# NLP Review
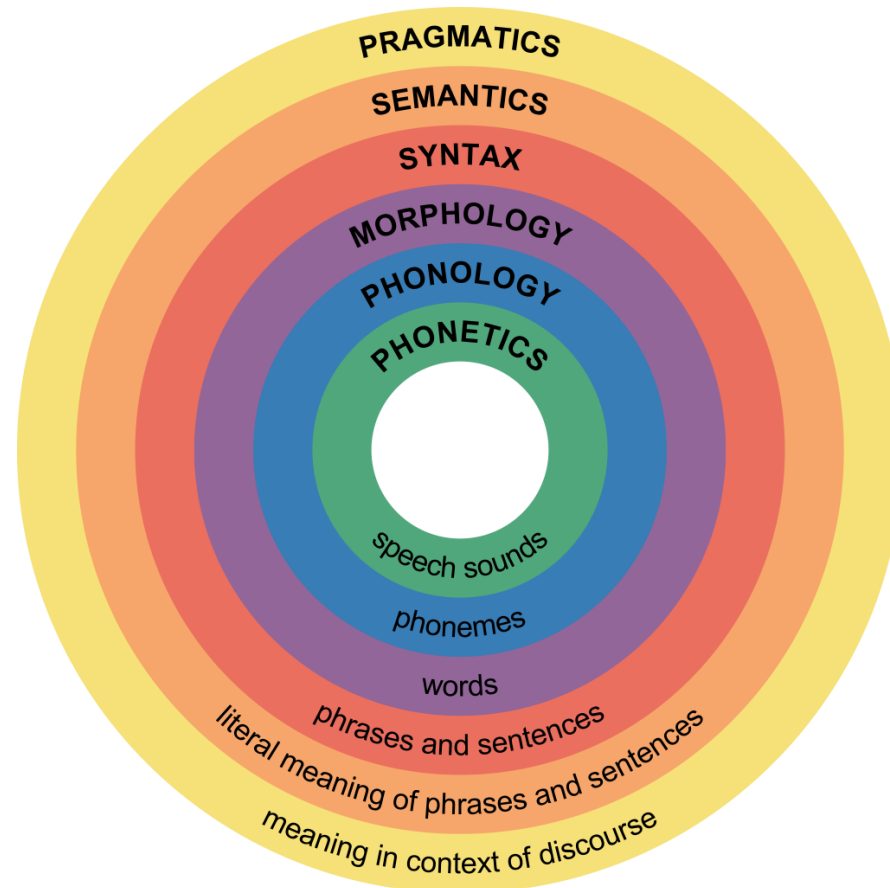
CMSC 473/673 - NATURAL LANGUAGE PROCESSING

*Slides modified from Dr. Frank Ferraro*

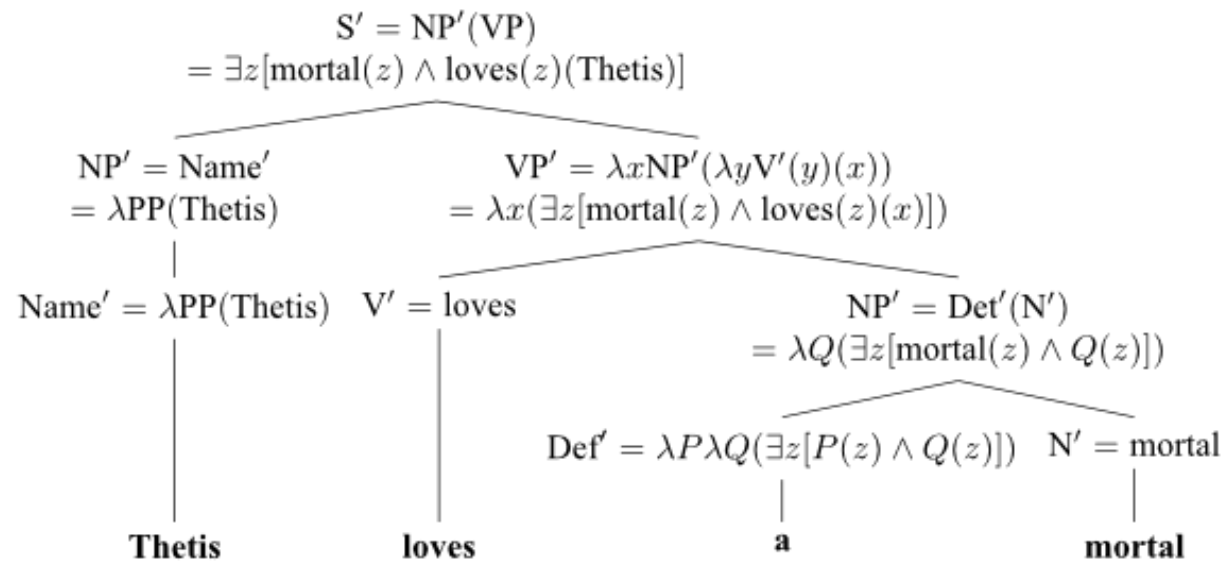# Linguistics

The study of language



PRAGMATICS — meaning in context of discourse
SEMANTICS — literal meaning of phrases and sentences
SYNTAX — phrases and sentences
MORPHOLOGY — words
PHONOLOGY — phonemes
PHONETICS — speech sounds

https://en.wikipedia.org/wiki/Morphology_(linguistics)#/media/File:Major_levels_of_linguistic_structure.svg

# Semantics

Meaning

$$S' = NP'(VP)$$
$$= \exists z[\mathrm{mortal}(z) \wedge \mathrm{loves}(z)(\mathrm{Thetis})]$$

$$NP' = \mathrm{Name}'$$
$$= \lambda P P(\mathrm{Thetis})$$

$$VP' = \lambda x NP'(\lambda y V'(y)(x))$$
$$= \lambda x(\exists z[\mathrm{mortal}(z) \wedge \mathrm{loves}(z)(x)])$$

$$\mathrm{Name}' = \lambda P P(\mathrm{Thetis}) \quad V' = \mathrm{loves}$$

$$NP' = \mathrm{Det}'(N')$$
$$= \lambda Q(\exists z[\mathrm{mortal}(z) \wedge Q(z)])$$

$$\mathrm{Def}' = \lambda P \lambda Q(\exists z[P(z) \wedge Q(z)]) \quad N' = \mathrm{mortal}$$
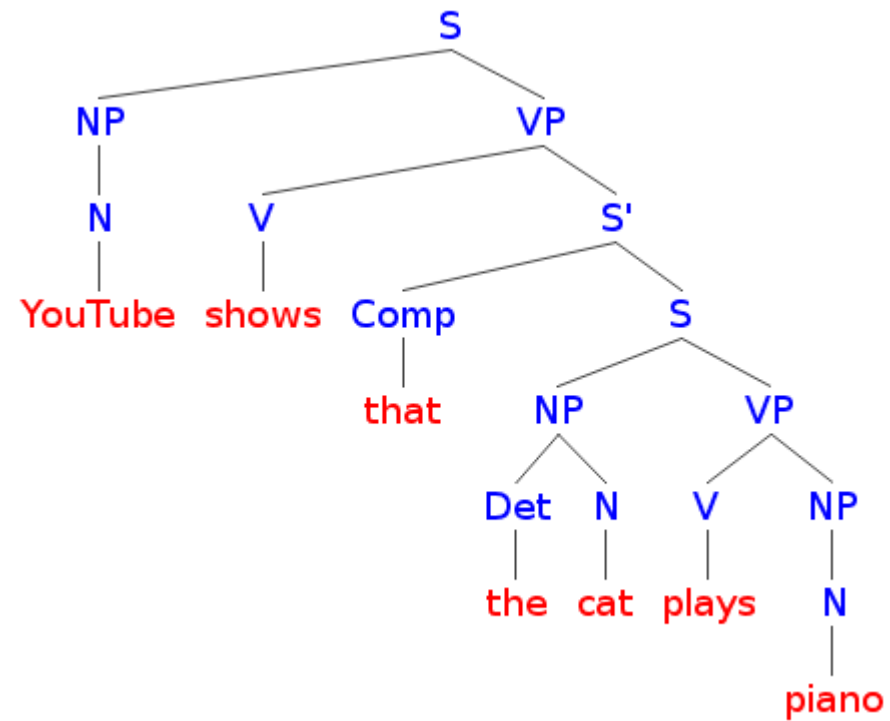
**Thetis**     **loves**     **a**     **mortal**

# Syntax

Grammar



https://allthingslinguistic.com/post/100617668093/how-to-draw-syntax-trees-part-3-type-1-a
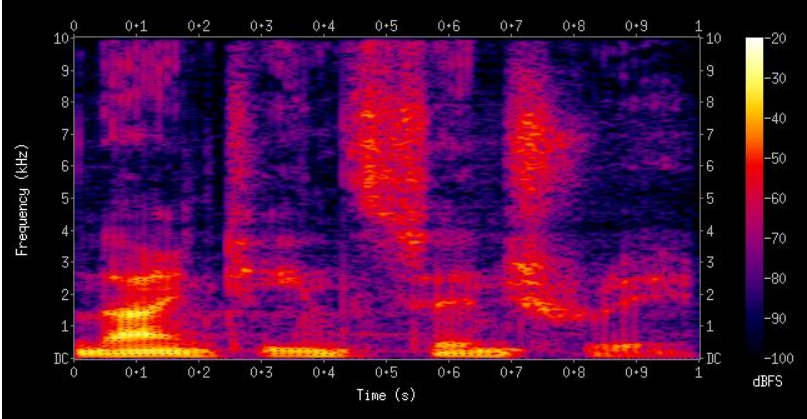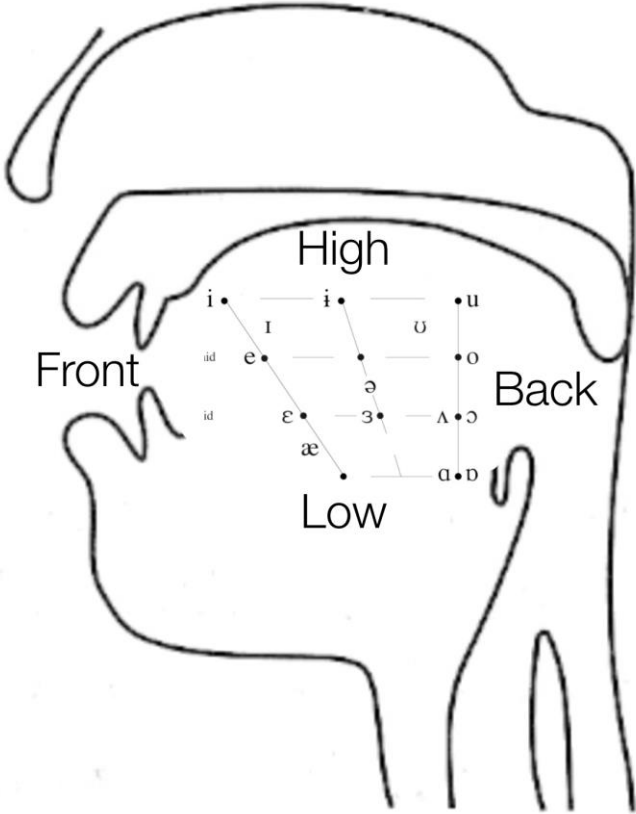
# Phonology

Processing of sounds


https://upload.wikimedia.org/wikipedia/commons/a/a5/Tsunami_by_hokusai_19th_century.jpg

tsunami

⬇

sunami

| /ðɪs/ *this* | DEP | *CODA | MAX |
|---|---|---|---|
| a. ☞ [dɪs] | | * | |
| b. ☞ [dɪ] | | | * |
| c. [dɪ.sə] | *! | | |

https://pubs.asha.org/doi/10.1044/0161-1461%282001/022%29

# Phonetics

Physical production/understanding of sounds



High

Front

Back

Low

# ML/NLP Framework



**instances**

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

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

θ

θ

1.352
36.26
262.4
925
...

