

Data, Annotation, Evaluation

Jonathan May

August 23, 2025

1 History of Methodologies

1.1 Pre-Statistical (1650s/1950s through approx. 1980s)

- Mostly about modeling specific linguistic phenomena in a small number of sentences, sometimes using code
- Linguists/highly trained coders wrote down fine-grained detailed rules to capture various aspects, e.g. ‘ “swallow” is a verb of ingestion, taking an animate subject and a physical object that is edible...’
- Very time-consuming, expensive, limited coverage (brittle), but high precision
- Academically satisfying, but not good at producing systems beyond the demo phase

1.2 Statistical

- Empirical approach: learn by observing language as it’s used “in the wild”
- Many different names:
 - Corpus Linguistics
 - Empirical NLP
 - Statistical NLP
 - Nowadays we use the term ‘statistical’ to differentiate from post-rule-based but pre-neural, but for the purposes of this discussion, neural approaches are also statistical.
- Central tools:
 - corpus
 - thing to count with (i.e. statistics)
 - (later on) machine learning methodologies, software/hardware for helping with scale

- Advantages
 - Generalize patterns as they actually exist (i.e. bottom-up, not top-down)
 - Little need for knowledge (just count)
 - Systems are robust and adaptable (change domain by changing corpus)
 - Systems degrade more gracefully (corner cases captured in data)
 - Evaluations are (more) meaningful
- Limitations
 - Bound by data — can't model what you can't see – “I held the book with my arm stretched out and opened my hand. It (floated away), (fell to the ground)”
 - Big Data methods fail when the data is small or wrong – sometimes you want to try to translate Oromo news to English with 50,000 words of bible when you want 10m+words of news
 - More computationally expensive (but less human-expensive) (usually a good trade-off)
 - Methods don't have the same pattern-recognition and generalization abilities of humans learning (and putting into rule-based methods) which can lead to unintuitive brittleness... Even with very big LMs we can still get surprised by a sudden wrong answer a human would never make (and a rule-based model would just not even bother to try).

1.3 Corpus (pl: corpora): A collection of (natural language) text systematically gathered and organized in some manner

- Features:
 - Size
 - Balanced/domain
 - Written/Spoken
 - Raw/Annotated
 - Free/Pay
- Some famous (text) legacy examples – by today's standards many of these are trivially small:
 - Brown Corpus: 1m words balanced English text, POS tags
 - Wall Street Journal: 1m words English news text, syntax trees
 - Canadian Hansards: 10m words French/English parliamentary text, aligned at sentence level
 - Google Books Ngrams: 500B words

- Wikipedia: 2B words in 4.4m articles. Well connected to other languages and some implicit markup
 - Common Crawl: Attempts to crawl the web since 2008, and is regularly updated. Petabytes of data, billions of pages, about 1 trillion words. Needs a lot of cleaning, dedup, etc to do anything useful (“a large amount of documents whose content are mostly unintelligible”) (GPT-2 paper, citing an earlier google paper). Often the foundation of other collections
 - The Pile: 800GB combination of 22 high quality datasets (you may not need to do so much cleaning)
- Larger corpora and their claims
 - Corpus for GPT-3: 500b tokens of Common Crawl + Wikipedia + private data sets ‘Books1/2’ (described but not shared)
 - Corpus for GPT-4: around 13T tokens (not described in any real detail)
 - Corpus for Llama/Llama 2/3: 1.4T/2T/15T tokens of ‘publicly available’ data (Llama 1 describes specifics, others don’t)
 - Corpus for Gemini: public and private data. Very few details!
 - Dolma (<https://allennai.org/dolma>): 3T, later reduced to 2.3T tokens, completely open, well documented.
 - Dear students who read ahead: any new ones that I forgot?
 - We will revisit these corpora when we get to LLM-specific training.

