

What is NLP Anyway?

- Technology to handle human language (usually text) using computers
- Aid **human-machine communication** (e.g. question answering, dialog, code generation)
- Aid **human-human communication** (e.g. machine translation, spell checking, assisted writing)
- **Analyze/understand language** (e.g. syntactic analysis, text classification, entity/relation recognition/linking)
- We now use NLP many times a day, sometimes without knowing it!

NLP can Answer Questions



You

Who is the current president of Carnegie Mellon University?



ChatGPT

I did a [quick search](#) for more information and here's what I found.

The current president of Carnegie Mellon University is Farnam Jahanian. He has been serving in this role since July 1, 2017, first as an interim president and then, from March 8, 2018, as the president [\[1\]](#).

NLP cannot Answer Questions



You

How many layers are included in gpt-3.5-turbo's architecture?



ChatGPT

GPT-3.5-Turbo, which is an optimized version of GPT-3.5 for faster responses and lower computational requirements, doesn't have a specific "layer" structure like the traditional GPT-3 models. The "Turbo" variant is designed for interactive applications, providing faster response times by optimizing certain aspects of the model's architecture and operation. However, it's important to note that the underlying model, GPT-3.5, has 175 billion parameters. In the context of GPT models, "layers" typically refer to transformer layers used in the architecture, but for GPT-3.5-Turbo, the focus is more on operational efficiency than on a specific layer count.

Retrieved Jan 5, 2024

NLP can Translate Text

緊急事態宣言から「まん延防止等重点措置」に移行した大阪府では、飲食店での酒類提供が一部解禁された。ただ、提供には府が認証する「ゴールドステッカー」の申請が必須。申請には43項目にのぼる感染対策をクリアする必要があり、飲食店からは「ハードルが高すぎる」との悲鳴が上がっている。「項目が40個以上もあって多すぎるし、ネットでの手続きも難しい。本当に、何から何までややこしい」



In Osaka Prefecture, which has transitioned from a state of emergency to "priority measures to prevent the spread of the disease," the ban on serving alcoholic beverages at restaurants has been partially lifted. However, in order to provide the service, it is necessary to apply for a "gold sticker" certified by the prefecture. To apply, it is necessary to clear 43 infection control measures, and restaurants have complained that the hurdles are too high. "There are over 40 items, which is too many, and it's difficult to complete the procedures online. It's really complicated in every way."

Front page news from Asahi Shimbun, translated by Google Jan 5, 2024

NLP cannot Translate Text

به لام 3 توپژهر به سه روکایه تی زانای سه به خوئی به به ردبوو
گریگوری پاول له شاری بالتیمور له ویلایه تی میریلاند له مانگی 3 ی
وهک سی T. rex سالی 2022 دا ئاماره یان به وه کرد که پیویسته
جوړ بناسریت.

However, three researchers, led by independent fossil scientist Gregory Paul of Baltimore, Maryland, argued in March 2022 that T. rex should be recognized as three species.

که به واتای "پاشای مارمیله کی درنده" دیت، T. rex جگه له جوړی
سه باری ئه وه 2 جوړی تریان پیشنیار کرد

In addition to the T. rex species, which means "king of ferocious lizards", they also proposed two other species.