# 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>

**instances**

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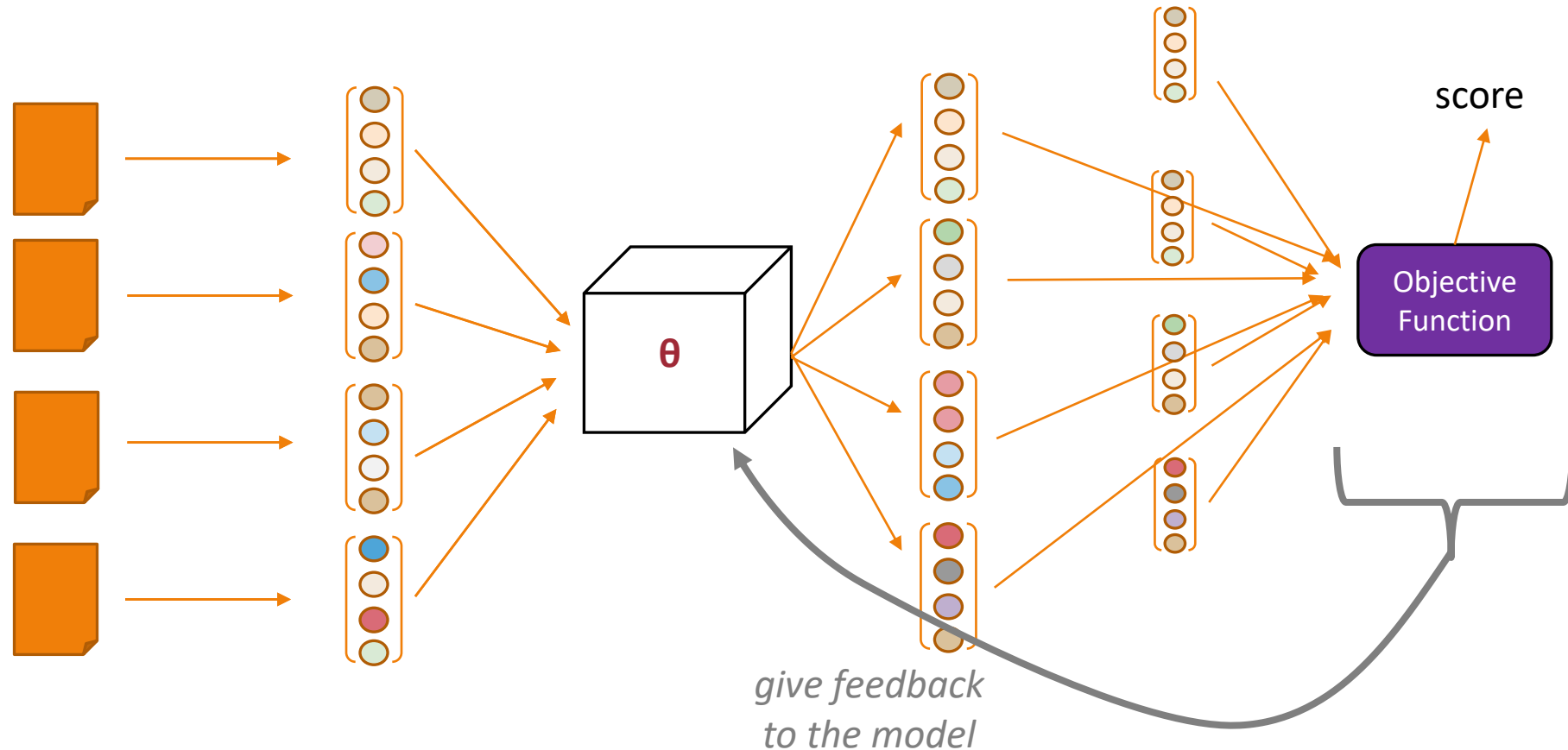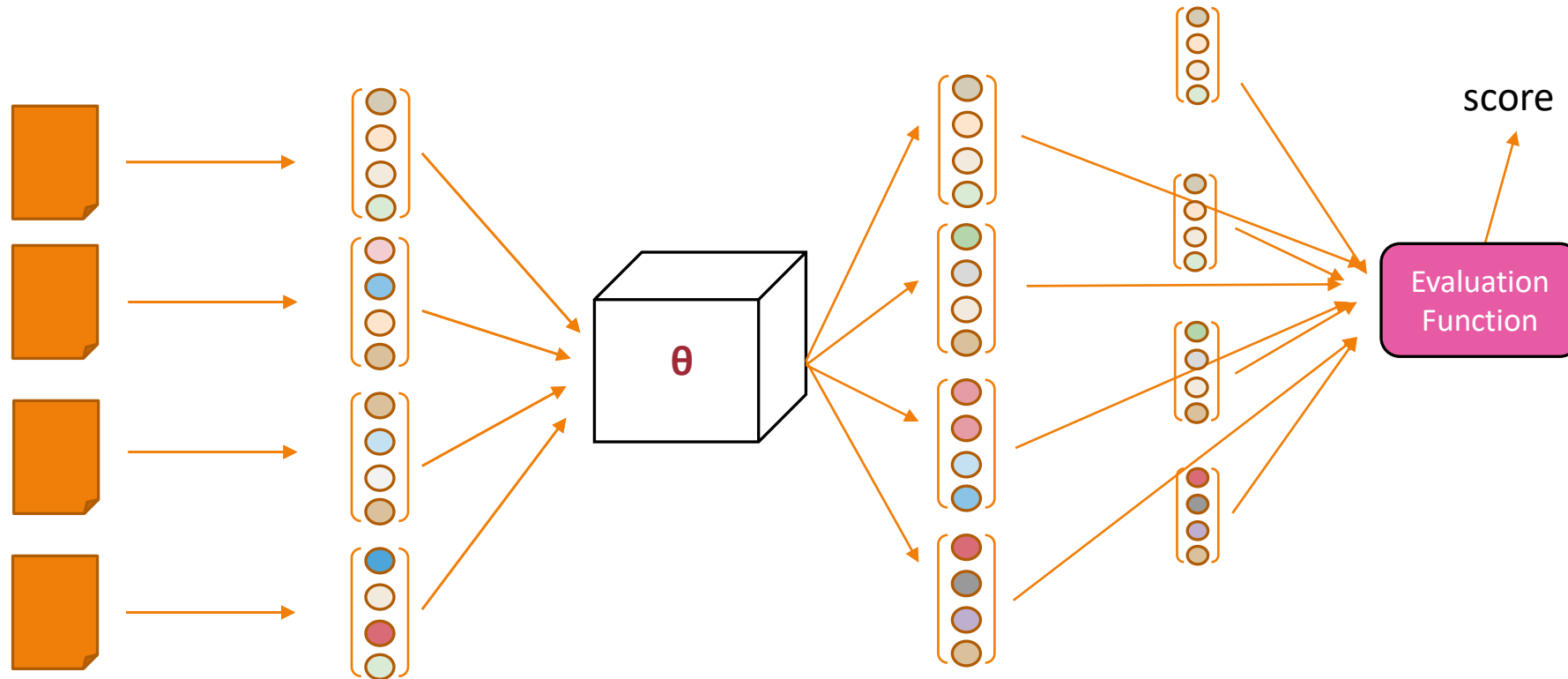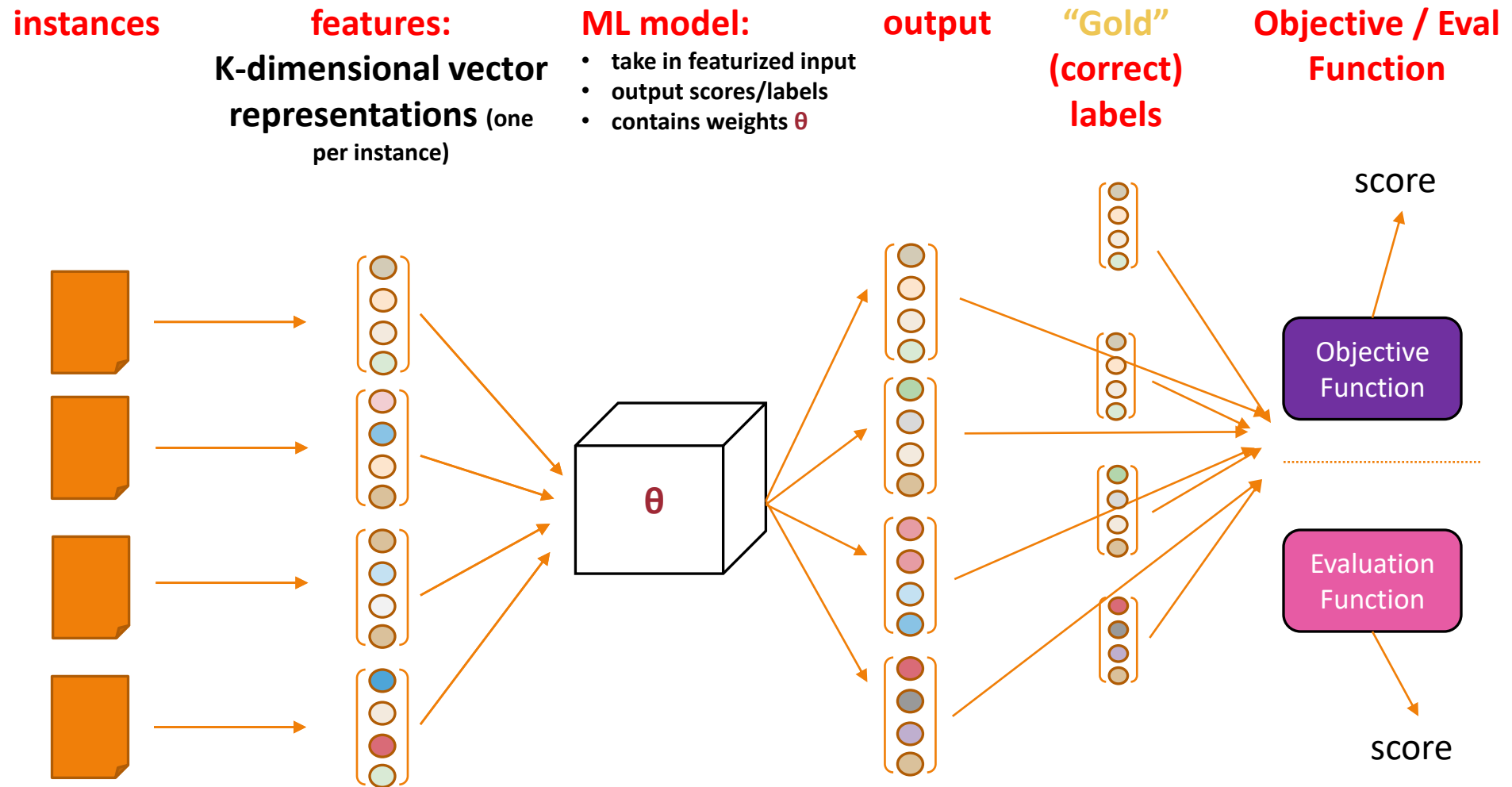**K-dimensional vector representations** (one per instance)

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

**output**

**"Gold" (correct) labels**

**Objective Function**

θ

score

Objective Function

*give feedback to the model*

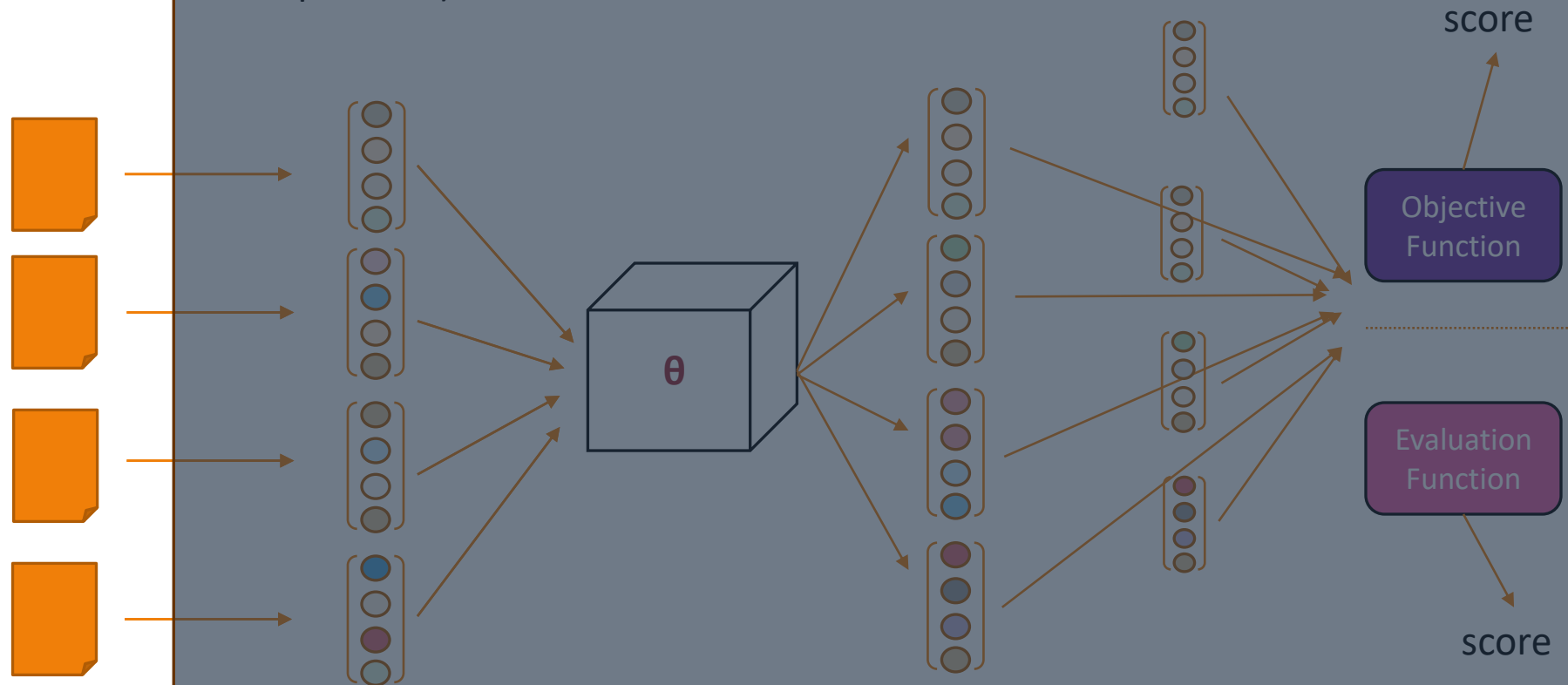# ML/NLP Framework for Prediction

**instances**

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

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

**output**

**"Gold" (correct) labels**

**Objective / Eval Function**

score

θ

Objective Function

Evaluation Function

score

# Where does the data come from?

Corpus (plural: corpora)
- Literally a "body" of text

Languages with few corpora are called "low-resource languages"
- This might not mean the language is endangered!

We can collect corpora in a few different ways:
- Curation: data tagged & organized by experts
- Internet: data "scraped" from open-access sources (Wikipedia, Reddit)
  - Or data collected with permission from closed sources (Facebook, texts) – more rare
- Elicitation: carefully getting participants to produce language (lab studies, crowdsourcing, field studies)
- Pre-existing corpora

**!** Facebook has gotten into trouble several times for using data or manipulating people's feeds without their permission

# Benchmarking

Tasks!