1.4 corpus processing

It’s helpful to get familiar with this for quick-and-dirty checking of your data, though it’s less necessary to *learn* it anymore; you can always ask your favorite LM to give you ‘one liner gnu commands’ e.g. But knowing you can do this at all can be a revelation; access to a big data set is daunting and paging through the data may not be revealing in the right way.

word counts and ngram counts

10 most frequent words in the text

```
sed 's/ /\n/g' sawyr11.txt | sort | uniq -c | sort -k1nr | head
```

10 most frequent words in the text after removing blank lines

```
sed 's/ /\n/g' sawyr11.txt | grep -v "^$" | sort | uniq -c | sort -k1nr | head
```

10 most frequent bigrams (2 word sequences) in the text

```
sed 's/ /\n/g' sawyr11.txt | grep -v "^$" > ts.words
```

```
tail -n+2 ts.words > ts.2pos
```

```
paste ts.words ts.2pos | sort | uniq -c | sort -k1nr | head
```

count number of words/ngrams, number of word/ngram types, number of

1-count word/ngram types

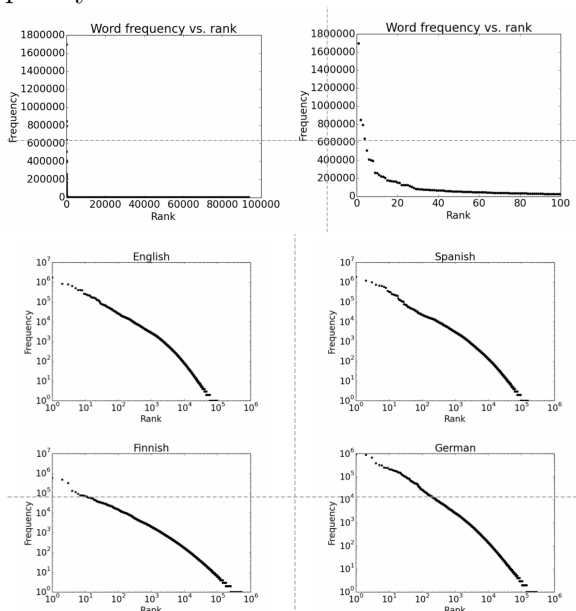
```

# words in the text without blank lines, one word per line, saved to a file
# (for convenience)
sed 's/ /\n/g' sawyr11.txt | grep -v "^$" > ts.words
# number of word tokens
wc -l ts.words
# number of word types
sort ts.words | uniq | wc -l
# number of one-count words
sort ts.words | uniq -c | awk '$1==1{print}' | wc -l
# number of two-word sequence (bigram) tokens
# (based on the answer to the number of word tokens you should know this
# without running the command...)
paste ts.words <(tail -n+2 ts.words) | wc -l
# number of bigram types
paste ts.words <(tail -n+2 ts.words) | sort | uniq | wc -l
# number of one-count bigrams
paste ts.words <(tail -n+2 ts.words) | sort | uniq -c | awk '$1==1{print}' | wc -l

```

1.5 Zipf's law

Take a naturally occurring corpus (in this case, of text). Count frequency of words. Order words by frequency, to form ranks. Then the rank of a word is inversely proportional to its frequency.



For frequency f and rank k , $f \approx \frac{1}{k^a}$ for some constant a . Thus $-a \approx \frac{\log(f)}{\log(k)}$.
Consequences:

- There will always be a lot of infrequent or unseen words

- This is true at all levels of linguistic structure (and all naturally occurring data sets, e.g. foraging patterns, solar flare intensities, crater sizes); see figure below which shows distribution of GPT3 data by language
- So we have to find clever ways of generalizing so we can get reasonable estimates for things we haven't seen often enough



1.6 Word Segmentation

- motivation from finnish but also for out of vocabulary issues
- morphology based
- unsupervised model-based
- unsupervised simple count-based (BPE): BPE example using “BITTER BUTLER BITE BUTTER BITER”

Segmentation will reduce the zipfian effect, since it reduces the vocabulary size. However, it will not eliminate it!

Nearly all of NLP benefits from a dispassionate analysis of how well the models we build perform the analysis or generation of NL data we are trying to have them analyze/generate. Aside from avoiding inherent bias and blindness toward a system's quality (which can lead to trouble down the road), there are circumstances where optimizing on the metric can be a good way of improving performance. But understanding what correctness and incorrectness is and how to measure these values can be tricky. Here we try to cover some of the key methods of evaluation in NLP. Evaluating these models means having labeled data for them; we'll discuss strategies for collection too.