به واتای "ئیمپراتوری مارمیله کی درنده دیت T. imperator

T. imperator means "emperor of the savage lizard

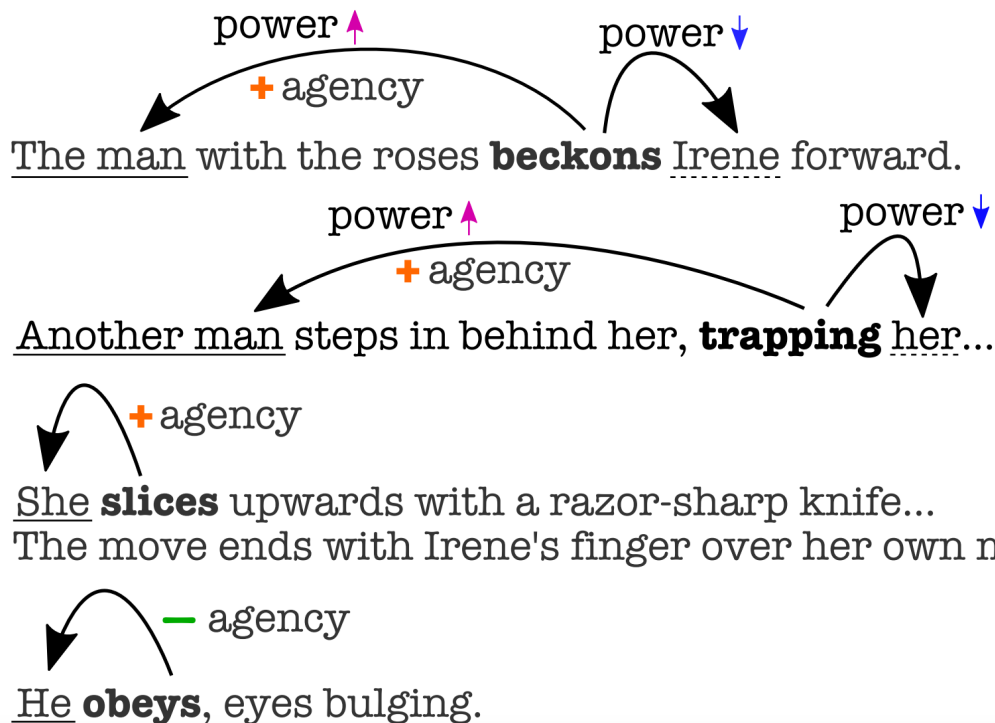
به واتای "شارنی مارمیله کی درنده T. regina

T. regina means "Queen of the ferocious snail."

Front page news from Voice of America Kurdish, translated by Google Jan 5, 2024

Language Analysis can Aid Scientific Inquiry

- e.g. *computational social science*, answering questions about society given observational data
- example: "do movie scripts portray female or male characters with more power or agency?" [Sap+ 2017]



Frame	β	gender
<i>agency</i> (AG)=+	-0.951	M**
<i>power</i> (AG>TH)	-0.468	M**
<i>agency</i> (AG)=-	0.277	F**
<i>power</i> (AG<TH)	<i>not sig.</i>	

Language Analysis Fails at Basic Tasks

First sentence of first article in NY Times Aug 29., 2021, recognized by Stanford CoreNLP

Hurricane Ida battered Louisiana on Sunday making landfall as a Category 4 storm, delivering an onslaught of harsh winds, floodwaters and power outages and threatening to assail Baton Rouge and New Orleans as one of the most devastating storms to strike the region since Hurricane Katrina.

Annotations from Stanford CoreNLP:
- CAUSE_OF_DEATH: Hurricane, storm, storms, Hurricane Katrina
- STATE_OR_PROVINCE: Louisiana
- DATE: 2021-08-29, Sunday
- NUMBER: 4.0, 4, 1.0
- CITY: Baton Rouge, New Orleans
- ORGANIZATION: Baton Rouge (highlighted in orange)

recognized by spaCy

Hurricane Ida ORG battered Louisiana GPE on Sunday DATE making landfall as a Category 4 storm, delivering an onslaught of harsh winds, floodwaters and power outages and threatening to assail Baton Rouge GPE and New Orleans GPE as one of the most devastating storms to strike the region since Hurricane Katrina.

In this Class, we Ask:

- What goes into building the state-of-the-art NLP systems that **work uncannily well at some tasks**?
- Where and why do current state-of-the-art NLP systems still **fail**?
- How can we **make appropriate improvements** and **achieve whatever we want to do with NLP**?

NLP System Building Overview

A General Framework for NLP Systems

- Create a function to map an **input** X into an **output** Y , where X and/or Y involve language.

Input X

Text

Text

Text

Text

Image

Output Y

Continuing Text

Text in Other Language

Label

Linguistic Structure

Text

Task

Language Modeling

Translation

Text Classification

Language Analysis

Image Captioning

Methods for Creating NLP Systems

- **Rules:** Manual creation of rules

```
def classify(x: str) -> str:
    sports_keywords = ["baseball", "soccer", "football", "tennis"]
    if any(keyword in x for keyword in sports_keywords):
        return "sports"
    else:
        return "other"
```

- **Prompting:** Prompting a language model w/o training

If the following sentences is about "sports"
reply "sports". Otherwise reply "other".