If you want people to work on your problem, make it easy for them to get started and to measure their progress. Provide:

- Test data, for evaluating the final systems
- Development data, for measuring whether a change to the system helps, and for tuning parameters
- An evaluation metric (formula for measuring how well a system does on the dev or test data)
- A program for computing the evaluation metric
- Labeled training data and other data resources
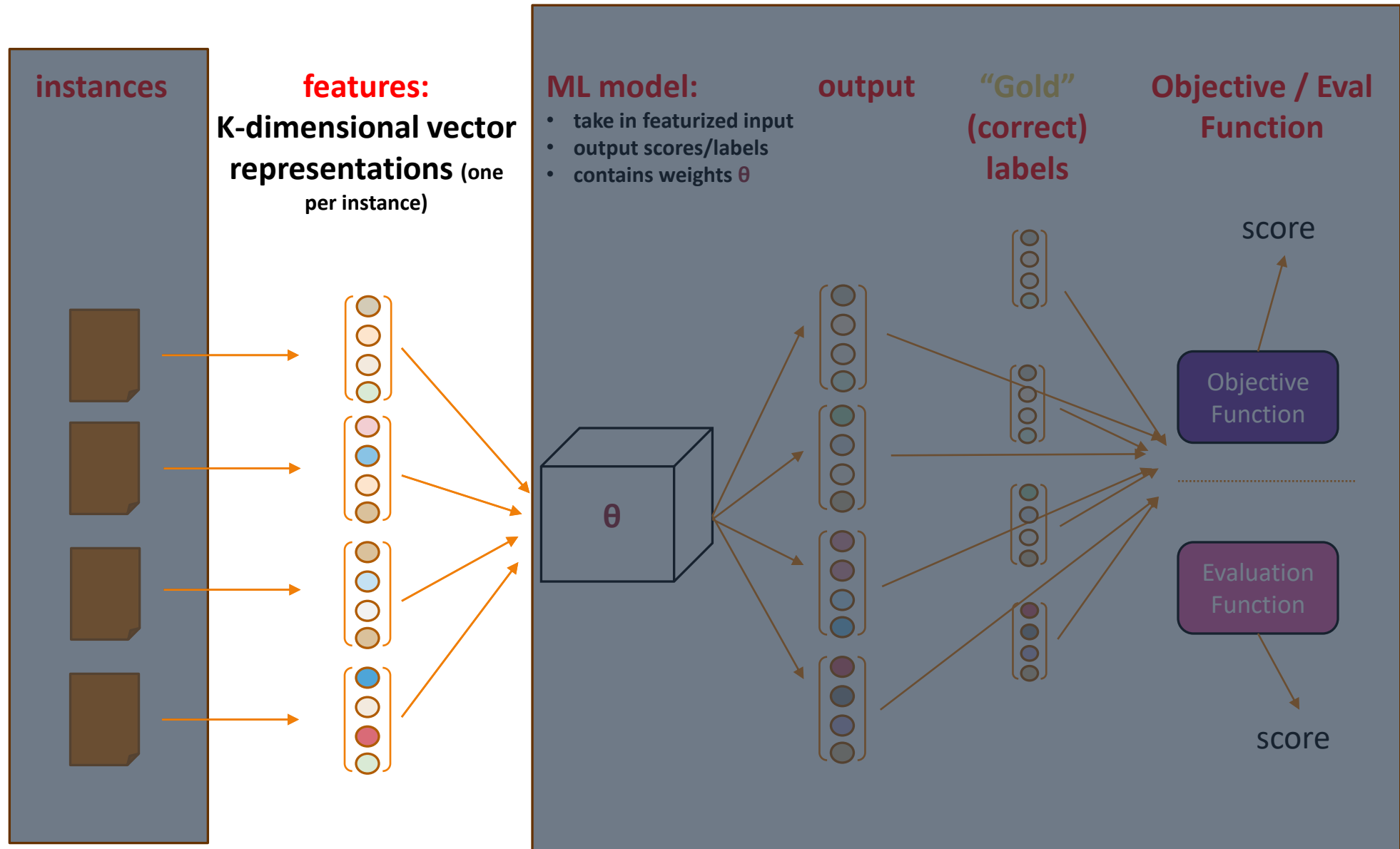- A prize? – with clear rules on what data can be used

# Tokens vs Types

The film got a great opening and the film went on to become a hit .


**Vocabulary:** the words (items) you know

**Type:** an element of the vocabulary.

**Token:** an instance of that type in running text.

**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

# ML Term: "Featurization"

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

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

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

Each of these is a feature
(/feature function)

# Overview of Featurization

Common goal: probabilistic classifier p(y | x)

Often done by defining **features** between x and y that are meaningful

- Denoted by a **general vector of K features**

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

Features can be thought of as "soft" rules

- E.g., POSITIVE sentiments tweets *may* be more likely to have the word "happy"

# Representing Linguistic Information

| | | |
|---|---|---|
| **User-defined** | Bag of words / one-hot encoding | Assign each word to some index i, where $0 \leq i < V$ |
| | | Represent each word w with a V-dimensional **binary** vector $e_w$, where $e_{w,i} = 1$ and 0 otherwise |
| **Model-produced** | Dense embedding | Let E be some *embedding size* (often 100, 200, 300, etc.) |
| | | Represent each word w with an E-dimensional **real-valued** vector $e_w$ |

# Bag-of-words

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)

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

Amazon acquired MGM in 2022, taking over a sprawling library that includes more than 4,000 feature films and 17,000 television shows. The tech behemoth also earned the rights to distribute all the Bond movies, but the new deal solidifies the company's oversight of Bond's big-screen future.

TECH

NOT TECH

*feature extraction*

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

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

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

Core assumption: the label can be predicted from counts of individual word types

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

Amazon acquired MGM in 2022, taking over a sprawling library that includes more than 4,000 feature films and 17,000 television shows. The tech behemoth also earned the rights to distribute all the Bond movies, but the new deal solidifies the company's oversight of Bond's big-screen future.

TECH

NOT TECH

*feature extraction*

| feature $f_i(x)$ | value |
|---|---|
| Amazon | 1 |
| acquired | 1 |
| behemoth | 1 |
| Bond | 2 |
| ... | |
| sniffle | 0 |
| ... | |

Core assumption: the label can be predicted from counts of individual word types

# How have we represented words?

Each word is a distinct item
- ◦ Bijection between the strings and unique integer ids:
- ◦ "cat" --> 3, "kitten" --> 792 "dog" --> 17394

- ◦ Are "cat" and "kitten" similar?

Equivalently: "One-hot" encoding
- ◦ Represent each word type w with a vector the size of the vocabulary
- ◦ This vector has V-1 zero entries, and 1 non-zero (one) entry

# One-Hot Encoding Example

Let our vocab be {a, cat, saw, mouse, happy}

V = # types = 5

Assign:

| a | 4 |
|---|---|
| cat | 2 |
| saw | 3 |
| mouse | 0 |
| happy | 1 |

How do we represent "cat?"

$$e_{\text{cat}} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

How do we represent "happy?"

$$e_{\text{happy}} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

# Useful Terminology: n-gram

Within a larger string (e.g., sentence),
a contiguous sequence of n items (e.g., words)

**Sometimes people use a bag of n-grams!**

Colorless green ideas sleep furiously

| n | Commonly called | History Size (Markov order) | Example n-gram ending in "furiously" |
|---|---|---|---|
| 1 | unigram | 0 | furiously |
| 2 | bigram | 1 | sleep furiously |
| 3 | trigram (3-gram) | 2 | ideas sleep furiously |
| 4 | 4-gram | 3 | green ideas sleep furiously |
| n | n-gram | n-1 | $w_{i-n+1} \dots w_{i-1} w_i$ |

# Representing Linguistic Information

| User-defined | Bag of words / one-hot encoding | Assign each word to some index i, where $0 \leq i < V$ |
| | | Represent each word w with a V-dimensional **binary** vector $e_w$, where $e_{w,i} = 1$ and 0 otherwise |

| Model-produced | Dense embedding | Let E be some *embedding size* (often 100, 200, 300, etc.) |
| | | Represent each word w with an E-dimensional **real-valued** vector $e_w$ |

# A Dense Representation (E=2)

[0.00315225, 0.00315225, 0.00547597, 0.00741556, 0.00912817, 0.01068435, 0.01212381, 0.01347162, 0.01474487, 0.0159558 ]

# Distributional Representations

A dense, "low"-dimensional vector representation

Many values are not 0 (or at least less sparse than one-hot)

Up till ~2013: E could be any size

2013-present: E << vocab

An E-dimensional vector, often (but not always) real-valued

These are also called
- **embeddings**
- **Continuous representations**
- **(word/sentence/...) vectors**
- **Vector-space models**

# (Some) Properties of Embeddings

## Capture "like" (similar) words

| target: | Redmond | Havel | ninjutsu | graffiti | capitulate |
|---|---|---|---|---|---|
| | Redmond Wash. | Vaclav Havel | ninja | spray paint | capitulation |
| | Redmond Washington | president Vaclav Havel | martial arts | grafitti | capitulated |
| | Microsoft | Velvet Revolution | swordsmanship | taggers | capitulating |

## Capture relationships



vector(*'king'*) – vector(*'man'*) + vector(*'woman'*) ≈ vector('queen')

vector(*'Paris'*) - vector(*'France'*) + vector(*'Italy'*) ≈ vector('Rome')

Mikolov et al. (2013)

# Three Common Kinds of Embedding Models

1. Co-occurrence matrices
2. Matrix Factorization: Singular value decomposition/Latent Semantic Analysis, Topic Models
3. Neural-network-inspired models (skip-grams, CBOW)

# Shared Intuition

Model the meaning of a word by "embedding" in a vector space

The meaning of a word is a vector of numbers

Contrast: word meaning is represented in many computational linguistic applications by a vocabulary index ("word number 545") or the string itself

# Three Common Kinds of Embedding Models

1. Co-occurrence matrices
2. Matrix Factorization: Singular value decomposition/Latent Semantic Analysis, Topic Models
3. Neural-network-inspired models (skip-grams, CBOW)

# Co-occurrence Matrix

*words*

Acquire basic contextual statistics (often counts) for each word type v  via *correlate*.

*correlates*

j

i

Per-correlated word statistics

# Co-occurrence Matrix

Acquire basic contextual statistics (often counts) for each word type v  via *correlate*:

For example:

   documents
   ◦ Record how often a word occurs in each document

*words*

j

*correlates*
i

Per-correlated word statistics

\# correlates = # documents

# Co-occurrence Matrix

Acquire basic contextual statistics (often counts) for each word type v via *correlate*:

For example:

documents

surrounding context words

◦ Record how often v occurs with other word types u

*words*

*correlates*

Per-correlated word statistics

j

i

# correlates =
# word types

# Co-occurrence Matrix

Acquire basic contextual statistics (often counts) for each word type v  via *correlate*:

For example:

documents

surrounding context words

linguistic annotations (POS tags, syntax)

…

*words*

*correlates*

j

i

Per-correlated word statistics

# "You shall know a word by the company it keeps!" Firth (1957)

document (↓)-word (→) count matrix

|  | battle | soldier | fool | clown |
|---|---|---|---|---|
| *As You Like It* | 1 | 2 | 37 | 6 |
| *Twelfth Night* | 1 | 2 | 58 | 117 |
| *Julius Caesar* | 8 | 12 | 1 | 0 |
| *Henry V* | 15 | 36 | 5 | 0 |

*basic bag-of-words counting*



I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

| | |
|---|---|
| it | 6 |
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| yet | 1 |
| would | 1 |
| whimsical | 1 |
| times | 1 |
| sweet | 1 |
| satirical | 1 |
| adventure | 1 |
| genre | 1 |
| fairy | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| ... | ... |

# "You shall know a word by the company it keeps!" Firth (1957)

document (↓)-word (→) count matrix

|  | battle | soldier | fool | clown |
|---|---|---|---|---|
| *As You Like It* | 1 | 2 | 37 | 6 |
| *Twelfth Night* | 1 | 2 | 58 | 117 |
| *Julius Caesar* | 8 | 12 | 1 | 0 |
| *Henry V* | 15 | 36 | 5 | 0 |

*Assumption: Two documents are similar if their vectors are similar*

# "You shall know a word by the company it keeps!" Firth (1957)

document (↓)-word (→) count matrix

|  | battle | soldier | fool | clown |
|---|---|---|---|---|
| *As You Like It* | 1 | 2 | 37 | 6 |
| *Twelfth Night* | 1 | 2 | 58 | 117 |
| *Julius Caesar* | 8 | 12 | 1 | 0 |
| *Henry V* | 15 | 36 | 5 | 0 |

*Assumption: Two words are similar if their vectors are similar???*

*Issue: Count word vectors are very large, sparse, and skewed!*

# "You shall know a word by the company it keeps!" Firth (1957)

context (↓)-word (→) count matrix

|  | apricot | pineapple | digital | information |
|---|---|---|---|---|
| aardvark | 0 | 0 | 0 | 0 |
| computer | 0 | 0 | 2 | 1 |
| data | 0 | 10 | 1 | 6 |
| pinch | 1 | 1 | 0 | 0 |
| result | 0 | 0 | 1 | 4 |
| sugar | 1 | 1 | 0 | 0 |

*Context: those other words within a small "window" of a target word*

a cloud [ computer stores digital data on ] a remote computer

# "You shall know a word by the company it keeps!" Firth (1957)

**context** (↓)-word (→) count matrix

| | apricot | pineapple | digital | information |
|---|---|---|---|---|
| aardvark | 0 | 0 | 0 | 0 |
| computer | 0 | 0 | 2 | 1 |
| data | 0 | 10 | 1 | 6 |
| pinch | 1 | 1 | 0 | 0 |
| result | 0 | 0 | 1 | 4 |
| sugar | 1 | 1 | 0 | 0 |

*Context: those other words within a small "window" of a target word*

*Assumption: Two words are similar if their vectors are similar*

*Issue: Count word vectors are very large, sparse, and skewed!*

# Pointwise Mutual Information (PMI): Dealing with Problems of Raw Counts

Raw word frequency is not a great measure of association between words

It's very skewed: "the" and "of" are very frequent, but maybe not the most discriminative

We'd rather have a measure that asks whether a context word is **particularly informative** about the target word.

