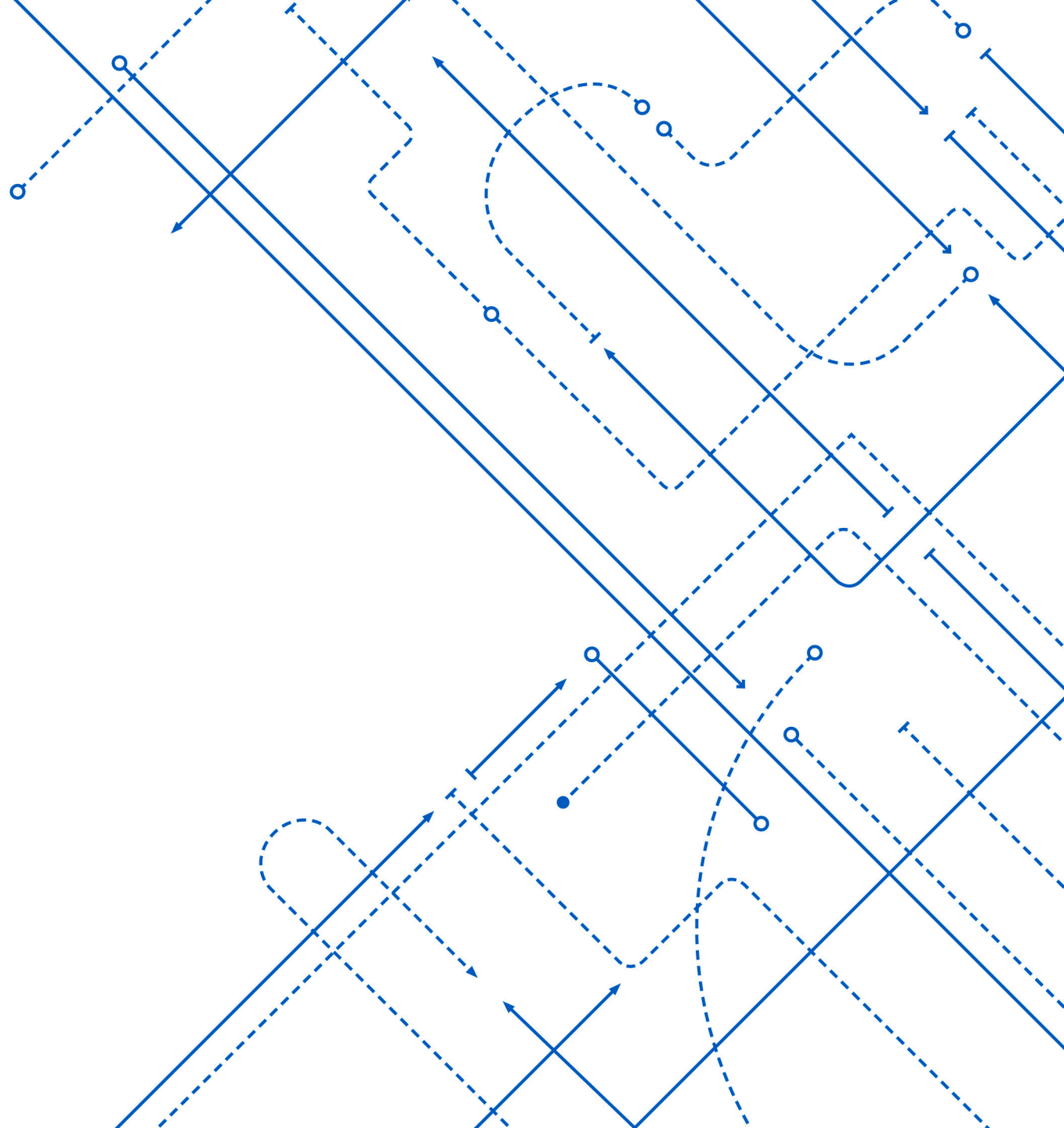
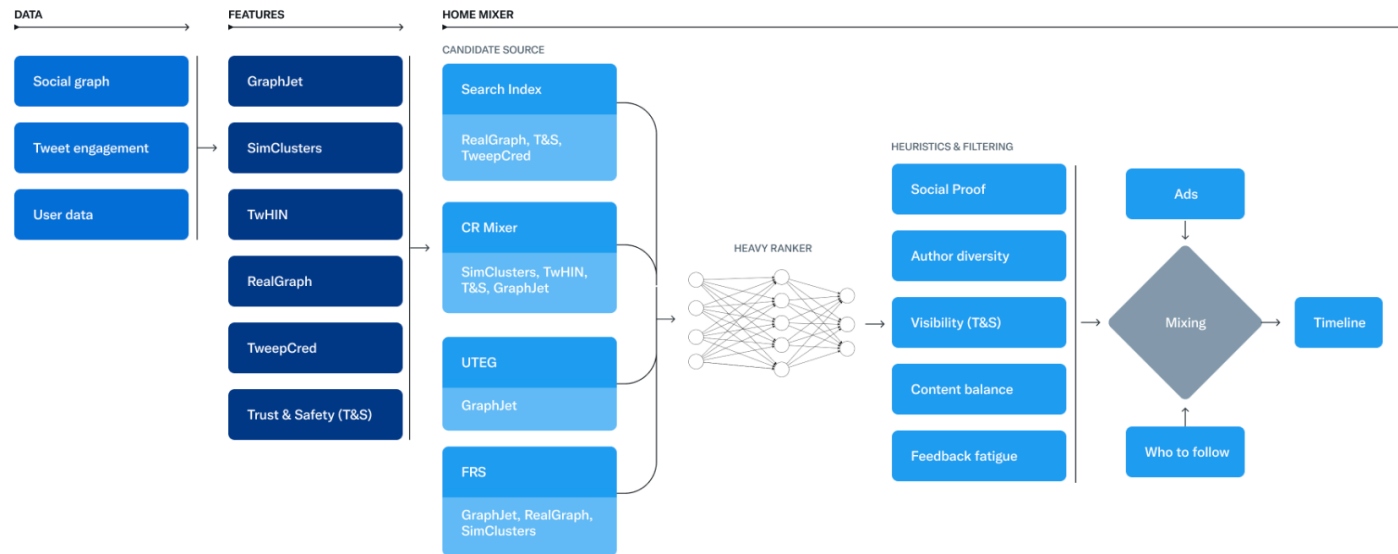