2 Quantitative Analysis

2.1 Development, Test, Blind

One way to build your system is to look at some labeled data, write code that tries to get the right labels on that data, see what labels it's getting are incorrect, refine, and repeat. At the end you will have a system that does well on the corpus you looked at. But this often means you will have neglected some phenomena present **outside** your training corpus.

Worse, you may engineer things so well on training that previously correctly labeled items outside of your training corpus are now mislabeled.

This is called *overfitting* and it usually means you are taking advantage of some eccentricity in your training data that does not generally hold (e.g. all sentences with an even number of words are positive sentiment).

To avoid overfitting it's a good idea to divide your data as follows: most (80-90%) is used to build your model (train corpus), a sample (10-20%) is used to periodically evaluate but not to build the model (dev corpus), and another sample (5-10%) is not looked at and only used for evaluation, very very seldomly (test corpus). You can track overfitting: plot a graph of training vs. dev performance over time; if dev starts to go down while train goes up, that's overfitting. How much to evaluate on test? It's an art, honestly. The more you do, the more you will probably overfit, and then you'll need another test corpus to verify this.

If you don't have a lot of data, you can do what's called *cross-validation*. You divide up your data and evaluate. Then you slide the dividing lines to create another *fold* of the data. Keep doing this and every chunk of data gets to be train, dev, and test. Then average the results. I think this has the tendency to lead to overfitting more quickly but it is not infrequently used.

The extreme version of cross-validation is *n-fold* where *n* is the number of items you have; this is called *leave-one-out* validation and allows the maximum amount of training data to be used (but I trust it the least).

A hybrid strategy, where cross-validation is used for train and dev, but the test set is held constant, is probably a good compromise.

The best scenario is when someone else holds on to a piece of the data for you and you only ever get to see it once, when you're totally done with your model. This corpus is called 'blind' and is usually only used in the context of shared tasks.

It's important to consider how you divide your data. Ideally data should be independent and identically distributed (IID). You might think random selection of labeled elements would be suitable. For a collection of sentiment-labeled movie reviews this is true. However, if your corpus is a set of documents, you don't want to have one sentence from a document in training and another in test. Words and phenomena tend to cluster in a document, since a document is about some topic, with topic-relevant words, and is generally written by one author, and will have a particular style. So best practice is to divide up along *document* boundaries randomly, but pay attention to the number of evaluated items that get added to each set.

Some notes on evaluation measures follow. The choice of evaluation measures is always subject to debate (look at MT metrics workshops) but here are some guidelines I like to use:

2.2 Accuracy

For strict classification, where each item receives a label from a fixed set of labels, and the distribution of labels is reasonably even (doesn't need to be all the way even, but shouldn't be 90% one class), simple accuracy = $\frac{\text{correct}}{\text{total}}$ is a perfectly good metric.

2.3 F-Measure

For cases where there is one ‘background’ label that predominates and a relatively few instances of a ‘content’ label (e.g. named entity recognition) F-measure, which combines precision and recall, is a better choice, and is calculated on all labels *except* the background label. Precision is $\frac{\text{correct}}{\text{hypothesis}}$, i.e., how much of what you predicted is correct. Recall is $\frac{\text{correct}}{\text{reference}}$, i.e. how much of what was correct did you predict. Using either one by itself can be misleading: why?

F1 is a *harmonic mean* of precision and recall. Specifically it’s $2 \cdot \frac{PR}{P+R}$. In general $F_\beta = (1 + \beta^2) \frac{PR}{\beta^2 P + R}$; F_2 is weighted to favor recall, and $F_{0.5}$ is weighted to favor precision. In certain circumstances one may be preferred (e.g. precision is less important if there will be a downstream selection task, while recall is less important if the task is to mine information from a very large corpus and the amount of information mined is more important than assurances it has all been mined).

What is described above can be more specifically described as *micro-averaged* F1, i.e. each non-background labeling is considered in making the tallies of what is correct, then the averages are calculated. One can also consider *macro-averaged* F1, where F1 is calculated for each kind of label (e.g. in named entity recognition, PERSON, then LOCATION) where all other labels are considered to be background. These F1s are then averaged together.