(Positive) Pointwise Mutual Information ((P)PMI)

**Pointwise mutual information**:

Do events x and y co-occur more than if they were independent?

probability words x and y occur together
(in the same context/window)

$$\mathrm{PMI}(x, y) = \log \frac{p(x, y)}{p(x)p(y)}$$

probability that word x occurs    probability that word y occurs

# Three Common Kinds of Embedding Models

1. Co-occurrence matrices
2. Matrix Factorization: Singular value decomposition/Latent Semantic Analysis, Topic Models
3. Neural-network-inspired models (skip-grams, CBOW)

# Word2Vec

Mikolov et al. (2013; NeurIPS): "Distributed Representations of Words and Phrases and their Compositionality"

Revisits the context-word approach

Learn a model p($c \mid w$) to predict a context word $c$ from a target word $w$

Learn two types of vector representations
- $h_c \in \mathbb{R}^E$: vector embeddings for each context word
- $v_w \in \mathbb{R}^E$: vector embeddings for each target word

$$p(c \mid w) \propto \exp(h_c^T v_w)$$

# Word2Vec

**context** (↓)-word (→) count matrix

|  | **apricot** | **pineapple** | **digital** | **information** |
|---|---|---|---|---|
| aardvark | 0 | 0 | 0 | 0 |
| computer | 0 | 0 | 2 | 1 |
| data | 0 | 10 | 1 | 6 |
| pinch | 1 | 1 | 0 | 0 |
| result | 0 | 0 | 1 | 4 |
| sugar | 1 | 1 | 0 | 0 |

*Context: those other words within a small "window" of a target word*

$$\max_{h,v} \sum_{c,w \text{ pairs}} \text{count}(c,w) \log p(c \mid w)$$

# Word2Vec

**context** (↓)-word (→) count matrix

| | apricot | pineapple | digital | information |
|---|---|---|---|---|
| aardvark | 0 | 0 | 0 | 0 |
| computer | 0 | 0 | 2 | 1 |
| data | 0 | 10 | 1 | 6 |
| pinch | 1 | 1 | 0 | 0 |
| result | 0 | 0 | 1 | 4 |
| sugar | 1 | 1 | 0 | 0 |

*Context: those other words within a small "window" of a target word*

$$\max_{h,v} \sum_{c,w \text{ pairs}} \text{count}(c,w) \left[ h_c^T v_w - \log(\sum_u \exp(h_u^T v_w))) \right]$$

# Example (Tensorflow)

The wide road shimmered in the hot sun.

`tf.keras.preprocessing.sequence.skipgrams`

↓

| (wide, road) | ⋯ | (road, shimmered) | (hot, sun) | ⋯ | (the, hot) |
|---|---|---|---|---|---|
| ( 2, 3) | ⋯ | (3, 4) | (6, 7) | ⋯ | (1, 6) |

`tf.random.log_uniform_candidate_sampler`
`(negative_samples = 4)`

↓                                        ↓

| (wide, road) | | (wide, sun) | (wide, hot) | (wide, temperature) | (wide, code) |
|---|---|---|---|---|---|
| ( 2, 3) | | (2, 7) | (2,6) | (2, 23) | (2, 2196) |

`concat and add label (pos:1/neg:0)`
↓

| (wide, road) | (wide, sun) | (wide, hot) | (wide, temperature) | (wide, code) |
|---|---|---|---|---|
| (2, 3) | (2, 7) | (2,6) | (2, 23) | (2, 2196) |
| 1 | 0 | 0 | 0 | 0 |

`build context words and labels for all vocab words`
↓

| Word | Context words | | | | | | Labels | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 3 | 7 | 6 | 23 | 2196 | ⇒ | 1 | 0 | 0 | 0 | 0 |
| 23 | 12 | 6 | 94 | 17 | 1085 | ⇒ | 1 | 0 | 0 | 0 | 0 |
| 84 | 784 | 11 | 68 | 41 | 453 | ⇒ | 1 | 0 | 0 | 0 | 0 |
| ⋮ | | | | | | | | | | | |
| V | 45 | 598 | 1 | 117 | 43 | ⇒ | 1 | 0 | 0 | 0 | 0 |

# Word2Vec Vectors are Weights of a NN

# FastText

P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, "**Enriching Word Vectors with Subword Information**," *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 135–146, 2017, doi: 10.1162/tacl_a_00051.

Main idea: learn **character n-gram embeddings** for the target word (not context) and modify the word2vec model to use these

Pre-trained models in 150+ languages
◦ https://fasttext.cc

# FastText Details

Main idea: learn **character n-gram embeddings** and for the target word (not the context) modify the word2vec model to use these

Original word2vec:

$$p(c\,|w) \propto \exp(h_c^T v_w)$$

FastText:

$$p(c\,|w) \propto \exp\left(h_c^T\left(\sum_{\text{n−gram } g \text{ in } w} z_g\right)\right)$$

# FastText Details

Main idea: learn **character n-gram embeddings** and for the target word (not the context) modify the word2vec model to use these

$$p(c \,|w) \propto \exp\left( h_c^T \left( \Sigma_{\text{n−gram } g \text{ in } w} \, z_g \right) \right)$$

*decompose*

`fluffy` ➔ `fl flu luf uff ffy fy`

Sub-word units like this have become an important part of today's NLP work!

# FastText Details

Main idea: learn **character n-gram embeddings** and for the target word (not the context) modify the word2vec model to use these

$$p(c \mid w) \propto \exp\left( h_c^T \left( \sum_{\text{n–gram } g \text{ in } w} z_g \right) \right)$$

*decompose*

`fluffy` → `fl  flu  luf  uff  ffy  fy`



← 

*To deterministically compute word embeddings*

*Learn n-gram embeddings*

# Contextual Word Embeddings

Word2vec-based models are not context-dependent

Single word type → single word embedding

If a single word type can have different meanings…

bank, bass, plant,…

… why should we only have one embedding?

Entire task devoted to classifying these meanings:
**Word Sense Disambiguation**

# Contextual Word Embeddings

Growing interest in this

Off-the-shelf is a bit more difficult
- Download and run a model
- Can't just download a file of embeddings

Two to know about (with code):
- ELMo: "Deep contextualized word representations" Peters et al. (2018; NAACL)
- https://allennlp.org/elmo
- BERT: "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding" Devlin et al. (2019; NAACL)
  - https://github.com/google-research/bert

# Evaluating Vector Embeddings

# Cosine: Measuring Similarity

Given 2 target words v and w how similar are their vectors?

Dot product or inner product from linear algebra

$$\text{dot-product}(\vec{v}, \vec{w}) = \vec{v} \cdot \vec{w} = \sum_{i=1}^{N} v_i w_i = v_1 w_1 + v_2 w_2 + \ldots + v_N w_N$$

◦ High when two vectors have large values in same dimensions, low for orthogonal vectors with zeros in complementary distribution

Correct for high magnitude vectors $\dfrac{\vec{a} \cdot \vec{b}}{|\vec{a}||\vec{b}|}$

# Cosine Similarity

Divide the dot product by the length of the two vectors

$$\frac{\vec{a} \cdot \vec{b}}{|\vec{a}||\vec{b}|}$$

This is the cosine of the angle between them

$$\vec{a} \cdot \vec{b} = |\vec{a}||\vec{b}| \cos \theta$$

$$\frac{\vec{a} \cdot \vec{b}}{|\vec{a}||\vec{b}|} = \cos \theta$$

# Cosine Similarity

Divide the dot product by the length of the two vectors

$$\frac{\vec{a} \cdot \vec{b}}{|\vec{a}||\vec{b}|}$$

This is the cosine of the angle between them

$$\vec{a} \cdot \vec{b} = |\vec{a}||\vec{b}| \cos \boldsymbol{\theta}$$

$$\frac{\vec{a} \cdot \vec{b}}{|\vec{a}||\vec{b}|} = \cos \boldsymbol{\theta}$$



convertible • Chevrolet • Ford •
car • cargo capacity truck • off-road •
fuel efficiency • towing •
80°
Origin

*https://upload.wikimedia.org/wikipedia/commons/2/23/CosineSimilarity.png*

# Example: Word Similarity

$$\cos(x, y) = \frac{\sum_i x_i y_i}{\sqrt{\sum_i x_i^2} \sqrt{\sum_i y_i^2}}$$

| | Dim. 1 | Dim. 2 | Dim. 3 |
|---|---|---|---|
| apricot | 2 | 0 | 0 |
| digital | 0 | 1 | 2 |
| information | 1 | 6 | 1 |

cosine(apricot,information) = $\dfrac{2 + 0 + 0}{\sqrt{4 + 0 + 0}\sqrt{1 + 36 + 1}} = 0.1622$

cosine(digital,information) = $\dfrac{0 + 6 + 2}{\sqrt{0 + 1 + 4}\sqrt{1 + 36 + 1}} = 0.5804$

cosine(apricot,digital) = $\dfrac{0 + 0 + 0}{\sqrt{4 + 0 + 0}\sqrt{0 + 1 + 4}} = 0.0$

# Cosine Similarity Range



- Angle θ close to 0
- Cos(θ) close to 1
- **Similar vectors**

- Angle θ close to 90
- Cos(θ) close to 0
- **Orthogonal vectors**

https://www.learndatasci.com/glossary/cosine-similarity/

**instances**

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

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

**θ**

**output**    **"Gold" (correct) labels**    **Objective / Eval Function**

score

Objective
Function

Evaluation
Function

score

# 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**.

θ

1.352
36.26
262.4
925
…

Input → Model → Output

# Types of models

## CLASSIFICATION

Model outputs comes from a finite set of values

Discrete result

*Examples*:
- What type of animal is this a picture of?
- Predicting the weather (sunny, cloudy, or rainy?)
- Ranking: Is this result *better* than this result?

## REGRESSION

Model outputs are continuous values

Continuous result

*Examples*:
- How far will I move if I drive my motors at this speed for 1 second?
- Predicting the weather (temperature)
- Ranking: *how good* is this result?

# Types of models



Classification                    Regression

# What are some other examples of these?

## CLASSIFICATION

Tone tagging

Sentiment classification

Named entity recognition

## REGRESSION

Quantity/scale of how much it sounds like a specific author

Numerical sentiment value

Political "score" from document

Likelihoods

Predicted Goodreads score

# Classification

# Modeling

Classification/
Text Processing

$$P(y \mid x)$$

Language
Model (LM) /
Generation

$$P(w_t | w_{t-1}, w_{t-2} \ldots)$$

A language model is used to **generate** the next word(s) given a history of words.

# 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 |

# Maxent Models for Classification: Discriminatively or …

Directly model the posterior

$$p(Y \mid X) = \mathbf{maxent}(X; Y)$$

Discriminatively trained classifier

"Discriminative classifiers like logistic regression instead learn what features from the input are most useful to discriminate between the different possible classes."
SLP, ch. 4

# Bayes' Rule

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

Likelihood — $P(X|Y)$

Prior — $P(Y)$

Posterior — $P(Y|X)$

**Posterior:**
probability of event Y with <u>knowledge that X has occurred</u>
NLP pg. 478

**Likelihood:**
probability of event X given that Y <u>has occurred</u>
NLP pg. 478

**Prior:**
probability of event X occurring (regardless of what other events happen)
NLP pg. 478

# Bayes' Rule

$$P(c|d) = \frac{P(d|c) \cdot P(c)}{P(d)}$$

$$P\left( \text{Entailed} \,\middle|\, \begin{array}{l} \text{s: Michael Jordan, coach Phil Jackson and the star} \\ \text{cast, including Scottie Pippen, took the Chicago} \\ \text{Bulls to six National Basketball Association} \\ \text{championships.} \\ \text{h: The Bulls basketball team is based in Chicago.} \end{array} \right)$$

$$= \frac{P\left( \begin{array}{l} \text{s: Michael Jordan, coach Phil Jackson and the star} \\ \text{cast, including Scottie Pippen, took the Chicago} \\ \text{Bulls to six National Basketball Association} \\ \text{championships.} \\ \text{h: The Bulls basketball team is based in Chicago.} \end{array} \,\middle|\, \text{Entailed} \right) \cdot P\left( \text{Entailed} \right)}{P\left( \begin{array}{l} \text{s: Michael Jordan, coach Phil Jackson and the star} \\ \text{cast, including Scottie Pippen, took the Chicago} \\ \text{Bulls to six National Basketball Association} \\ \text{championships.} \\ \text{h: The Bulls basketball team is based in Chicago.} \end{array} \right)}$$

# Maxent Models for Classification: Discriminatively or Generatively Trained

Directly model the posterior

$$p(Y \mid X) = \mathbf{maxent}(X; Y)$$

**Discriminatively** trained classifier

Model the posterior with Bayes rule

$$p(Y \mid X) \propto p(X \mid Y)p(Y)$$

**Generatively** trained classifier with maxent-based language model

# Maximum Entropy (Log-linear) Models For Discriminatively Trained Classifiers

$$p(y \mid x) = \text{maxent}(x, y)$$

Modeled jointly!

# Core Aspects to Maxent Classifier p(y|x)

We need to define:

- **features** $f(x)$ from x that are meaningful;

- **weights** $\theta$ (at least one per feature, often one per feature/label combination) to say how important each feature is; and

- a way to **form probabilities** from $f$ and $\theta$

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

Amazon acquired MGM in 2022, taking over a sprawling library that includes more than 4,000 feature films and 17,000 television shows. The tech behemoth also earned the rights to distribute all the Bond movies, but the new deal solidifies the company's oversight of Bond's big-screen future.