NLP & Social Media



Twitter Recommendation Algorithm

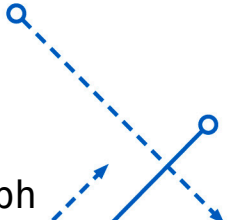
The Twitter Recommendation Algorithm is a set of services and jobs that are responsible for constructing and serving the Home Timeline. For an introduction to how the algorithm works, please refer to our [engineering blog](#). The diagram below illustrates how major services and jobs interconnect.



These are the main components of the Recommendation Algorithm included in this repository:

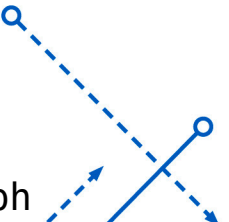
```
elon twitter algorithm Islanders at Hurricanes >>
(
  "author_is_elon",
  candidate =>
    candidate
      .getOrElse(AuthorIdFeature, None).contains(candidate.getOrElse(DDGStatsElonFeature, 0L))),
(
  "author_is_power_user",
  candidate =>
    candidate
      .getOrElse(AuthorIdFeature, None)
      .exists(candidate.getOrElse(DDGStatsVitsFeature, Set.empty[Long]).contains)),
(
  "author_is_democrat",
  candidate =>
    candidate
      .getOrElse(AuthorIdFeature, None)
      .exists(candidate.getOrElse(DDGStatsDemocratsFeature, Set.empty[Long]).contains)),
(
  "author_is_republican",
  candidate =>
    candidate
      .getOrElse(AuthorIdFeature, None)
      .exists(candidate.getOrElse(DDGStatsRepublicansFeature, Set.empty[Long]).contains)),

```



A Weakly Supervised Classifier and Dataset of White Supremacist Language

Anonymous ACL submission



Kenny and Atri's quick summary

- Most people: how would this be used and what were the implications of that?
- Most people liked the idea of weak supervision and the explicit incorporation of anti-racist text
- Common question: Could OpenAI have used something like this? Did they?
- Technical Questions
 - How is text different from images/video?
 - What is a topic / topic model?
 - Various questions about the data

Rest of today / some of Wednesday

- Differentiating various ways of accessing data
- An ad-hoc, mostly question-driven discussion about NLP and evaluation

Where does data come from?

How to get more labeled training data?

