 10.35 | *copper* (in ±$k$ words) | ⇒ lɛd |
| 10.28 | lead *time* | ⇒ li:d |
| 10.24 | lead *levels* | ⇒ lɛd |
| 10.16 | lead *poisoning* | ⇒ lɛd |
| 8.55 | *big* lead | ⇒ li:d |
| 8.49 | *narrow* lead | ⇒ li:d |
| 7.76 | *take/V* + lead | ⇒ li:d |
| 5.99 | lead , *NOUN* | ⇒ lɛd |
| 1.15 | lead *in* | ⇒ li:d |

○ ○ ○

*slide courtesy of D. Yarowsky (modified)*

# Text-to-Speech Synthesis

**Problem:**

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

$\Rightarrow$ *l$\epsilon$d* (as in *lead mine*)   or

$\Rightarrow$ *li:d* (as in *lead role*)

**Training Data:**

| Pronunciation | Context |
|---|---|
| **(1) l$\epsilon$d** | ... it monitors the *lead* levels in drinking ... |
| " " | ... conference on *lead* poisoning in ... |
| " " | ... strontium and *lead* isotope zonation ... |
| **(2) li:d** | ... maintained their *lead* Thursday over ... |
| " " | ... to Boston and *lead* singer for Purple ... |
| " " | ... Bush a 17-point *lead* in Texas , only 3 ... |

**Test Data:**

| Pronunciation | Context |
|---|---|
| ??? | ... median blood *lead* concentration was .. |
| ??? | ... his double-digit *lead* nationwide . The ... |

# Token Classification

Word pronunciation

Accent restoration

Word sense disambiguation (WSD)
within or across languages

...

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

$F_1 = [f_{1,1} , f_{1,2} , ... f_{1,m}]$

$\theta$

actual class $c_j$

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

score

Objective Function

Evaluation Function

score

prediction

# 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

# Example: Part of Speech Tagging

We could treat tagging as a token classification problem
- ◦ Tag each word independently given features of context
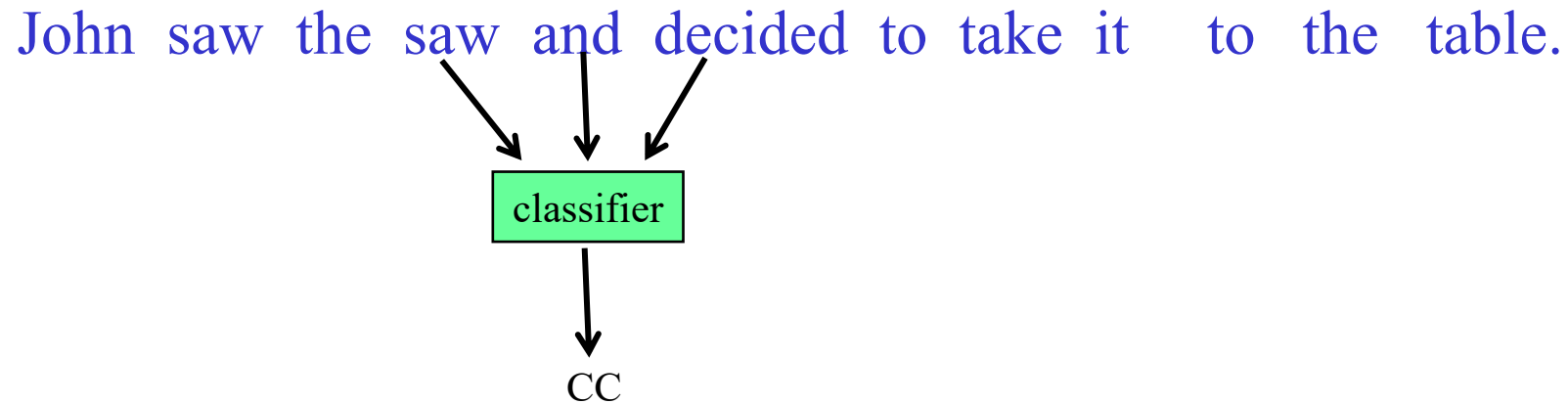- ◦ And features of the word's spelling (suffixes, capitalization)

# Sequence Labeling as Classification

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).
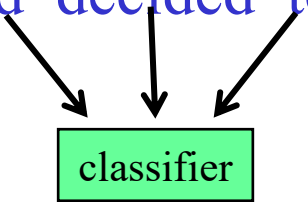
John saw the saw and decided to take it to the table.