Micro-F1 is the same as accuracy when there is no background class; Macro-F1 can be done in these situations too if there is class imbalance (i.e. one dominant class is mostly correct but you want to weight each class evenly instead of each item) Example:

		Gold				
		None	Person	Location	Company	Total
Hypothesis	None	N/A	0	5	10	(15)
	Person	0	200	10	0	210
	Location	0	5	40	0	45
	Company	5	0	0	10	15
	Total	(5)	205	55	20	

Micro-averaged Precision: $\frac{200+40+10}{270} = 0.926$

Micro-averaged Recall: $\frac{200+40+10}{280} = .893$

Micro-averaged F1: $2 \cdot \frac{.926 \cdot .893}{.926 + .893} = 2 * .827 / 1.819 = .909$

Macro-averaged precision: $\frac{\frac{200}{210} + \frac{40}{45} + \frac{10}{15}}{3} = .836$

Macro-averaged recall: $\frac{\frac{200}{205} + \frac{40}{55} + \frac{10}{20}}{3} = .734$

Macro-averaged F1: $2 \cdot \frac{.836 \cdot .734}{.836 + .734} = 2 * .614 / 1.539 = .782$

2.4 Granularity

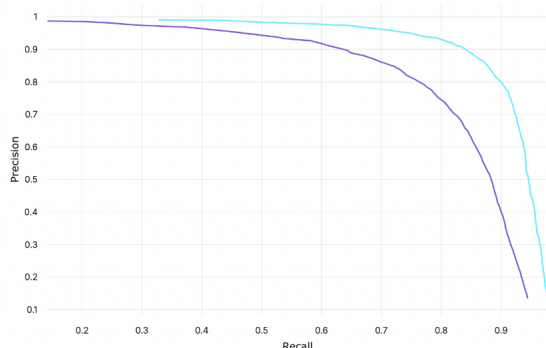
It’s important to consider what item is ultimately judged correct or incorrect. Even if each word or sentence is labeled, it may make more sense for a sequence of adjacent labels to be examined to consider whether an item is correct or not.

2.5 Rank-based evaluation

Sometimes you get to return more than one label, or you get to label your items in preferential order (e.g. information retrieval). Sometimes your results are actually pretty bad but you want to convey that your algorithm is better than random (e.g. unsupervised bilingual lexicon induction). There are a few strategies for evaluating this kind of data:

Precision/Recall@N: Consider up to N ranked items. Careful – the way this metric is implemented can vary! In e.g. information retrieval, $P@N$ means you consider N items per query and calculate precision over the $N*Q$ total items retrieved. However in e.g. bilingual lexicon induction, $P@N$ means if *any* of the N items retrieved is correct you consider the entire item is correct, i.e. precision is calculated over Q total items retrieved. If $P@N$ monotonically increases with N , it's the latter.

Precision@R: Especially in retrieval like cases, you can control the tradeoff between more precision and more recall but adjusting your threshold of returning an item. This would numerically show the precision at a fixed recall. This can also be plotted to determine ideal operating points. A calculation of the *area under the curve (AUC)* is a good single number that summarizes this tradeoff (higher is better).



2.6 Edit distance/word error rate

For structured output, especially text that is generated, the idea of ‘how much work does it take to fix this’ is very relevant, especially if the task will be, say, a first pass before human correction (this is very common in, say, the translation service industry). Given the number of substitutions S , deletions D , and insertions I made over a text of length n , the error rate (WER or TER for task T) is:

$$WER = \frac{S + D + I}{n}$$

Variants assign different cost to different operations or operations involving different components. For example, inserting/removing determiners might only cost 0.4 of an edit, while doing so for content words could cost 2.3. These values are set experimentally.

2.7 Intrinsic Vs. Extrinsic Analysis

Intrinsic: evaluate the task on its own merits. E.g. parse F1, POS tag accuracy. Evaluation ‘close’ to the task. Advantage: directly evaluates the model you’re building. Disadvantages:

does it matter?

Extrinsic: evaluate the task as it plugs into some other task (e.g. parsing in the service of summarization). Need to hold the other technology constant. There are different levels of this (narrow = pos tagging for parsing; wide = machine translation for surgery outcome). This also brings up the question of **component** vs. **end-to-end** evaluation. Should you use a sequence of noisy components (more realistic, exposes problems of interaction, can be tough to tell what actually causes the problem) or use ‘gold’ data at every point in the pipeline before your tool (better for narrow debugging, doesn’t give a realistic picture)? Both have their value. Don’t let something out before testing end to end, though.