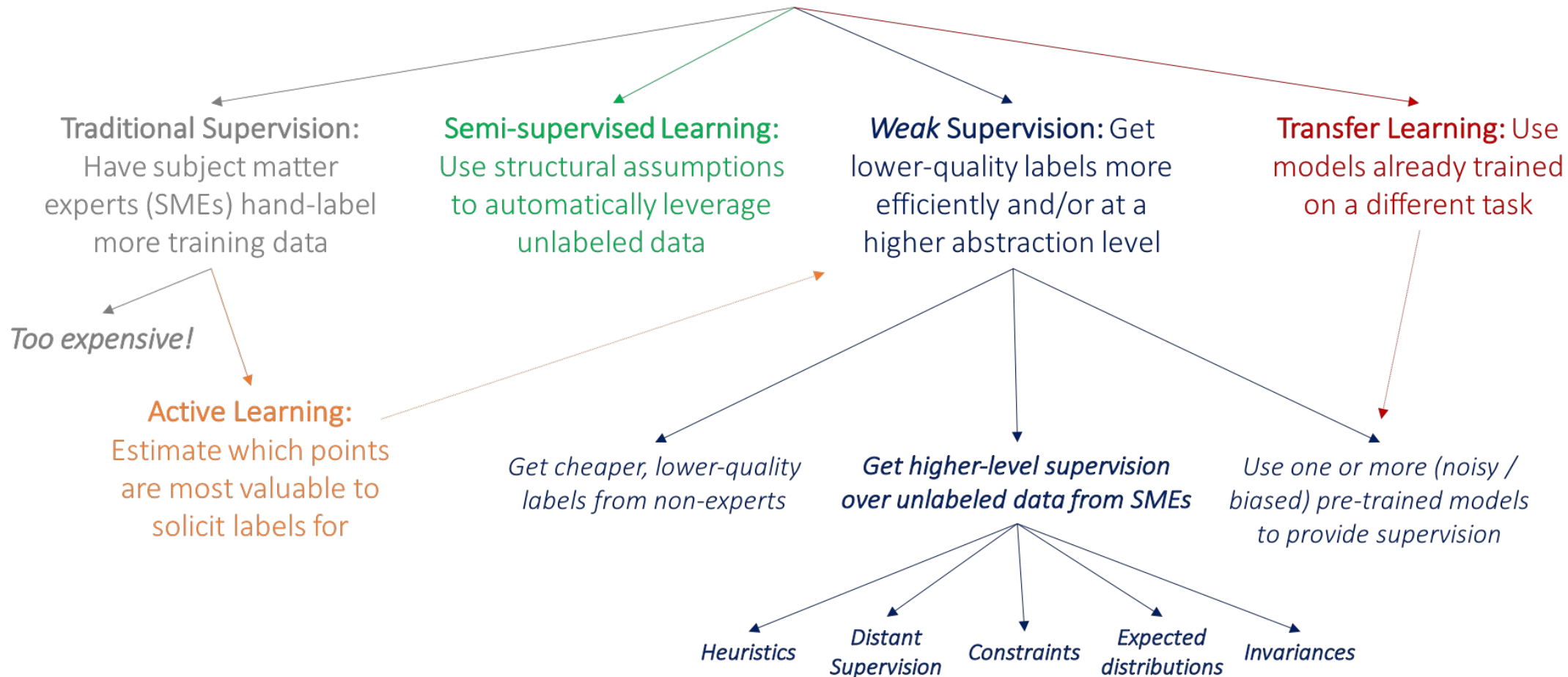
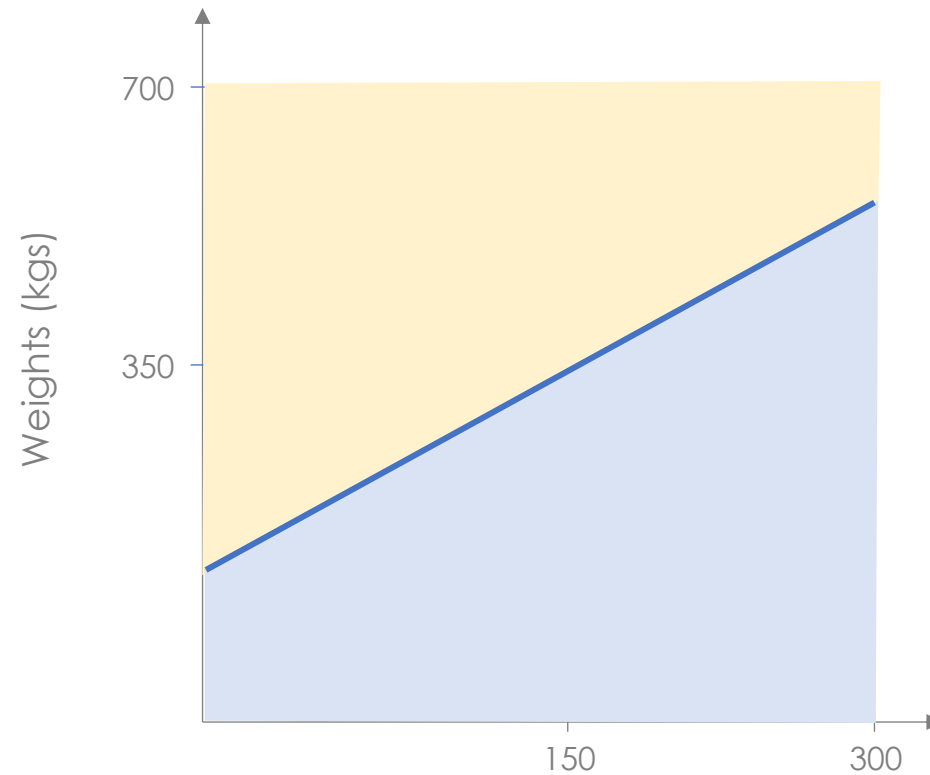