# 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.

classifier

VBD

# 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.

classifier

DT

# 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.

classifier

NN

# Sequence Labeling as Classification

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).
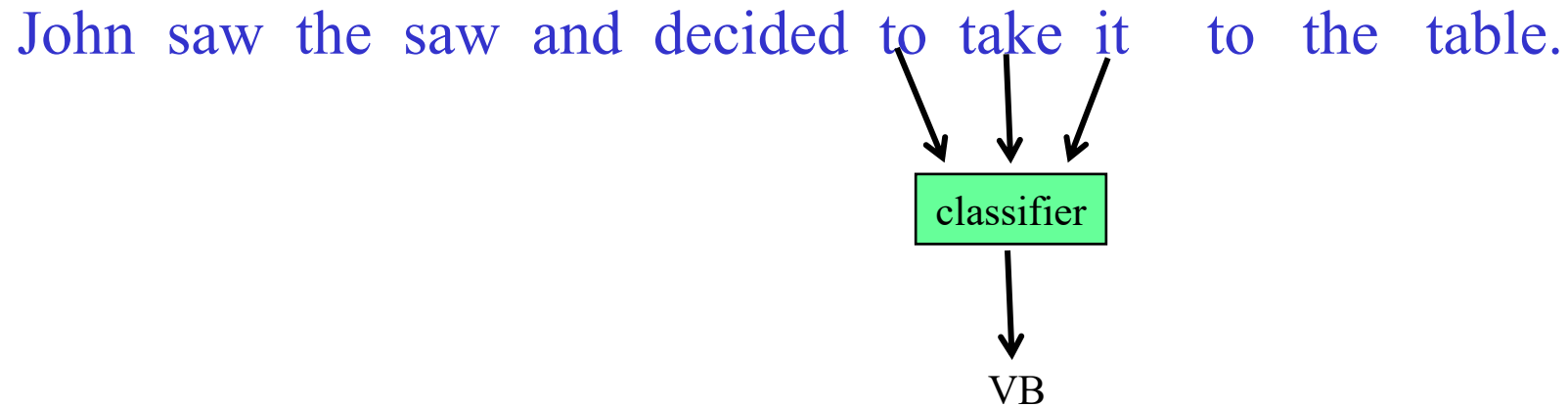
John saw the saw and decided to take it to the table.



classifier

CC

# 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.

classifier

VBD

# Sequence Labeling as Classification

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).
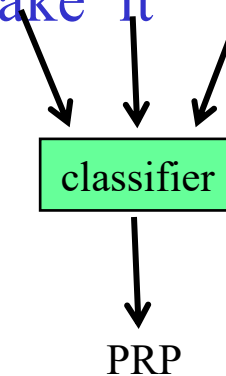
John  saw  the  saw  and  decided  to  take  it    to   the   table.

classifier

TO

# Sequence Labeling as Classification

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

# Sequence Labeling as Classification

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

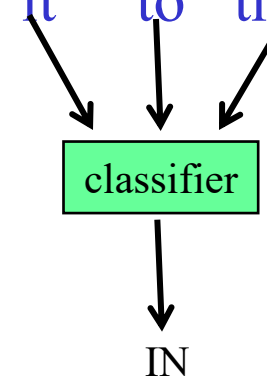John   saw   the   saw   and   decided   to   take   it   to   the   table.

classifier

PRP

# Sequence Labeling as Classification

Classify each token independently but use as input features, information about the surrounding tokens (sliding window).
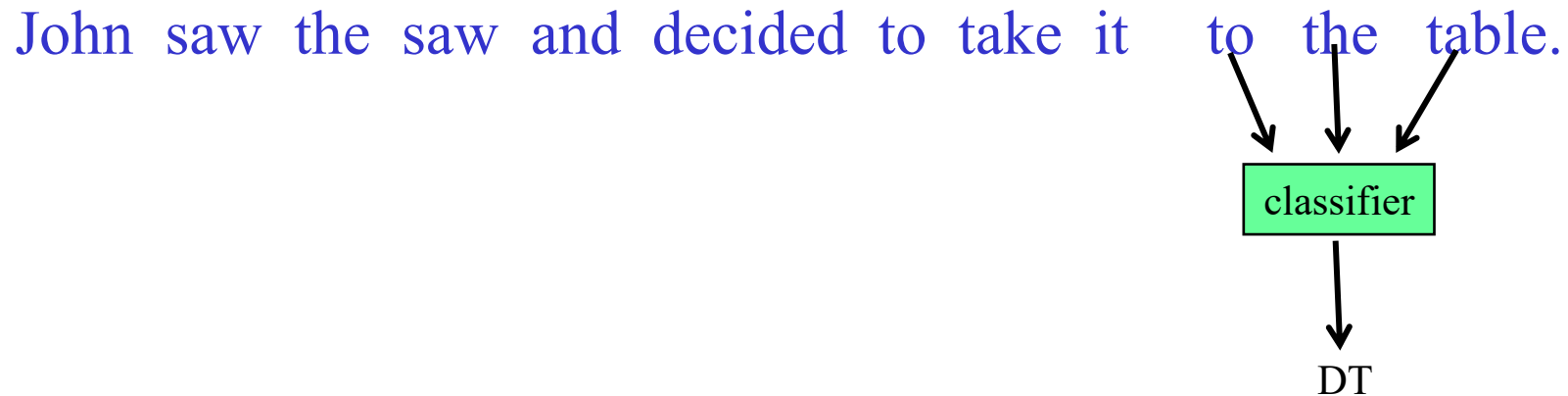
John  saw  the  saw  and  decided  to  take  it    to   the   table.

classifier

IN

# 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.

classifier

DT