2.8 Human judgments

The gold standard! Just like with machines, you have to tell humans how to evaluate. And if the way you tell them to evaluate is too ‘weird’ you may not get results you like.

The major advantage is you can ask for judgement calls that are specifically something machines can’t¹ do. For example:

- Given one reference translation, ask them to edit a machine translation until it *means the same thing* as the reference.
- Evaluate the sentiment of a text on a n -point scale
- Is a response sarcasm?
- Were you talking to a human?
- Did you enjoy this experience?
- ...

There are a few drawbacks:

- Humans are slow compared to machines.
- Humans are expensive compared to machines.
- Humans are inconsistent.

The last one is maybe trickiest. The same human may give two different annotations. Or two humans may be internally consistent but disagree with each other. There are ways to improve agreement:

- Have several annotate and take the most frequent label (hopefully small label set)
- Have several annotate, then talk together and resolve differences
- Have several annotate, then have a super annotator resolve differences.

¹We should discuss whether the latest versions of LLMs can in fact do these tasks

Ultimately, though, you want to track *inter-annotator agreement* (IAA):

Basic idea: ask n humans to annotate the same data. Check for their overlap. There are several metrics, most based on this basic equation: $\frac{P(a)-P(e)}{1-P(e)}$ for actual agreement a and expected agreement e . 1 for perfect agreement. 0 for chance agreement. Can be negative if agreement is worse than chance.

- Cohen’s κ (only good for 2 humans): See below. What is a good value? Wide disagreement. If you have to choose, at least 0.8 (but some might say 0.4).
- Scott’s π – same as Cohen’s κ but take squared arithmetic means to determine $P(e)$. Less informative than Cohen’s since it assumes equal distribution of responses.
- Fleiss’s Kappa (generalizes to n humans)...maybe 0.5 is ok? Generalization of Scott’s π .
- Krippendorff’s α – more general. Also allows for missing data, partial agreement. Quite complicated and computationally intensive to calculate.

Demo of Cohen’s κ :

		B	B	B	total
		pos	neut	neg	
A	pos	54	28	3	85
A	neut	31	18	23	72
A	neg	0	21	72	93
	total	85	67	98	250

A used ‘pos’ 85 times of 250 annotations = .34. B also. So for pos, $P(e) = .34 \times .34 = .1156$. Similarly for neutral, $72/250 = .288$, $67/250 = .268$, $.288 \times .268 = .077$, and for neg, $93/250 = .372$, $98/250 = .392$, $.372 \times .392 = .146$. Then, $P(e) = .1156 + .077 + .146 = .339$. $P(a) = \frac{54+18+72}{250} = .576$. Then $\kappa = \frac{.576-.339}{1-.339} = .359$.

For Scott’s π but instead we calculate $P(e)$ as $(\frac{85+85}{500})^2 + (\frac{72+67}{500})^2 + (\frac{93+98}{500})^2 = .3388$ so $\kappa = \frac{.576-.3388}{1-.3388} = .3587$. These are close because the responses are relatively equally distributed

Note: IAA not really helpful for, e.g., translation. There you just want to collect many different responses and use them all to evaluate.

Which to use? Often several are used together. Cohen’s over Scott’s, but Fleiss’ if needed, α if extra mechanism needed. Statisticians might know more.

2.9 LLM as Judge

Before we get to where LLMs come from, we of course all know we can use pre-existing models (that are generally not trained to do the task you want annotation from) as judges. And of course you can ask anything you might ask a human. There are similar cost/time considerations as with humans, but as noted above models are likely to be cheaper and faster than humans. They can be inconsistent too, though in rather different ways than humans are. Should you use them as a complete replacement for human judgements?

2.10 Statistical Significance

Ok, you have some results, automatic or human. But can you rely on them? Questions you can ask yourself:

- Are the judgements actually measuring something real?
- Are they something that we care about?
- Is it from the domain/genre that we care about?
- Is it from the right distribution?
- Are there enough examples that we can trust it?