TECH

NOT TECH

**w**: weights

| feature | weight |
|---------|--------|
| Amazon | .43 |
| acquired | 0.025 |
| behemoth | 0.008 |
| Bond | -0.0001 |
| … | |

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

| feature $f_i(x)$ | value |
|------------------|-------|
| Amazon | 1 |
| acquired | 1 |
| behemoth | 1 |
| Bond | 2 |
| … | |
| sniffle | 0 |
| … | |

# Maxent Modeling

$$p(\ \textsc{Tech}\ |\ \boxed{\begin{array}{l}\text{Amazon acquired MGM in 2022, taking} \\ \text{over a sprawling library that includes} \\ \text{more than 4,000 feature films and} \\ \text{17,000 television shows. The tech} \\ \text{behemoth also earned the rights to} \\ \text{distribute all the Bond movies, but the} \\ \text{new deal solidifies the company's} \\ \text{oversight of Bond's big-screen future.}\end{array}}\ ) \propto$$

$$\exp(\ \begin{array}{l}\text{weight}_{1,\text{ Tech}} * \text{applies}_1(\text{📄})\ \ + \\ \text{weight}_{2,\text{ Tech}} * \text{applies}_2(\text{📄})\ \ + \\ \text{weight}_{3,\text{ Tech}} * \text{applies}_3(\text{📄})\ \ + \\ \qquad\qquad\qquad ...\end{array}\ ))$$

K different
weights…

for K different
features

multiplied and then summed

# Maxent Modeling

$$p(\ \text{TECH}\ |\ \boxed{\begin{array}{l}\text{Amazon acquired MGM in 2022, taking}\\\text{over a sprawling library that includes}\\\text{more than 4,000 feature films and}\\\text{17,000 television shows. The tech}\\\text{behemoth also earned the rights to}\\\text{distribute all the Bond movies, but the}\\\text{new deal solidifies the company's}\\\text{oversight of Bond's big-screen future.}\end{array}}\ )\propto$$

$$\exp(\ \theta^{T}_{\text{TECH}}\ f(\boxed{})\ )\qquad \begin{bmatrix}.31\ -.5\ .1\ .002\ .522\ \dots\end{bmatrix} \times \begin{bmatrix}1\\1\\1\\2\\0\\\dots\end{bmatrix}$$

K different weights…  for K different features  multiplied and then summed

# Maxent Classifier, schematically

# Maxent Modeling

Amazon acquired MGM in 2022, taking over a sprawling library that includes more than 4,000 feature films and 17,000 television shows. The tech behemoth also earned the rights to distribute all the Bond movies, but the new deal solidifies the company's oversight of Bond's big-screen future.

$$p(\ \text{T\scriptsize ECH}\ |\ \boxed{\ }\ ) \propto$$

$$\frac{1}{Z}\exp(\ \theta_{\text{T\scriptsize ECH}}^{T}\,f(\ \text{📄}\ )\ )$$

# Normalization for Classification

$$\text{Z} =$$

$$\sum_{\text{label } j} \exp\left( \theta_j^T f(\text{📄}) \right)$$

$$p(y \mid x) \propto \exp(\theta_y^T f(x))$$

*classify doc x with label y in one go*

# Normalization for Classification (long form)

$$Z =$$

$$\sum_{\text{label } j} \exp(\quad \begin{aligned} &\text{weight}_{1,\,j} * \text{applies}_1(\text{📄}) \quad + \\ &\text{weight}_{2,\,j} * \text{applies}_2(\text{📄}) \quad + \quad ) \\ &\text{weight}_{3,\,j} * \text{applies}_3(\text{📄}) \quad + \\ &\qquad\qquad \ldots \end{aligned}$$

$$p(y \mid x) \propto \exp(\theta_y^T f(x))$$

*classify doc x with label y in one go*

# Multi-label Maxent Classifier, schematically



$f(x)$          $y$

$\theta_{\text{entailed}}$

$\theta_{\text{contra}}$

$\theta_{\text{neutral}}$

$p(y = \text{entailed}|x) \propto$
$\exp(\theta_{\text{entailed}} f(x))$

$p(y = \text{contra}|x) \propto$
$\exp(\theta_{\text{contra}} f(x))$

$p(y = \text{neutral}|x) \propto$
$\exp(\theta_{\text{neutral}} f(x))$

output:
$i$ = argmax score$_i$
class $i$

# Final Equation for Logistic Regression

$$p(y \mid x) = \frac{\exp(\theta_y^T f(x))}{\sum_{y'} \exp(\theta_{y'}^T f(x))}$$

$$p(Y \mid x) = \text{softmax}(\theta f(x))$$

# Maxent Models for Classification: Discriminatively or Generatively Trained

Directly model the posterior

$$p(Y \mid X) = \mathbf{maxent}(X; Y)$$

**Discriminatively** trained classifier

---------------------------------

Model the posterior with Bayes rule

$$p(Y \mid X) \propto p(X \mid Y)p(Y)$$

**Generatively** trained classifier with maxent-based language model

# Bayes' Rule

Likelihood          Prior

$$P(Y|X) = \frac{\overbrace{P(X|Y)}^{} \cdot \overbrace{P(Y)}^{}}{P(X)}$$

Posterior

It's harder to model P(Y|X) directly
since it might be that we only see
that set of features once!

# Bayes' Rule → Naïve Bayes Assumption

Bayes
$$\hat{c} = \underset{c \in C}{\mathrm{argmax}}\, P(c|d) = \underset{c \in C}{\mathrm{argmax}}\, \frac{P(d|c) \cdot P(c)}{P(d)}$$

$$\hat{c} = \underset{c \in C}{\mathrm{argmax}}\, P(c|d) = \underset{c \in C}{\mathrm{argmax}}\, \frac{P(d|c) \cdot P(c)}{\cancel{P(d)}}$$

We can make this assumption because P(d) stays the same regardless of the class!

Naïve
Bayes
$$\hat{c} = \underset{c \in C}{\mathrm{argmax}}\, P(c|d) \approx \underset{c \in C}{\mathrm{argmax}}\, P(d|c) \cdot P(c)$$

# Bayes' Rule → Naïve Bayes Assumption

Bayes

$$\hat{c} = \underset{c \in C}{\mathrm{argmax}}\, P(c|d) = \underset{c \in C}{\mathrm{argmax}}\, \frac{P(d|c)\cdot P(c)}{P(d)}$$

Naïve
Bayes

$$\hat{c} = \underset{c \in C}{\mathrm{argmax}}\, P(c|d) \approx \underset{c \in C}{\mathrm{argmax}}\, P(d|c) \cdot P(c)$$

Naïve bayes is **generative** because we are sort of assuming this is how the data point is generated: pick a class $c$ and then generate the words by sampling from P(d|c)

*SLP 4.1*

# Modeling

Classification/
Text Processing

$$P(y \mid x)$$

Language
Model (LM) /
Generation

$$P(w_t \mid w_{t-1}, w_{t-2} \dots)$$

A language model is used to **generate** the next word(s) given a history of words.

# Language Models

Maximum likelihood (MLE): simple counting

Other count-based models
◦ E.g., Laplace smoothing (add-1, add-λ)

Maxent n-gram models ← **Featureful LMs**

Neural n-gram models ← **Feedforward LMs**

Recurrent/autoregressive NNs ← **Precursor to modern LMs**

Transformers ← **Modern LMs**

# Language Models

**Maximum likelihood (MLE): simple counting**

Other count-based models
- E.g., Laplace smoothing (add-1, add-λ)

Maxent n-gram models     ←————————— Featureful LMs

Neural n-gram models     ←————————— Feedforward LMs

Recurrent/autoregressive NNs   ←——————— Precursor to modern LMs

Transformers     ←—————————————— Modern LMs

# Key Idea: Probability Chain Rule

$$p(w_1, w_2, \ldots, w_S) =$$
$$p(w_1)p(w_2 \mid w_1)p(w_3 \mid w_1, w_2) \cdots p(w_S \mid w_1, \ldots, w_{S-1}) =$$
$$\prod_i^S p(w_i \mid w_1, \ldots, w_{i-1})$$

# N-Grams

Maintaining an entire *joint* inventory over sentences could be too much to ask

Store "smaller" pieces?

*p(Colorless green ideas sleep furiously) =*
*p(Colorless) \**
*p(green | Colorless) \**
*p(ideas | Colorless green) \**
*p(sleep | Colorless green ideas) \**
*p(furiously | Colorless green ideas sleep)*

# N-Grams

*p(furiously | Colorless green ideas sleep)*

How much does "Colorless" influence the choice of "furiously?"

Remove history and contextual info

p(furiously | Colorless green ideas sleep) **≈**

p(furiously | ideas sleep)

# N-Grams

p(Colorless green ideas sleep furiously) =
p(Colorless) *
p(green | Colorless) *
p(ideas | Colorless green) *
p(sleep | ~~Colorless~~ green ideas) *
p(furiously | ~~Colorless green~~ ideas sleep)

# Trigrams

p(Colorless green ideas sleep furiously) =
p(Colorless) *
p(green | Colorless) *
p(ideas | Colorless green) *
p(sleep | green ideas) *
p(furiously | ideas sleep)

# Trigrams

p(Colorless green ideas sleep furiously) =
p(Colorless) *
p(green | Colorless) *
p(ideas | Colorless green) *
p(sleep | green ideas) *
p(furiously | ideas sleep)

# Trigrams

p(Colorless green ideas sleep furiously) =
p(Colorless | *<BOS> <BOS>*) *
p(green | *<BOS>* Colorless) *
p(ideas | Colorless green) *
p(sleep | green ideas) *
p(furiously | ideas sleep)

*Consistent notation*: Pad the left with <BOS> (beginning of sentence) symbols

# Trigrams

p(Colorless green ideas sleep furiously) =
p(Colorless | *<BOS> <BOS>*) *
p(green | *<BOS>* Colorless) *
p(ideas | Colorless green) *
p(sleep | green ideas) *
p(furiously | ideas sleep) *
p(*<EOS>* | sleep furiously)

*Consistent notation*: Pad the left with <BOS> (beginning of sentence) symbols
*Fully proper distribution*: Pad the right with a single <EOS> symbol

# N-Gram Terminology

| n | Commonly called | History Size (Markov order) | Example |
|---|---|---|---|
| 1 | unigram | 0 | p(furiously) |
| 2 | bigram | 1 | p(furiously \| sleep) |
| 3 | trigram (3-gram) | 2 | p(furiously \| ideas sleep) |
| 4 | 4-gram | 3 | p(furiously \| green ideas sleep) |
| n | n-gram | n-1 | $p(w_i \mid w_{i-n+1} \ldots w_{i-1})$ |

# Count-Based N-Grams (Unigrams)

word type

$$p(z) = \frac{count(z)}{\sum_v count(v)}$$

word type

# Count-Based N-Grams (Unigrams)

word type

$$p(\text{z}) \propto \frac{count(\text{z})}{W}$$

number of tokens observed

# Count-Based N-Grams (Unigrams)

The film got a great opening and the film went on to become a hit .

| Word (Type) z | Raw Count count(z) | Normalization | Probability p(z) |
|:---:|:---:|:---:|:---:|
| The | 1 | | 1/16 |
| film | 2 | | 1/8 |
| got | 1 | | 1/16 |
| a | 2 | | 1/8 |
| great | 1 | | 1/16 |
| opening | 1 | | 1/16 |
| and | 1 | | 1/16 |
| the | 1 | 16 | 1/16 |
| went | 1 | | 1/16 |
| on | 1 | | 1/16 |
| to | 1 | | 1/16 |
| become | 1 | | 1/16 |
| hit | 1 | | 1/16 |
| . | 1 | | 1/16 |

# Count-Based N-Grams (Trigrams)

order matters in
conditioning

order matters in
*count*

$$p(\text{z}|\text{x},\text{y}) \propto count(\text{x},\text{y},\text{z})$$

Count of the
sequence of items
"x y z"

# Count-Based N-Grams (Trigrams)

order matters in
conditioning

order matters in
*count*

$$p(\text{z}|\text{x},\text{y}) \propto count(\text{x},\text{y},\text{z})$$

count(x, y, z) ≠ count(x, z, y) ≠ count(y, x, z) ≠ …

# Count-Based N-Grams (Trigrams)

$$p(z|x,y) \propto count(x,y,z)$$

$$= \frac{count(x,y,z)}{\sum_v count(x,y,v)}$$

# Count-Based N-Grams (Lowercased Trigrams)

the film got a great opening and the film went on to become a hit .

| Context: x y | Word (Type): z | Raw Count | Normalization | Probability: p(z \| x y) |
|---|---|---|---|---|
| the film | the | 0 | | 0/2 |
| the film | film | 0 | 2 | 0/2 |
| the film | got | 1 | | 1/2 |
| the film | went | 1 | | 1/2 |
| ... | | | | |
| a great | great | 0 | | 0/1 |
| a great | opening | 1 | 1 | 1/1 |
| a great | and | 0 | | 0/1 |
| a great | the | 0 | | 0/1 |
| ... | | | | |

# Language Models

Maximum likelihood (MLE): simple counting

Other count-based models

- ◦ E.g., **Laplace smoothing (add-1, add-λ)**

Maxent n-gram models ← Featureful LMs

Neural n-gram models ← Feedforward LMs

Recurrent/autoregressive NNs ← Precursor to modern LMs

Transformers ← Modern LMs

# 0s Are Not Your (Language Model's) Friend

$$p(\text{item}) \propto count(\text{item}) = 0 \rightarrow$$
$$p(\text{item}) = 0$$

0 probability → item is *impossible*

     0s annihilate: x*y*z*0 = 0

Language is creative:

     new words keep appearing

     existing words could appear in known contexts

How much do you trust your data?