{X}

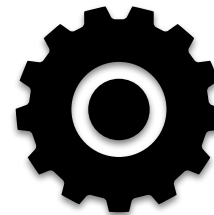
→ LM

- **Fine-tuning:** Machine learning from paired data $\langle X, Y \rangle$

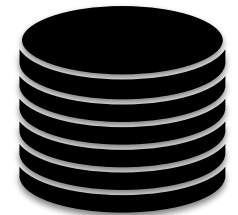
I love to play baseball.	sports
The stock price is going up.	other
He got a hat-trick yesterday.	sports
He is wearing tennis shoes.	other



Training



Model



Data Requirements for System Building

- **Rules/prompting based on intuition:**
No data needed, but also no performance guarantees
- **Rules/prompting based on spot-checks:**
A small amount of data with input X only
- **Rules/prompting with rigorous evaluation:**
Development set with input X and output Y (e.g. 200-2000 examples). Additional held-out test set also preferable.
- **Fine-tuning:**
Additional train set. More is often better — constant accuracy increase when data size doubles.

Dev

Test

Train



Let's Try to Make a Rule-based NLP System!

Example Task:

Review Sentiment Analysis

- Given a review on a reviewing web site (X), decide whether its label (Y) is positive (1), negative (-1) or neutral (0)

I hate this movie →
positive
neutral
negative

I love this movie →
positive
neutral
negative

I saw this movie →
positive
neutral
negative

A Three-step Process for Making Predictions

- **Feature extraction:** Extract the salient features for making the decision from text
- **Score calculation:** Calculate a score for one or more possibilities
- **Decision function:** Choose one of the several possibilities

Formally

- **Feature Extraction:** $f(x) = \mathbf{w} \cdot \mathbf{x}$
- **Score Calculation:** binary, multi-class
 $s = \mathbf{w} \cdot \mathbf{x}$
- **Decision:** $\hat{y} = \text{decide}(s)$

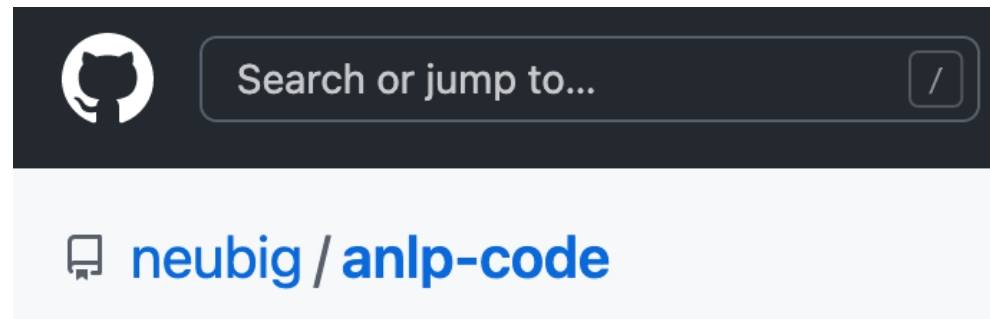
Sentiment Classification Code Walk!

<https://github.com/neubig/anlp-code/tree/main/01-simpleclassifier>

- See code for all major steps:
 1. Featurization
 2. Scoring
 3. Decision rule
 4. Accuracy calculation
 5. Error analysis

Now Let's Look at Data

<https://github.com/neubig/anlp-code>



data/sst-sentiment-text-threeclass

- Remember: don't look at "test"!

Now Let's Improve!

1. What's going wrong with my system?
→ Look at error analysis
2. Modify the system (featurization, scoring function, etc.)
3. Measure accuracy improvements, accept/reject change
4. Repeat from 1
5. Finally, when satisfied with dev accuracy, evaluate on test!

Some Difficult Cases

Low-frequency Words

The action switches between past and present , but the material link is too **tenuous** to anchor the emotional connections that **purport** to span a 125-year divide .

negative

Here 's yet another studio horror franchise **mucking** up its storyline with **glitches** casual fans could correct in their sleep .

negative

Solution?: Keep working till we get all of them? Incorporate external resources such as sentiment dictionaries?

Conjugation

An operatic , sprawling picture that 's **entertainingly** acted , **magnificently** shot and gripping enough to sustain most of its 170-minute length .

positive

It 's basically an **overlong** episode of Tales from the Crypt .

negative

Solution?: Use the root form and POS of word?

Note: Would require morphological analysis.