The last question is something we can answer. See also section 4.4.3 of Eisenstein and the Berg-Kirkpatrick reading this narrative is taken from:

Let's say we have two classifiers, A and B. A is better than B on some test set x by $\delta(x)$. Null hypothesis H_0 : A is not actually better than B. If true, how likely is it that A would be better than B on some new set x' by at least $\delta(x)$? If $P(\delta(x') > \delta(x)|H_0) < 0.05^2$ then we can say with 95% confidence that H_0 is rejected and indeed A is better than B. 0.05 is called the p -value.

There are a variety of methods for establishing p -value, but one easy way that works for lots of metrics and situations where we don't have limitless test data is the following bootstrap approach:

1. Draw b bootstrap samples y of size n from x with replacement.
2. Let s be 0
3. For each y , if $\delta(y) > 2\delta(x)$, increment s
4. $p \approx s/b$

What's going on here? Since the samples are all drawn from the test set we in fact want to show how often A is better than expected, and it's expected to beat B by $\delta(x)$ since we already know it does. In the original presentation of this work, $b = 10^6$ which showed stable behavior of p calculation. This was done over a wide variety of task types.

3 Annotated Task Corpora

If you want to know how well you're doing on a task you should compare yourself to examples of the task done correctly. A catch-all term for this is an 'annotation' (also a 'labeling') of data. Here are some examples:

²Actually a random variable \mathbf{X} is used, not x'

task	input	label
POS tagging	The boy ran home	DT NN VBD NN
constituency parsing	The boy ran home	(S (NP DT NN) (VP VBD NN))
dependency parsing	The boy ran home	2-det 3-nsubj 0-root 3-advmod
sentiment	The boy ran home	neutral
translation into french	The boy ran home	Le garçon a couru à la maison

Some other kinds of annotation (the scope is limitless):

- phonetic: how was a word spoken, intoned, where were pauses taken, what words were stressed
- semantic: mark words as they're used with senses, draw semantic relatedness graphs
- pragmatics/discourse: what role does each turn play (e.g. acknowledgement, request for feedback, acceptance, marker for new phase). Is a statement a thesis, antithesis, elaboration, rebuttal, justification, etc? what are the discourse units? Resolve anaphora – what do pronouns refer to? markers for attribution (something someone else claimed happened), For more look at rhetorical structure theory (RST)

4 How to Gather Annotation

The lack of data has stymied many projects but I encourage you to think of it as an opportunity, not a roadblock! Advice: “If you want to get a lot of citations, publish a corpus.” (Philipp Koehn, prof. JHU. USC alum. Lots of released corpora.)

4.1 Found

Sometimes there is natively an annotation already in existence, though possibly not in quite the right format. A simple case is that text reviews of movies often come with star ratings, which can be turned into positive/negative sentiment. A more esoteric kind of found annotation would be, say, using text spans in wikipedia that also contain page links that have info boxes used for celebrities to determine the presence of a person mention in NE tagging. Such methods can be quite powerful but are also quite noisy; they are often referred to as ‘silver standard’ for this reason.

Beware of using a found annotation that itself was programmatically generated or that you programmatically generated from the data. This is circular reasoning and will lead you to either create a trivial data set (e.g. list of descriptions of Japanese comic characters with annotation of whether they were originally anime or manga; labels generated by finding which word comes first in the description) or one that is impossible to label (sentiment labeled data; labels are chosen by running sentiment analysis). It sounds absurd but is easier to accidentally create than you might think!

4.2 DIY

You can (and should) annotate some data yourself before trying to get other people to do it. This will help you develop your guidelines and give you a sense of how difficult the task is. However, it will still be an overestimate because you may try to write down annotation guidelines (see below), but there will be hidden assumptions that won't come out until you try to have someone else do the same annotation you're doing.

4.3 Ask a Friend

Just doing your own annotation is unwise; you could do consistent annotation but nobody would be able to follow on since you won't need to write a comprehensive standard. You also will inherently be observing any test data you produce, so your systems will likely overfit. Asking someone to try to do the annotation (advisor, colleague in your lab) is a good first test of your approach. It also gives you a way to judge IAA. If they are good enough and the annotation is simple enough you can make a sufficiently sized corpus without too much effort, you can have your friend's annotations serve as a test set, and maybe they can be a coauthor of your paper.