# Add-λ estimation

Other names: Laplace smoothing, Lidstone smoothing

Pretend we saw each word λ more times than we did

$$p(z) \propto count(z) + \lambda$$

Add λ to all the counts

# Add-λ estimation

Other names: Laplace smoothing, Lidstone smoothing

Pretend we saw each word λ more times than we did

Add λ to all the counts

$$p(\mathrm{z}) \propto count(\mathrm{z}) + \lambda$$

$$= \frac{count(\mathrm{z}) + \lambda}{\sum_v (count(\mathrm{v}) + \lambda)}$$

# Add-λ estimation

Other names: Laplace smoothing, Lidstone smoothing

Pretend we saw each word λ more times than we did

Add λ to all the counts

$$p(\text{z}) \propto count(\text{z}) + \lambda$$

$$= \frac{count(\text{z}) + \lambda}{W + V\lambda}$$

# tokens     # types

# What are the tri-grams for "The film , a hit !"

| Trigrams | MLE p(trigram) | UNK-ed trigrams | Smoothed p(trigram) |
|---|---|---|---|
| <BOS> <BOS> The | 1 | <BOS> <BOS> The | 2/17 |
| <BOS> The film | 1 | <BOS> The film | 2/17 |
| The film , | 0 | The film <UNK> | 1/17 |
| film , a | 0 | film <UNK> a | 1/16 |
| , a hit | 0 | <UNK> a hit | 1/16 |
| a hit ! | 0 | a hit <UNK> | 1/17 |
| hit ! <EOS> | 0 | hit <UNK> <EOS> | 1/16 |

# Language Models

Maximum likelihood (MLE): simple counting

Other count-based models
◦ E.g., Laplace smoothing (add-1, add-λ)

**Maxent n-gram models** ⟵ Featureful LMs

Neural n-gram models ⟵ Feedforward LMs

Recurrent/autoregressive NNs ⟵ Precursor to modern LMs

Transformers ⟵ Modern LMs

# Text Generation
## as *Classification Problem*?

I could eat so many delicious _____

I could eat so many juicy _____

| Types | Probability |
|-----------|-------------|
| apples | .03 |
| sandwiches | .02 |
| pineapples | .004 |
| houses | .00002 |
| … | … |

?

# Maxent Models as Featureful n-gram Language Models

p(Colorless green ideas sleep furiously | Label) =
p(Colorless | Label, <BOS>) * … * p(<EOS> | Label , furiously)

Model each n-gram term with a maxent model

$$p(x_i \mid y, x_{i-N+1:i-1}) = \text{maxent}(y, x_{i-N+1:i-1}, x_i)$$

*generatively trained:*
*learn to model (class-specific) language*

# Language Model with Maxent n-grams

label

$$p_n(\text{📄}|y) = \prod_{i=1}^{M} \text{maxent}(y, x_{i-n+1:i-1}, x_i)$$

n-gram

Iterate through all possible output vocab types $x'$---just like in count-based LMs

$$= \prod_{i=1}^{M} \frac{\exp(\theta_{x_i}^T f(y, x_{i-n+1:i-1}))}{\sum_{x'} \exp(\theta_{x'}^T f(y, x_{i-n+1:i-1}))}$$

# Count-based Language Models

*given some context...*

$w_{i-3}$    $w_{i-2}$    $w_{i-1}$

*compute beliefs about what is likely...*

$$p(w_i|\, w_{i-3}, w_{i-2}, w_{i-1}) \propto count(w_{i-3}, w_{i-2}, w_{i-1}, w_i)$$

*predict the next word*

$w_i$

# Maxent Language Models

*given some context...*

$w_{i-3}$    $w_{i-2}$    $w_{i-1}$

*compute beliefs about what is likely...*

$$p(w_i|\, w_{i-3}, w_{i-2}, w_{i-1}) = \mathrm{softmax}(\theta_{w_i} \cdot f(w_{i-3}, w_{i-2}, w_{i-1}))$$

*predict the next word*

$w_i$

# 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 θ

**output**

**"Gold" (correct) labels**

**Objective Function**

θ

score

Objective Function

*give feedback to the model*

# Types of Learning

**SUPERVISED LEARNING**

**UNSUPERVISED LEARNING**

cat

dog

dog

hamster

# Types of Learning

## SUPERVISED LEARNING

Data has feedback (labels)

Data consists of input-output pairs

Learn mapping from input to output

*Examples:*
- Dataset classification
- How likely is it that this person will get into a car accident?

## UNSUPERVISED LEARNING

No explicit feedback in data

Learn patterns directly from data

*Examples*:
- Clustering
- Do these people fall under multiple groups?

# What are some other examples of these?

**SUPERVISED LEARNING**

- Machine translation
- Object segmentation (vision)
- Document classification

**UNSUPERVISED LEARNING**

- Clustering
- Language modeling

# Types of Algorithms

|  | **Supervised Learning** | **Unsupervised Learning** |
|---|---|---|
| **Discrete** | classification or categorization | clustering |
| **Continuous** | regression | dimensionality reduction |

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



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 Function

θ

score

Objective Function

*give feedback to the model*

$$p_\theta(y \mid x)$$ probabilistic model

$$F(\theta; x, y)$$ **objective**

# Primary Objective: Likelihood

Goal: *maximize* the score your model gives to the training data it observes

This is called the **likelihood of your data**

In classification, this is p(label | 📄)

For language modeling, this is p(word | history of words)

# Objective = Full Likelihood? (Classification)

Our goal probability

Our maxent equation

$$\prod_i p_\theta(y_i|x_i) \propto \prod_i \exp(\theta_{y_i}^T f(x_i))$$

These values can have very small magnitude ➜ underflow

Differentiating this product could be a pain

# Logarithms

(0, 1] ➔ (-∞, 0]

Products ➔ Sums

    log(ab) = log(a) + log(b)

    log(a/b) = log(a) – log(b)

Inverse of exp

    log(exp(x)) = x

How might you find the log of this?

$$\prod_i p_\theta(y_i | x_i)$$

# Log-Likelihood (Classification)

Wide range of (negative) numbers
Sums are more stable

$$\log \prod_i p_\theta(y_i|x_i) = \sum_i \log p_\theta(y_i|x_i)$$

*Products* ➔ *Sums*

*log(ab) = log(a) + log(b)*

*log(a/b) = log(a) − log(b)*

# *Maximize* Log-Likelihood (Classification)

$$\frac{\exp(\theta_y^T f(x))}{\sum_{y'} \exp(\theta_{y'}^T f(x))}$$

$$\log \prod_i p_\theta(y_i|x_i) = \sum_i \log p_\theta(y_i|x_i)$$

Differentiating this becomes nicer (even though Z depends on θ)

*Inverse of exp*
*log(exp(x)) = x*

$$= \sum_i \theta_{y_i}^T f(x_i) - \log Z(x_i)$$

# Log-Likelihood (Classification)

Wide range of (negative) numbers
Sums are more stable

$$\log \prod_i p_\theta(y_i|x_i) = \sum_i \log p_\theta(y_i|x_i)$$

$$= \sum_i \theta_{y_i}^T f(x_i) - \log Z(x_i)$$

$$= F(\theta)$$

# Equivalent Version 2:
## *Minimize* Cross Entropy Loss

loss uses y (random variable), or model's probabilities $\ell^{\mathrm{xent}}(\overrightarrow{y^*}, p(y|x))$

**Cross entropy:**
How much $\hat{y}$ differs from the true $y$

$$\ell^{\mathrm{xent}}\left(\overrightarrow{y^*}, y\right) = -\sum_{k} \overrightarrow{y^*}[k] \log p(y = k|x)$$

index of "1" indicates correct value

$$\begin{pmatrix} 0 \\ 0 \\ \dots \\ 1 \\ \dots \\ 0 \end{pmatrix}$$

one-hot vector

# Classification Log-likelihood $\cong$ Cross Entropy Loss

**Log Likelihood**

$$F(\theta) = \sum_i \theta_{y_i}^T f(x_i) - \log Z(x_i)$$

objective is concave

**Cross Entropy Loss**

$$\ell^{\text{xent}}(\overrightarrow{y^*}, y) = -\sum_k \overrightarrow{y^*}[k] \log p(y = k|x)$$

objective is convex

# Preventing Extreme Values

Likelihood on its own can lead to overfitting and/or extreme values in the probability computation

$$F(\theta) = \sum_i \theta_{y_i}^T f(x_i) - \log Z(x_i)$$

Learn the parameters based on some (fixed) data/examples

# Regularization: Preventing Extreme Values

$$F(\theta) = \left( \sum_i \theta_{y_i}^T f(x_i) - \log Z(x_i) \right) - R(\theta)$$

With fixed/predefined features, the values of $\theta$ determine how "good" or "bad" our objective learning is

- Augment the objective with a **regularizer**
- This regularizer places an inductive bias (or, prior) on the general "shape" and values of $\theta$

# (Squared) L2 Regularization

$$R(\theta) = \|\theta\|_2^2 = \sum_k \theta_k^2$$

# Regularization:
# Preventing Extreme Values

$$F(\theta) = \left( \sum_i \theta_{y_i}^T f(x_i) - \log Z(x_i) \right) - \sum_k \theta_k^2$$

# ML/NLP Framework for Learning

# Optimizing F(θ)



F(θ)

θ

# Optimizing F(θ)

# Optimizing F(θ)

# What if you can't find the roots? Follow the derivative



F'(θ)
*derivative of F wrt θ*

F(θ)

θ

θ*

# What if you can't find the roots? Follow the derivative

Set t = 0
Pick a starting value $\theta_t$
Until converged:
1. Get value $z_t = F(\theta_t)$

$F'(\theta)$
*derivative of F wrt θ*

$z_0$

$F(\theta)$

$\theta_0$

$\theta*$

$\theta$

# What if you can't find the roots? Follow the derivative

Set t = 0
Pick a starting value $\theta_t$
Until converged:
1. Get value $z_t = F(\theta_t)$
2. Get derivative $g_t = F'(\theta_t)$

F'($\theta$)
*derivative of F wrt $\theta$*

$z_0$

F($\theta$)

$g_0$

$\theta_0$

$\theta^*$

$\theta$

# What if you can't find the roots? Follow the derivative

Set t = 0

Pick a starting value $\theta_t$

Until converged:

1. Get value $z_t = F(\theta_t)$
2. Get derivative $g_t = F'(\theta_t)$
3. Get scaling factor (learning rate) $\rho_t$
4. Set $\theta_{t+1} = \theta_t + \rho_t * g_t$
5. Set t += 1

$F'(\theta)$
*derivative of F wrt θ*

$z_0$

$F(\theta)$

$g_0$

$\theta$

$\theta_0 \rightarrow \theta_1$

$\theta^*$

# What if you can't find the roots? Follow the derivative

Set t = 0
Pick a starting value $\theta_t$
Until converged:
1. Get value $z_t = F(\theta_t)$
2. Get derivative $g_t = F'(\theta_t)$
3. Get scaling factor (learning rate) $\rho_t$
4. Set $\theta_{t+1} = \theta_t + \rho_t * g_t$
5. Set t += 1

$F'(\theta)$
*derivative of F wrt $\theta$*

$z_1$

$z_0$

$F(\theta)$

$g_0$

$g_1$

$\theta$

$\theta_0 \rightarrow \theta_1 \rightarrow \theta_2 \quad \theta^*$

# What if you can't find the roots? Follow the derivative

Set t = 0
Pick a starting value $\theta_t$
Until converged:
1. Get value $z_t = F(\theta_t)$
2. Get derivative $g_t = F'(\theta_t)$
3. Get scaling factor (learning rate) $\rho_t$
4. Set $\theta_{t+1} = \theta_t + \rho_t * g_t$
5. Set t += 1

$F'(\theta)$
*derivative of F wrt θ*

$F(\theta)$

$z_0$, $z_1$, $z_2$, $z_3$

$g_0$, $g_1$, $g_2$

$\theta_0 \rightarrow \theta_1 \rightarrow \theta_2 \rightarrow \theta_3 \; \theta^*$

$\theta$

# Gradient = Multi-variable derivative

K-dimensional input

$$\nabla_\theta F\left(\theta\right) = \left(\frac{\partial F}{\partial \theta_1}, \frac{\partial F}{\partial \theta_2}, \ldots, \frac{\partial F}{\partial \theta_K}\right)$$

K-dimensional output

# Gradient Ascent

# ML/NLP Framework for Prediction

# Getting Labels from the Classifier

Given X, our classifier produces a score for each possible label

$$p(\bullet|X) \text{ vs. } p(\circ|X)$$

Can turn a probability ("regression") model into a classification model

$$\text{best label} = \arg\max_{\text{label}} P(\text{label}|\text{example})$$

# Example of `argmax`

Amazon acquired MGM in 2022, taking over a sprawling library that includes more than 4,000 feature films and 17,000 television shows. The tech behemoth also earned the rights to distribute all the Bond movies, but the new deal solidifies the company's oversight of Bond's big-screen future.

| | |
|---|---|
| POLITICS | .002 |
| **MOVIES** | **.48** |
| SPORTS | .0001 |
| TECH | .39 |
| HEALTH | .0001 |
| FINANCE | .05 |
| … | |

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

# Determining how good a model is: Baselines

# Central Question: How Well Are We Doing?

**Classification**
**Regression**
**Clustering**

*the **task**: what kind of problem are you solving?*

- Precision, Recall, F1
- Accuracy
- Log-loss
- ROC-AUC
- …

- (Root) Mean Square Error
- Mean Absolute Error
- …

- Mutual Information
- V-score
- …

This does not have to be the same thing as the loss function you optimize

# Evaluating Classification

# Classification Evaluation:
# the 2-by-2 contingency table

Assumption 1: There are two classes/labels

Assumption 2: ⬤ is the "positive" label

# Classification Evaluation: the 2-by-2 contingency table

| *What label does our system predict? (↓)* | What is the actual label? | |
| --- | --- | --- |
| | Actual Target Class ("●") | Not Target Class ("○") |
| **Selected/ Guessed ("●")** | True Positive (TP) ● Actual ● Guessed | False Positive (FP) ○ Actual ● Guessed |
| **Not selected/ not guessed ("○")** | False Negative (FN) ● Actual ○ Guessed | True Negative (TN) ○ Actual ○ Guessed |

# Contingency Table (out of table form)

Query:
Articles about dogs

Simple model classifies based on presence of "dog" or "dogs"

**FN**

Adopting an Animal at the Pound

**TP**

10 Tips for Grooming your Dog

What Foods to Hide from Your Dog During the Holidays

Why Huskies Howl

Who's a good boy?