Image source: <http://ai.stanford.edu/blog/weak-supervision/>

A linear model for classification

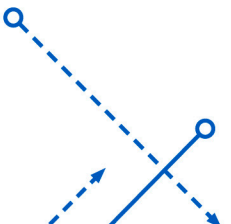


Linear models

The models we considered above (as we have seen before) are called *linear models* (because you can literally draw a line to present them in 2D). In particular, in the case of two input variables, a linear model is **completely defined** once you specify the line as which of the two sides is the positive side (and the other side automatically becomes the negative side).

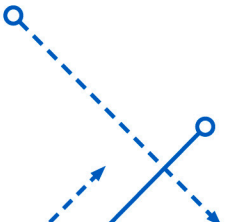
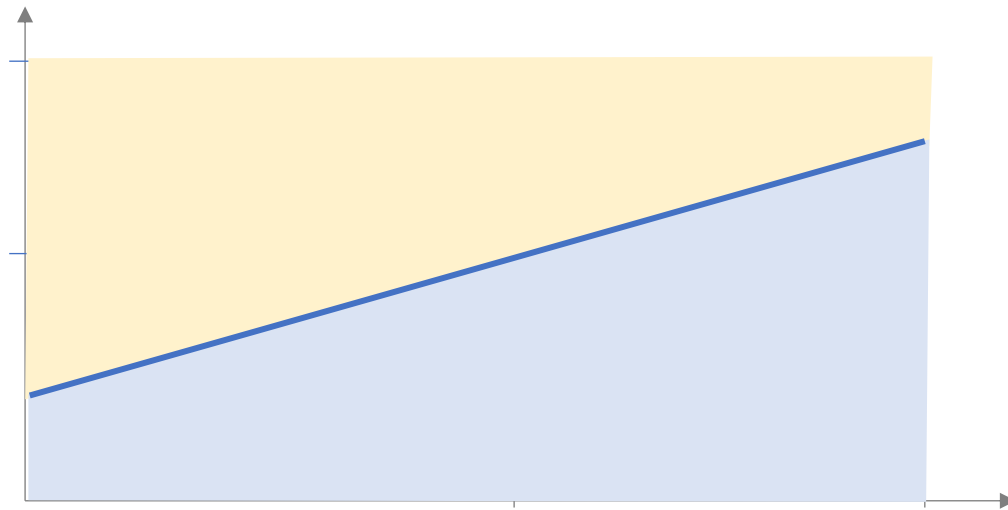
Sample task: stance detection

Stance detection: The task of determining whether someone is for or against a particular thing. We'll focus on “stance towards 4/574”



Sample task: stance detection

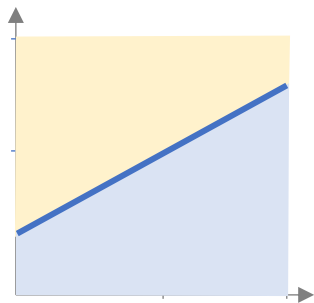
Stance detection: The task of determining whether someone is for or against a particular thing. We'll focus on “stance towards 4/574”