4.4 Hire

Once you want to get serious about annotation you'll want to parallelize and increase your throughput rate and diversity of annotator. Direct hiring (we can sometimes bring in MS or undergrad student workers) has the main advantage that you can have fairly tight control over your workers' outputs, you can have them use custom tools quite easily, and you can get and give a lot of feedback. If your annotation has special skill requirements (e.g. knowledge of a particular language) this may be the best way to go. The downside is that it can be tough to find dedicated annotators since most annotations are rather tedious. It can also be expensive; you'll be competing with other employers and your team is likely to be able to choose where to work.

4.5 Crowd

Amazon Mechanical Turk (which has declined in quality over time due to neglect) and other similar companies (lately we prefer Prolific) have made the job of hiring pools of annotators much easier, but there are a number of caveats when working with MTurk (etc). First, it is not in practice as cheap as you might imagine. Ethically you should pay a living wage; we use \$15/hr as a minimum (and probably should raise that since this number was established in maybe 2018). You will get questions about this when you publish and you may get pushback. You will also get lower quality work or no takers if your rate is too low. In practice you can't pay hourly, you have to pay by the piece. This means your annotators will have much less appetite for reading very long annotation guidelines. It also means you need to calibrate your pay by getting test annotation and estimating completion time.

Another major hassle is that there will be attempts to exploit you by not doing meaningful work. Although you are allowed to not pay for work that is not done, rejecting work already

done for low quality will earn you a bad reputation, which will lower the quality and number of workers. It is important to carefully vet workers by having them randomly annotate items you already have the answers to and then validate their responses against these. You can also have workers do an entire set you know the answer to, and if they do well, you can give them specific qualifications to do more work. For bots/non-workers, best bet is to simply not hire them again.

When possible, a more reliable way to align your intent with worker performance is to reach out in a personalized way. For MTurk, the ‘Turker Nation’ slack workspace (turkernation.slack.com) seems to be the dominant way to do this. Workers there sometimes discuss using other platforms; we have used Prolific (<https://www.prolific.com/>).

If you are doing research in a lab, annotation via crowd work may be considered human subjects research and require from one or more institutional review boards (IRBs). This is more of an issue for your PI than you but you should discuss it before proceeding.

Finally, your interaction with workers is more limited. The interface options are smaller and conveying subtle differences in annotation standards can be tricky.

Caveats aside, using crowd workers is actually quite useful and sometimes the only way to get annotation work done. It’s a good idea to try doing some crowd work to get a feel for it!

4.6 Generate?

What about using large scale language models that can be given instructions, just like humans, to simulate human annotation? This approach is increasingly being pursued.³ Is this an appropriate substitute? We will discuss this more throughout the class, but for the time being it’s important to note that there are tradeoffs that are very similar to the tradeoffs with using humans. The point of annotation and data collection is to distill some amount of knowledge into a model; if we are distilling human knowledge, this is more or less the definition of NLP. But if we are distilling a distillation, it’s unclear whose knowledge we are capturing. More practically, it’s important to monitor generated output for guardrail trips (“I’m sorry, I can’t do this because it would violate my ethics”) and bad performance, especially on “outlier” inputs (e.g. tricky or subtle cases like a negative sentiment expressed with sarcastic positivity). You may notice that the dimensions and options for how you annotate data are very similar to those for how you judge.

4.7 Inter-Annotator Agreement

While annotating you’ll want to get IAA in exactly the same way you’d get IAA for human evaluation scores per above.

5 Annotation Guidelines

It is generally a good idea to write down your intended rules for how to annotate in as plain a language as possible. Note that this need not and should not be essentially a computer

³This is a particularly hot topic in 2024 and so information and research here are subject to change.

program (otherwise it wouldn't be an interesting NLP task) but can often be quite detailed. The Part of Speech tagging guidelines for the Penn Treebank, for example, are 37 pages long! Guidelines are only effective if they're followed, of course, so you have to judge how much time your annotators will put into learning how to annotate. If you've hired them, with an hourly wage, longer manuals will probably be okay. If they're getting paid by the piece, e.g. in mechanical turk, they aren't going to spend time reading manuals and not getting paid.