Negation

This one is not nearly as dreadful as expected .

positive

Serving Sara does n't serve up a whole lot of laughs .

negative

Solution?: If a negation modifies a word, disregard it.

Note: Would probably need to do syntactic analysis.

Metaphor, Analogy

Puts a human face on a land most Westerners are unfamiliar with.

positive

Green might want to hang onto that ski mask , as robbery may be the only way to pay for his next project .

negative

Has all the depth of a wading pool .

negative

Solution?: ???

Other Languages

見事に視聴者の心を掴む作品でした。

positive

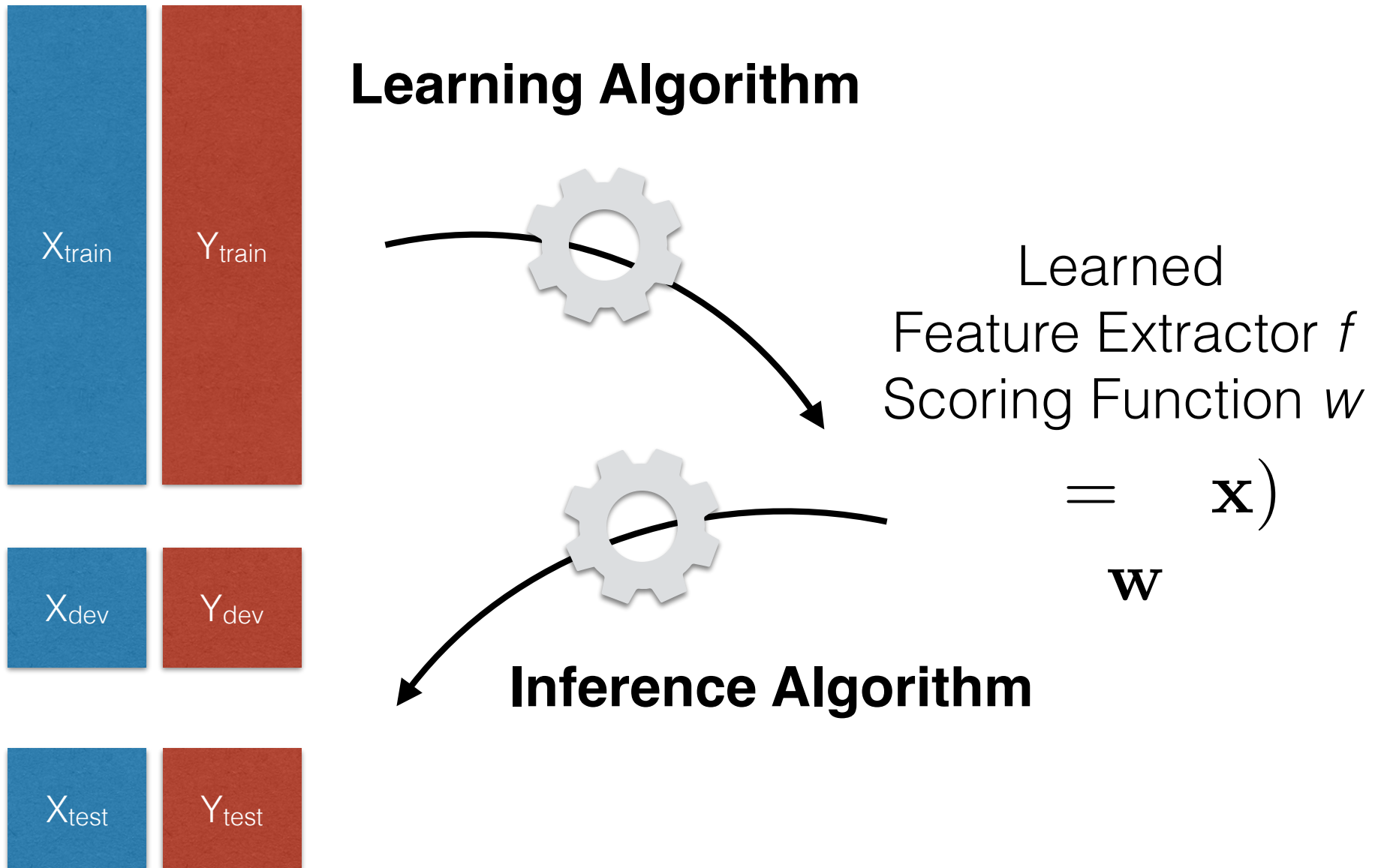
モンハンの名前がついてるからとりあえずモンハン要素を
ちょこちょこ入れればいいだる感が凄い。