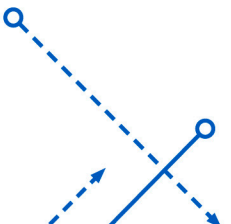
Sample task: stance detection

Stance detection: The task of determining whether someone is for or against a particular thing. We'll focus on "stance towards 4/574"

This class is garbage. The professor makes bad jokes and I can't read his handwriting

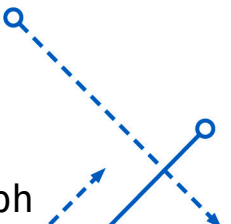


Arguably the greatest moments in human history have come when Kenny takes the floor for 4/574 each week



Stance detection in the real world

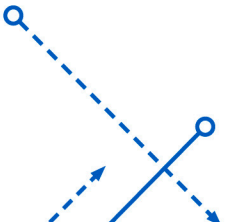
The prof's jokes are bad
and he can't make a quiz
without an error to save his
life but I occasionally learn
some stuff



Back to the example...

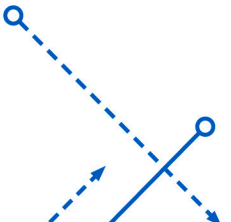
This class is garbage. The professor makes bad jokes and I can't read his handwriting

- How would you approach the task of stance detection? Specifically...
 - What would your **features** be?
 - How would you make decisions based on those features?
 - What **loss function** would you use?



Approach 1: Bag of words + Linear threshold classifier

1. Convert each course evaluation statement into a “bag of words” representation
2. Fix a weight for each word in terms of how having it in a sentence implies a positive/pro or negative/anti stance
3. Sum up the weights for all of the words to get a **score**
4. If the **score** is > 0 , predict “pro-5/474”, otherwise, predict “anti”



Approach 1, Step 1

This class is garbage. The professor makes bad jokes and I can't read his handwriting

The prof's jokes are bad and he can't make a quiz without an error to save his life but I occasionally learn some stuff

Arguably the greatest moments in human history have come when Kenny takes the floor for 4/574 each week

Approach 1, Step 2

- How might we get these values in the easiest possible way?
- ... later... how we can learn them from data

Word	Weight (w)
garbage	-5.0
bad	-3.0
can't	-0.5
bad	-3.0
error	-2.1
learn	3.0
greatest	4.0
Arguably, human, professor, handwriting, ...	0.0

Approach 1, Step 2

This class is **garbage**. The professor makes **bad** jokes and I **can't** read his handwriting

The prof's jokes are **bad** and he **can't** make a quiz without an **error** to save his life but I occasionally **learn** some stuff

Arguably the **greatest** moments in human history have come when Kenny takes the floor for 4/574 each week

Word	Weight (w)
garbage	-5.0
bad	-3.0
can't	-0.5
bad	-3.0
error	-2.1
learn	3.0
greatest	4.0
Arguably, human, professor, handwriting, ...	0.0

Approach 1, Step 3

The prof's jokes are **bad** and he can't make a quiz without an **error** to save his life but I occasionally **learn** some stuff

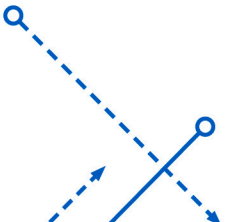
Word	Weight (w)
garbage	-5.0
bad	-3.0
can't	-0.5
bad	-3.0
error	-2.1
learn	3.0
greatest	4.0
Arguably, human, professor, handwriting, ...	0.0

Approach 1, Step 4

This class is **garbage**. The professor makes **bad** jokes and I **can't** read his **handwriting**

The prof's jokes are **bad** and he can't make a quiz without an **error** to save his life but I occasionally **learn** some stuff

Arguably the **greatest** moments in human history have come when Kenny takes the floor for 4/574 each week

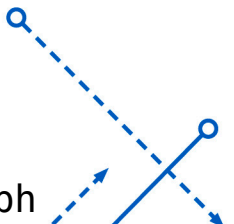


Cool!

- We have built our first classifier!
- **Quiz:** Did this classifier use (training) data at all?
- **How could it have used data to inform the model?**

...Put another way, how to learn word weights?

CIML, pg. 43



...Put another way, how to learn word weights?

An online algorithm to learn weights for the words...

The *perceptron* algorithm.

An early, well-known approach!