Here are some classic examples. From the (37 page) PTB PoS tagging guidelines:

Adjective, superlative—JJS

Adjectives with the superlative ending *-est* (as well as *worst*) are tagged as JJS. *Most* and *least* when used as adjectives, as in *the most or the least mail*, are also tagged as JJS. *Most* and *least* can also be tagged as JJS when they occur by themselves; see the entries for these words in Section 5. Adjectives with a superlative meaning but without the superlative ending *-est*, like *first*, *last* or *unsurpassed*, should simply be tagged as JJ.

Typically, part way through the annotation, tricky cases will be uncovered, then resolved, and added to the annotation standards, like was presumably done here:

Words that refer to languages or nations, like English or French, can be either adjectives (JJ) or proper nouns (NNP, NNPS).

EXAMPLES: English/JJ cuisine tends to be uninspired.

The English/NNPS tend to be uninspired cooks.

In prenominal position, such words are almost always adjectives (JJ).

Do not be led to tag such words as proper nouns just because they occur in idiomatic collocations.

EXAMPLES: Chinese/JJ cabbage; Chinese/JJ cooking

Welsh/JJ rarebit; Welsh/JJ poetry

However, note:

EXAMPLE: an English/NP sentence (cf. a sentence of English/NP)

6 Ethics

Increasingly, we as a community are coming to terms with the ethics of the (particularly, unlabeled) data collection that has fed our models, especially as the more recent models appear to be built using *all* the human-generated text that can be found on the 'open' web and beyond. Here is a somewhat ill-organized list of ethical issues/thoughts regarding data collection:

- In June 2024 Mustafa Suleyman (DeepMind co-founder, CEO of AI at MS) said⁴: "I think that with respect to content that's already on the open web, the social contract of that content since the '90s has been that it is fair use. Anyone can copy it, recreate with it, reproduce with it. That has been 'freeware,' if you like, that's been the

⁴<https://www.inc.com/kit-eaton/microsofts-ai-chief-says-content-is-fair-game-on-open-web.html>, <https://musically.com/2024/07/01/microsofts-ai-boss-causes-a-stir-with-his-fair-use-views/>

understanding.” This is at odds with US copyright law, which give you automatic rights over your material unless you explicitly yield them. Should it be allowed to scrape data? Should academics be given different rights than companies profiting off of others’ content? Will this lead to an increasing unavailability of content?

- Suleyman also said⁵ “There’s a separate category where a website, or a publisher, or a news organization had explicitly said ‘do not scrape or crawl me for any other reason than indexing me so that other people can find this content.’ That’s a grey area, and I think it’s going to work its way through the courts.” This refers to `robots.txt` file policies. It is generally a social compact and not a law. But should it be a law? And is it okay to ignore this file?
- If you try to gather all the language-usage data you can from the internet, particularly the text data, the properties of the data will be skewed in a way that is not perhaps reflective of the usage of text as a whole. It is largely (perhaps about 60%) in English, and skews toward authorship by so-called WEIRD (Western, Educated, Industrialized, Rich, and Democratic) people. This means that annotated data built on this data will be skewed, as will unsupervised models. What are the consequences of this? What can be done about it?
- You may have noticed that there is plenty of content to be found that exhibits what you might consider anti-social behaviors or attitudes, such as can be found on places like Reddit, 4chan, or in the comments of YouTube and TikTok videos. As with WEIRD data, incorporation of this data into your models could lead to behaviors or representations that may not be aligned with your own values. But at the same time, removing such content could skew models in a way that does not represent significant demographics. What is to be done here? Do you consciously exclude the views/language/behavior of a subpopulation that you find abhorrent, in the course of research that is not specifically targeting your or others’ ethics?
- Since around 2006 (when Google Translate first had a viable, large-scale translation product) machine-generated data has mixed with human-generated data on the web, raising fears of an Ourobours effect, where one kind is mistaken for the other. There is an ever greater potential now for machine-generated data to be regarded as human-generated. What are the potential negative consequences of this? What would happen if, say, Quora was no longer given any new human-generated content? If it is possible to watermark generated text, should this be mandated? What would the negative consequences of this be?

There are emerging ethical issues constantly, so I’d like to raise them in class and amend these notes for future use.

⁵<https://www.theverge.com/2024/6/28/24188391/microsoft-ai-suleyman-social-contract-freeware>