**TN**

15 Weird Animals You Can Find in North America

**FP**



It's raining cats and dogs

When Cats Rule the World

# Contingency Table Example

**Predicted:** ⚪ 🟠 🟠 🟠 ⚪ 🟠

**Actual:** 🟠 🟠 🟠 ⚪ ⚪ ⚪

| *What label does our system predict? (↓)* | *What is the actual label?* | |
|---|---|---|
| | Actual Target Class ("🟠") | Not Target Class ("⚪") |
| **Selected/ Guessed ("🟠")** | True Positive (TP) | False Positive (FP) |
| **Not selected/ not guessed ("⚪")** | False Negative (FN) | True Negative (TN) |

# Contingency Table Example

**Predicted:** ◯ ● ● ● ◯ ●

**Actual:** ● ● ● ◯ ◯ ◯

| *What label does our system predict? (↓)* | *What is the actual label?* | |
| --- | --- | --- |
| | Actual Target Class ("●") | Not Target Class ("◯") |
| **Selected/ Guessed ("●")** | True Positive (TP) = 2 | False Positive (FP) = 2 |
| **Not selected/ not guessed ("◯")** | False Negative (FN) = 1 | True Negative (TN) = 1 |

# Classification Evaluation: Accuracy, Precision, and Recall

**Accuracy**: % of items correct

$$\frac{TP + TN}{TP + FP + FN + TN}$$

|  | **Actually Target** | **Actually Not Target** |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

# Classification Evaluation: Accuracy, Precision, and Recall

**Accuracy**: % of items correct

$$\frac{TP + TN}{TP + FP + FN + TN}$$

**Precision**: % of selected items that are correct

$$\frac{TP}{TP + FP}$$

"**Precision** measures the percentage of the items that precision the system detected (i.e., the system labeled as positive) that are in fact positive (i.e., are positive according to the human gold labels"
SLP, ch. 4

|  | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

# Classification Evaluation: Accuracy, Precision, and Recall

**Accuracy**: % of items correct

$$\frac{TP + TN}{TP + FP + FN + TN}$$

**Precision**: % of selected items that are correct

$$\frac{TP}{TP + FP}$$

**Recall**: % of correct items that are selected

$$\frac{TP}{TP + FN}$$

"**Recall** measures the percentage of items actually present in the input that were correctly identified by the system." SLP, ch. 4

| | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

# Classification Evaluation: Accuracy, Precision, and Recall

**Accuracy**: % of items correct

$$\frac{TP + TN}{TP + FP + FN + TN}$$

**Precision**: % of selected items that are correct

$$\frac{TP}{TP + FP}$$

Min: 0 ☹
Max: 1 😀

**Recall**: % of correct items that are selected

$$\frac{TP}{TP + FN}$$

| | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

# Comparing Accuracy & F1

**Accuracy**: % of items correct

$$\frac{TP + TN}{TP + FP + FN + TN}$$

$$F_1 = \frac{2*P*R}{P+R} = \frac{2*TP}{2*TP+FP+FN}$$

When would you want to use accuracy vs F1?

Accuracy works better if the dataset is <u>balanced</u>

Accuracy takes everything in consideration

F-Score is focused on TP

|  | **Actually Target** | **Actually Not Target** |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

# P/R/F in a Multi-class Setting: Micro- vs. Macro-Averaging

**Macroaveraging**: Compute performance for each class, then average.

$$\text{macroprecision} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FP}_c} = \frac{1}{C}\sum_c \text{precision}_c$$

$$\text{macrorecall} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FN}_c} = \frac{1}{C}\sum_c \text{recall}_c$$

**Microaveraging**: Collect decisions for all classes, compute contingency table, evaluate.

$$\text{microprecision} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FP}_c} \qquad\qquad \text{microrecall} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FN}_c}$$

# Macro/Micro Example



Predicted "A"   Predicted "B"   Predicted "C"   Predicted "D"

True "A"   True "B"   True "C"   True "D"

https://www.evidentlyai.com/classification-metrics/multi-class-metrics

# Macro-Average

Predicted "A" | Predicted "B" | Predicted "C" | Predicted "D"

**Class A**

Recall: **87%.**
Precision: 72%.

**Class B**

Recall: **33%.**
Precision: **20%.**

**Class C**

Recall: **90%.**
Precision: **90%.**

**Class D**

Recall: **93%.**
Precision: **100%.**

True "A" | True "B" | True "C" | True "D"

**Macro-average**

Recall = (0.87 + 0.33 + 0.9 + 0.93)/4 = **0.76**
Precision = (0.72+0.2+0.9+1)/4=**0.71**

*https://www.evidentlyai.com/classification-metrics/multi-class-metrics*

Each *instance* has equal weight

# Micro-Average

All true positives

All false negatives

All false positives

13 — True positive "A"
1 — True positive "B"
9 — True positive "C"
14 — True positive "D"

2 — False negative "A"
4 — False negative "B"
1 — False negative "C"
1 — False negative "D"

2 — False positive "B"
5 — False positive "A"
1 — False positive "C"
False positive "D" : 0

$$\text{Precision} = \frac{13 + 1 + 9 + 14}{(13 + 1 + 9 + 14) + (2 + 5 + 1 + 0)} = 0.82$$
Micro-average

$$\text{Recall} = \frac{13 + 1 + 9 + 14}{(13 + 1 + 9 + 14) + (2 + 4 + 1 + 1)} = 0.82$$
Micro-average

Total TP    Total FP    Total FN

Predicted "A"    Predicted "B"    Predicted "C"    Predicted "D"

https://www.evidentlyai.com/classification-metrics/multi-class-metrics

# But how do we compute stats for multiple classes?

We already saw how the "polarity" affects the stats we compute…

Two main approaches. Either:

1. Compute "one-vs-all" 2x2 tables. OR

2. Generalize the 2x2 tables and compute per-class TP / FP / FN based on the diagonals and off-diagonals

# 1. Compute "one-vs-all" 2x2 tables

Predicted

Actual

| Look for ● | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

| Look for ○ | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

| Look for ▭ | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

# 1. Compute "one-vs-all" 2x2 tables

Predicted   ● ○ ▭ ● ○ ▭ ● ○ ▭

Actual   ● ○ ○ ▭ ○ ▭ ● ● ●

| Look for ● | Actually Target | Actually Not Target |
|---|---|---|
| Selected/Guessed | 2 | 1 |
| Not select/not guessed | 2 | 4 |

| Look for ○ | Actually Target | Actually Not Target |
|---|---|---|
| Selected/Guessed | 2 | 1 |
| Not select/not guessed | 1 | 5 |

| Look for ▭ | Actually Target | Actually Not Target |
|---|---|---|
| Selected/Guessed | 1 | 2 |
| Not select/not guessed | 1 | 5 |

NLP REVIEW

# 2. Generalizing the 2-by-2 contingency table

|  | | Correct Value | | |
|---|---|---|---|---|
| | | 🟠 | ◯ | ▭ |
| **Guessed Value** | 🟠 | # | # | # |
| | ◯ | # | # | # |
| | ▭ | # | # | # |

This is also called a **Confusion Matrix**

# 2. Generalizing the 2-by-2 contingency table

Predicted

Actual

| | Correct Value | | |
|---|---|---|---|
| Guessed Value | 2 | 0 | 1 |
| | 1 | 2 | 0 |
| | 1 | 1 | 1 |

This is also called a **Confusion Matrix**

# 2. Generalizing the 2-by-2 contingency table

Predicted

Actual

| | | Correct Value | | |
|---|---|---|---|---|
| | | ● | ○ | ▭ |
| **Guessed Value** | ● | A  2 | B  0 | C  1 |
| | ○ | D  1 | E  2 | F  0 |
| | ▭ | G  1 | H  1 | I  1 |

How do you compute $TP$ ● ?

# 2. Generalizing the 2-by-2 contingency table

Predicted

Actual

| | Correct Value | | |
|---|---|---|---|
| | ● | ○ | ▭ |
| **Guessed Value** ● | A 2 | B 0 | C 1 |
| ○ | D 1 | E 2 | F 0 |
| ▭ | G 1 | H 1 | I 1 |

How do you compute $FN_{●}$?

# 2. Generalizing the 2-by-2 contingency table

Predicted

Actual

| | Correct Value | | |
|---|---|---|---|
| Guessed Value | 2 (A) | 0 (B) | 1 (C) |
| | 1 (D) | 2 (E) | 0 (F) |
| | 1 (G) | 1 (H) | 1 (I) |

How do you compute $FP_{\square}$ ?

# Evaluating Generation

# Evaluating Language Models

*What is "correct?"*

*What is working "well?"*

**Extrinsic**: Evaluate LM in downstream task

Test an MT, ASR, etc. system and see which LM does better

Issue: Propagate & conflate errors

**Intrinsic**: Treat LM as its own downstream task

Use perplexity (from information theory)

# Perplexity: Average "Surprisal"

Lower is better : lower perplexity ➡ less surprised



p(word type | context)

word type

Less certain ➡
More surprised ➡
Higher perplexity



p(word type | context)

word type

More certain ➡
Less surprised ➡
Lower perplexity

# Perplexity

Lower is better : lower perplexity ➔ less surprised

$$\text{perplexity} = \exp(\text{avg crossentropy})$$

# Perplexity

Lower is better : lower perplexity ➔ less surprised

$$\text{perplexity} = \exp(\frac{-1}{M}\log p(w_1, \dots, w_M))$$

# Perplexity

Lower is better : lower perplexity ➡ less surprised

*e.g., n-gram history (n-1 items)*

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^{M} \log p(w_i \mid h_i)\right)$$

# Example perplexity for trigram model

"The film , a hit !"

| Trigrams | MLE p(trigram) |
|---|---|
| <BOS> <BOS> The | 1 |
| <BOS> The film | 1 |
| The film , | 0 |
| film , a | 0 |
| , a hit | 0 |
| a hit ! | 0 |
| hit ! <EOS> | 0 |
| **Perplexity** | Infinity |

$$\text{perplexity} = \exp(\frac{-1}{M} \sum_{i=1}^{M} \log p(w_i \mid h_i))$$

# Example perplexity for trigram model

| Trigrams | MLE p(trigram) | Smoothed p(trigram) |
|---|---|---|
| <BOS> <BOS> The | 1 | 2/17 |
| <BOS> The film | 1 | 2/17 |
| The film , | 0 | 1/17 |
| film , a | 0 | 1/16 |
| , a hit | 0 | 1/16 |
| a hit ! | 0 | 1/17 |
| hit ! <EOS> | 0 | 1/16 |
| **Perplexity** | Infinity | 13.59 |

"The film , a hit !"

$$\text{perplexity} = \exp(\frac{-1}{M} \sum_{i=1}^{M} \log p(w_i \mid h_i))$$

# 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

# 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*

# Document Classification

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Language Identification

Sentiment analysis

…

# Document Classification

Assigning subject categories, topics, or genres

Spam detection

Authorship identification

Language Identification

Sentiment analysis

…

# Document Classification

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
...

features $F_1$ extracted from document $d_1$

predicted class $c_1$ from $C$

actual class $c_1$

score

Objective Function

θ

Evaluation Function

score

# 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*

# p(class | token in context)

## 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 "plant" tokens

| 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 ... |

# p(class | token in context)

## 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 *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 ... |

# 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 |

# p(class | token in context)

## 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) lεd** | ... 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)*

# 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 |

<span style="color:red">generates a whole bunch of potential cues – use data to find out which ones work best</span>

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

○ ○ ○

*slide courtesy of D. Yarowsky (modified)*

# Token Classification

Word pronunciation

Word sense disambiguation (WSD) within or across languages

Accent restoration

**Other examples?**

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

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

$\theta$

actual class $c_j$

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

score

Objective Function

Evaluation Function

score

predicted class $\widehat{c_j}$

# 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

# Part of Speech Tagging



<BOS> John saw the saw and decided to take it to the table .

NNP VBD DT NN CC VBD TO VB PRP IN DT NN PUNCT

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

<BOS> John saw

John saw the

saw the saw

Sliding window

$\theta$

actual class $c_j$

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

score

Objective Function

Evaluation Function

score

predicted class $\widehat{c}_j$

# Machine Translation: Word Alignment

# Token Classification in a Sequence

Part of speech tagging

Word alignment



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

Sliding window

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

$\theta$

score

actual class $c_j$

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

Objective Function

Evaluation Function

score

predicted class $\hat{c}_j$

# 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)

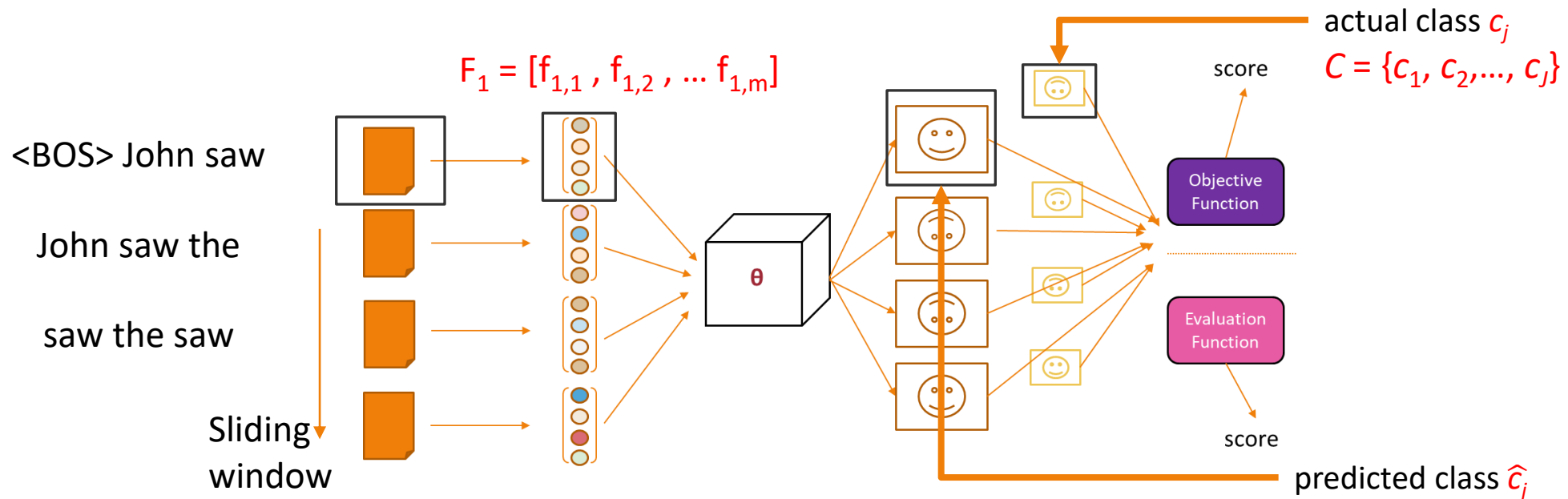6. Semantic annotation

# Named Entity Recognition

| TYPE | DESCRIPTION |
|------|-------------|
| PERSON | People, including fictional |
| NORP | Nationalities or religious or political groups |
| FACILITY | Buildings, airports, highways, bridges, etc |
| ORG | Companies, agencies, institutions, etc |
| GPE | Countries, cities, states |
| LOC | Non-GPE locations, mountain ranges, bodies of water |
| PRODUCT | Objects, vehicles, foods, etc (Not services) |
| EVENT | Named hurricanes, battles, wars, sports events, etc |
| WORK_OF_ART | Titles of books, songs, etc |
| LAW | Named documents made into laws |
| LANGUAGE | Any named language |
| DATE | Absolute or relative dates or periods. |
| TIME | Times smaller than a day |
| PERCENT | Percentage, including "%". |
| MONEY | Monetary values, including unit |
| QUANTITY | Measurements, as of weight or distance |
| ORDINAL | "first", "second", etc |
| CARDINAL | Numerals that do not fall under another type |

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.