IMO, can complicate understanding at this point in the class

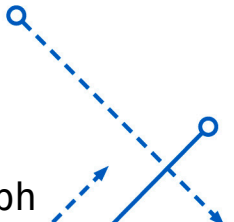
Algorithm 5 PERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```
1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$  // initialize weights
2:  $b \leftarrow 0$  // initialize bias
3: for  $iter = 1 \dots MaxIter$  do
4:   for all  $(x,y) \in \mathbf{D}$  do
5:      $a \leftarrow \sum_{d=1}^D w_d x_d + b$  // compute activation for this example
6:     if  $ya \leq 0$  then
7:        $w_d \leftarrow w_d + yx_d$ , for all  $d = 1 \dots D$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:     end if
10:  end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 
```

Algorithm 6 PERCEPTRONTEST($w_0, w_1, \dots, w_D, b, \hat{x}$)

```
1:  $a \leftarrow \sum_{d=1}^D w_d \hat{x}_d + b$  // compute activation for the test example
2: return SIGN( $a$ )
```

CIML, pg. 43



Will any of this be relevant soon?

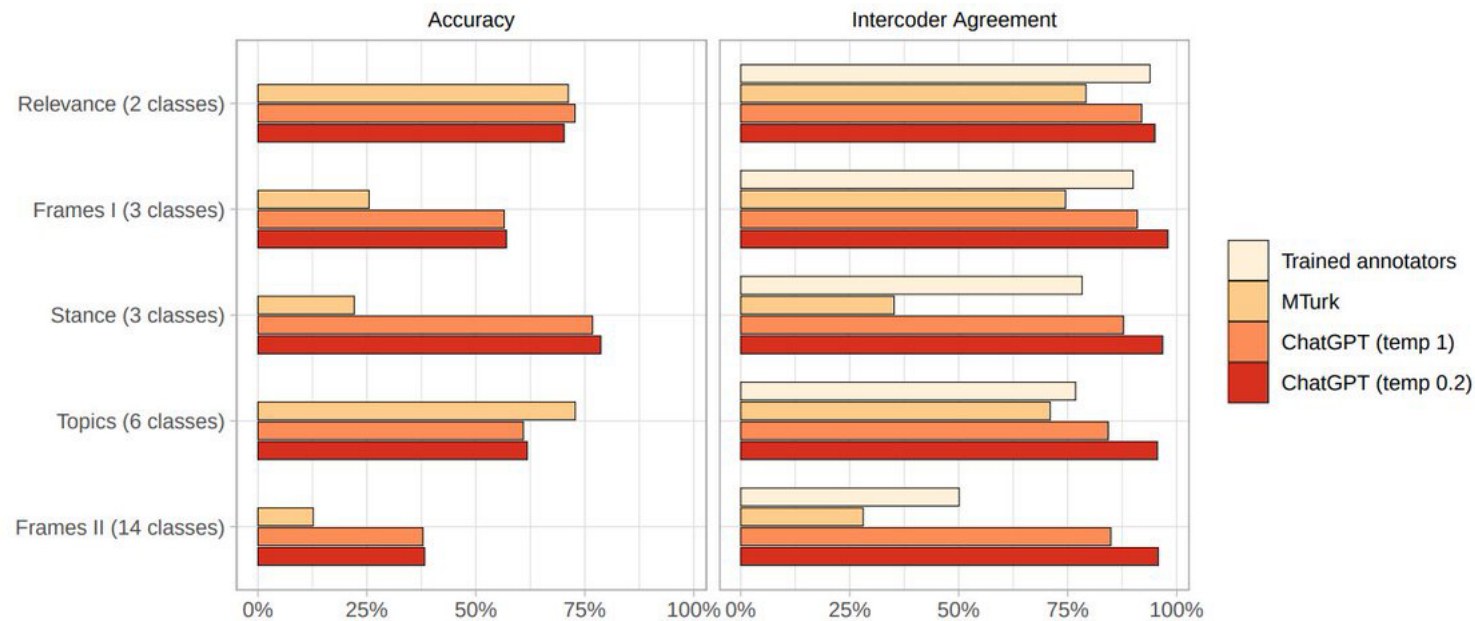


Figure 1: *ChatGPT zero-shot text annotation performance, compared to MTurk and trained annotators. ChatGPT's accuracy outperforms that of MTurk for four of the five tasks. ChatGPT's inter-coder agreement outperforms that of both MTurk and trained annotators in all tasks. Accuracy means agreement with the trained annotators.*

OK!

What questions do you have?!