negative

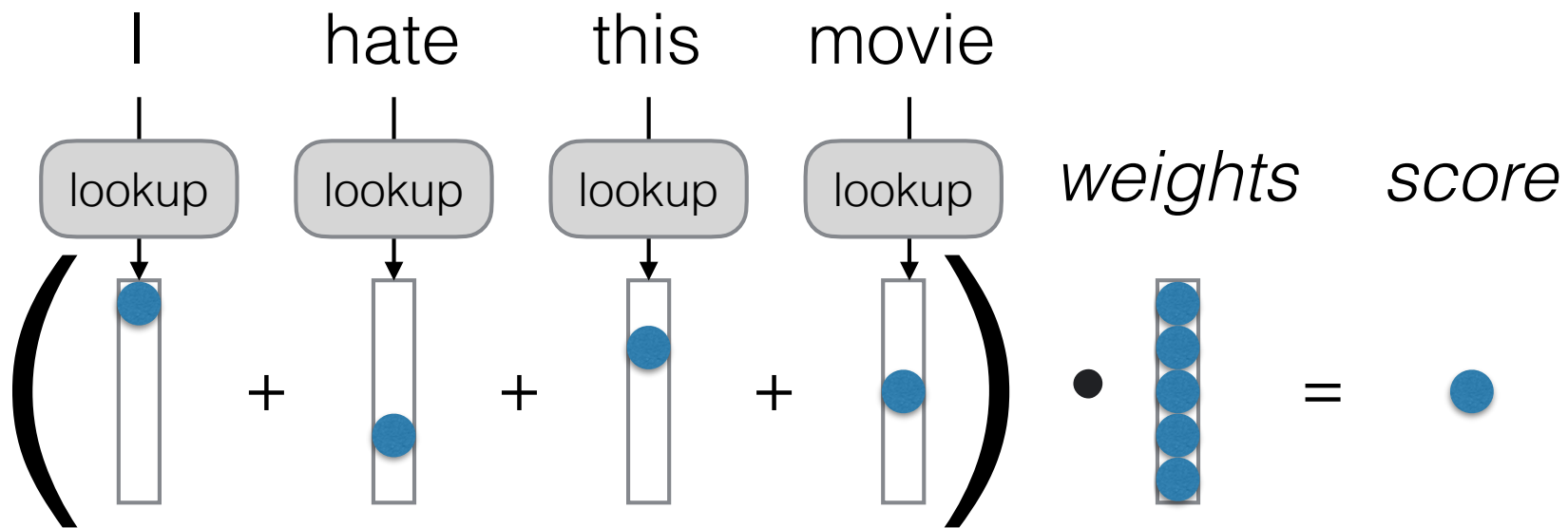
Solution?: Learn Japanese?

Machine Learning Based NLP

Machine Learning



A First Attempt: Bag of Words (BOW)



Features f are based on word identity, weights w learned

Which problems mentioned before would this solve?

What do Our Vectors Represent?

- **Binary classification:** Each word has a single scalar, positive indicating “yes” and negative indicating “no”
- **Multi-class classification:** Each word has its own 5 elements corresponding to [very good, good, neutral, bad, very bad]

Binary

love	2.4
hate	-3.5
nice	1.2
no	-0.2
dog	-0.3
...	...

Multi-class

	v. positive	positive	neutral	negative	v. negative
love	2.4	1.5	-0.5	-0.8	-1.4
hate	-3.5	-2.0	-1.0	0.4	3.2
nice	1.2	2.1	0.4	-0.1	-0.2
no	-0.2	0.3	-0.1	0.4	0.5
dog	-0.1	0.3	0.6	0.2	-0.2
...

Simple Training of BOW Models

- Use an algorithm called “structured perceptron”

```
feature_weights = {}
for x, y in data:
    # Make a prediction
    features = extract_features(x)
    predicted_y = run_classifier(features)
    # Update the weights if the prediction is wrong
    if predicted_y != y:
        for feature in features:
            feature_weights[feature] = (
                feature_weights.get(feature, 0) +
                y * features[feature]
            )
```

Full Example:

<https://github.com/neubig/anlp-code/tree/main/01-simpleclassifier>