# Example Use: Information Extraction

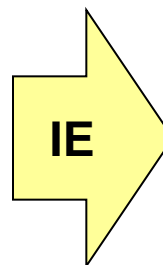**As a task:** | **Filling slots in a database from sub-segments of text.**

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…

**IE** ➡

| NAME | TITLE | ORGANIZATION |
|------|-------|--------------|
| Bill Gates | CEO | Microsoft |
| Bill Veghte | VP | Microsoft |
| Richard Stallman | founder | Free Soft.. |

Note:
IE is a task on its own but it can be an application of NER

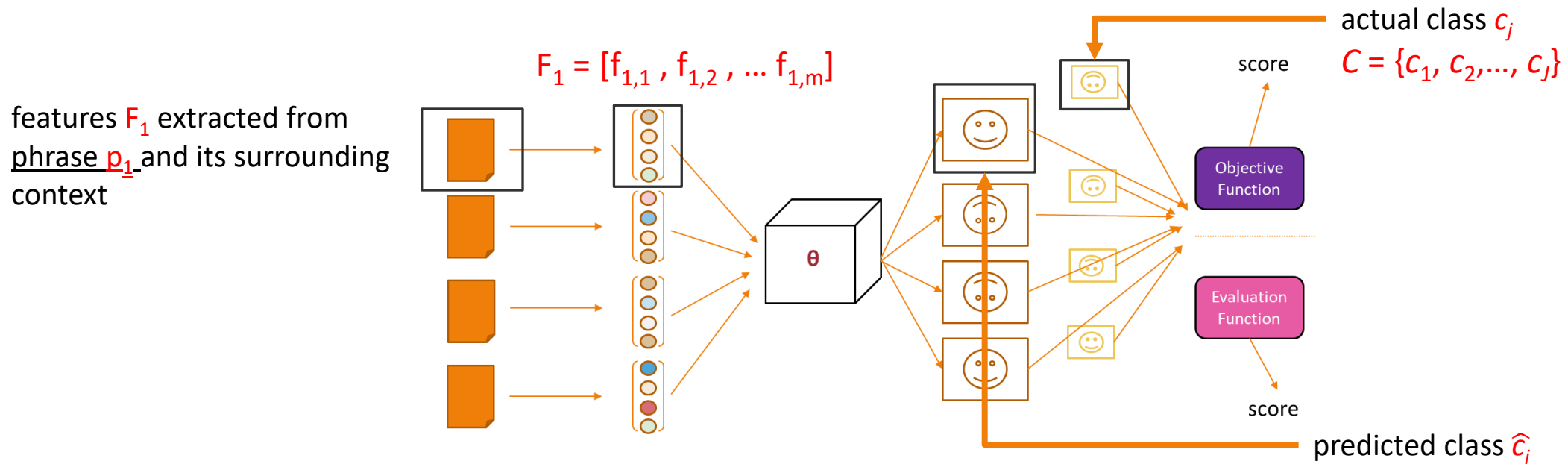*Slide from Chris Brew, adapted from slide by William Cohen*

# Chunking

Named entity recognition

Information extraction

Identifying idioms



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

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

$\theta$

score

actual class $c_j$

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

Objective Function

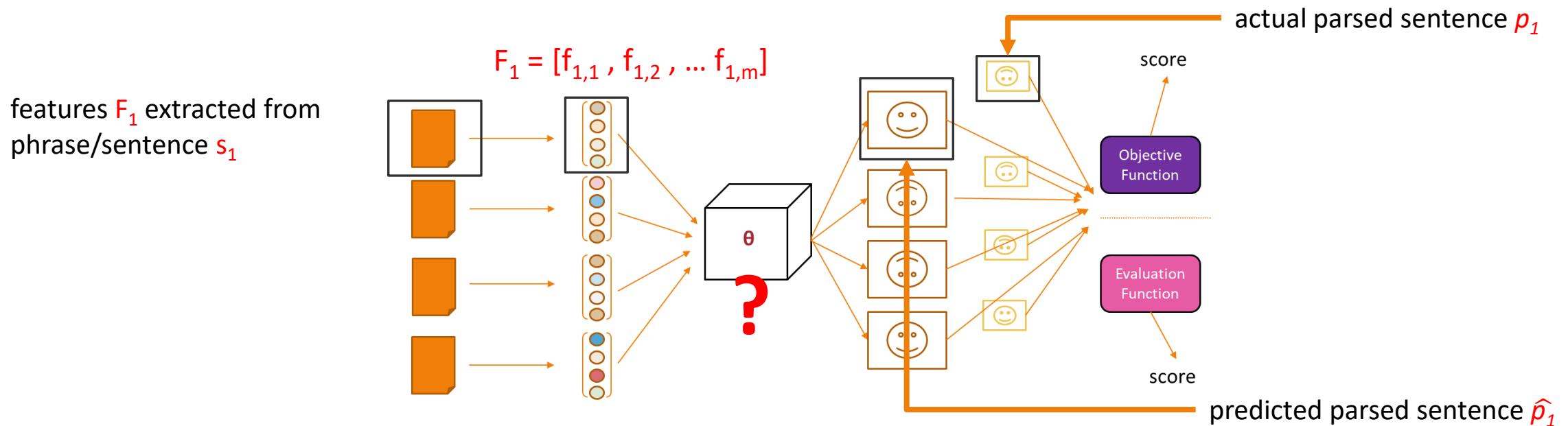Evaluation Function

score

predicted class $\hat{c}_j$

# 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

# Syntax Parsing



features $F_1$ extracted from phrase/sentence $s_1$

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

$\theta$

actual parsed sentence $p_1$

score

Objective Function

Evaluation Function

score

predicted parsed sentence $\hat{p_1}$

# Assign Structure (Parse) with a Context Free Grammar

S → NP VP          PP → P NP
NP → Det Noun     AdjP → Adj Noun
NP → Noun          VP → V NP
NP → Det AdjP     Noun → Baltimore
NP → NP PP         ...



Baltimore is a great city

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

*bracket notation*

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

*S-expression*

# 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")

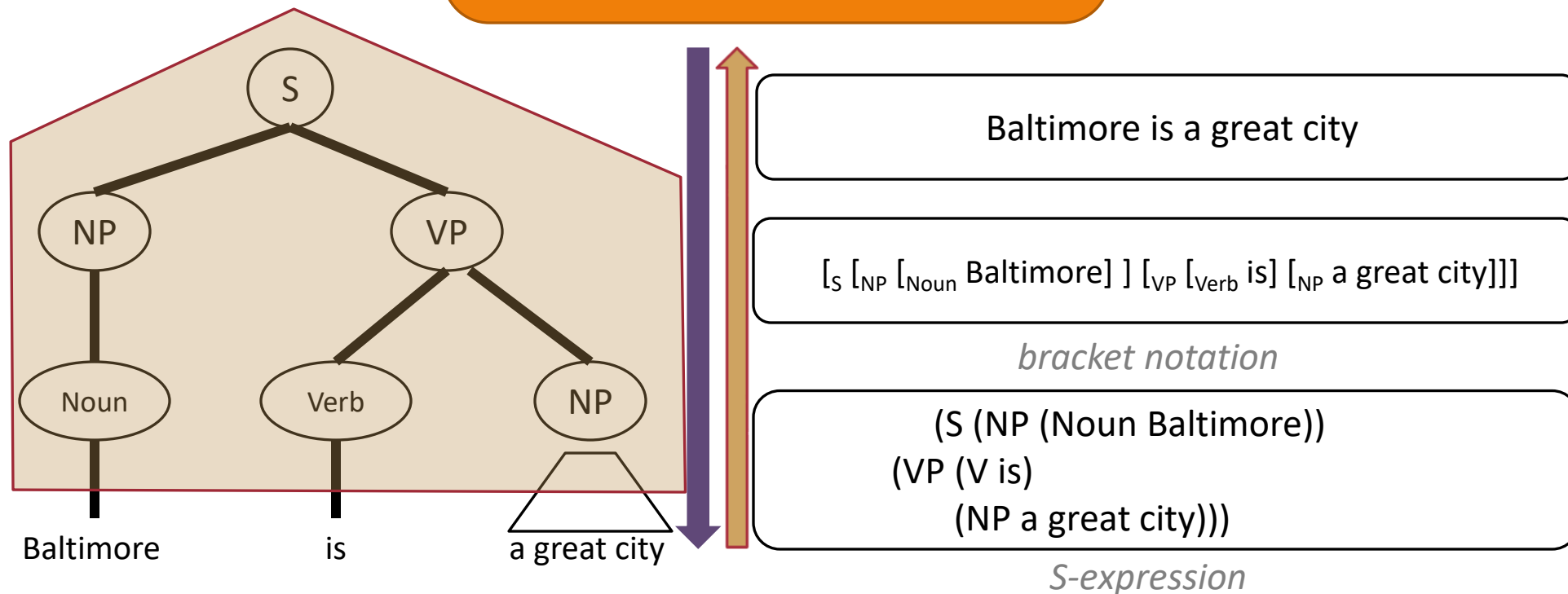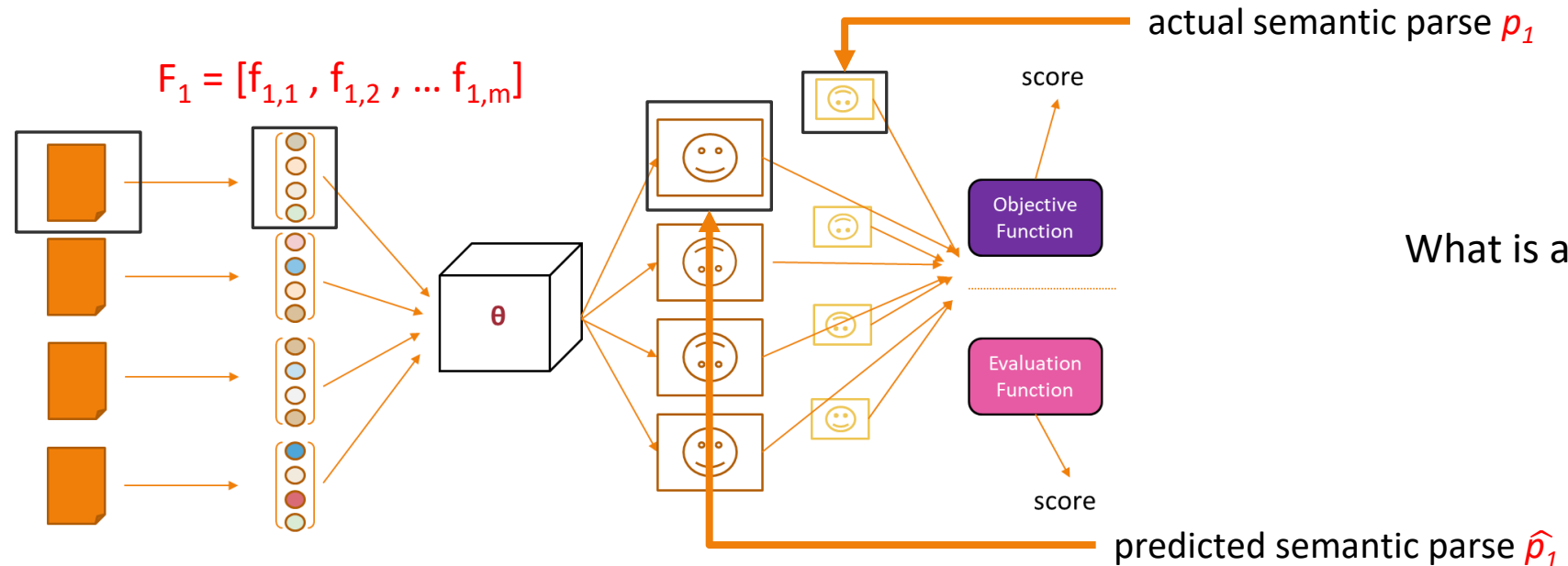5. Syntactic annotation (syntax parsing)

6. Semantic annotation

# Semantic Parsing

Semantic role labeling (SRL)



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

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

actual semantic parse $p_1$

score

Objective Function

What is a semantic parse?

Evaluation Function

score

predicted semantic parse $\hat{p}_1$

# Semantic Role Labeling (SRL)

For each <u>predicate</u> (e.g., verb)
1. find its arguments (e.g., NPs)
2. determine their **semantic roles**

> John <u>drove</u> Mary from Austin to Dallas in his Toyota Prius.

- ◦ **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

# Semantic Role Labeling (SRL)

For each <u>predicate</u> (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.

agent         patient        source    destination    instrument