

# Data collection & annotation

CMSC 723 / LING 723 / INST 725

Hal Daumé III [he/him]

19 Sep 2019

# Announcements, logistics

- HW2 update is posted  
**if you downloaded an old version, please get update**

# What part of your data is language?

- Input
- Output
- Supplemental annotation
- Other

# How do you collect this data

- Place microphones around the Irlbe center
- Download stuff from the web
  - Periodicals
  - Blogs / social media (twitter, weibo, facebook, etc.)
  - Discussion fora
  - Political debates / legal docs
- Emails
- Have people come in to a lab and produce language
- Crowdworkers
- Other?

How do you ensure you get the language (and variety) that you want?

- Choice of source
- Language id tools
- Heuristics based on common words

# Annotation types

- Document-level labels (classification) or scores (regression)
- Word- or span-level labels (tagging)
- Arbitrary labels (e.g., translation, captioning)

# Annotation guidelines

- How you convey to annotators what you want them to annotate
- Can be anywhere from 1 paragraph to  $\geq 104$  pages
  - e.g., [www.seas.upenn.edu/~pdtb/PDTBAPI/pdtb-annotation-manual.pdf](http://www.seas.upenn.edu/~pdtb/PDTBAPI/pdtb-annotation-manual.pdf)

# Example annotation guidelines

(h/t Naeemul Hassan, UMD Journalism, iSchool; CIKM'15, C+J'15; "ClaimBuster")



- **NFS: Non-factual sentence**

Subjective sentences (opinions, beliefs, declarations) and many questions fall under this category. These sentences do not contain any factual claim.

- *But I think it's time to talk about the future*
- *You remember the last time you said that?*

- **Unimportant Factual Sentence (UFS)**

They contain factual claims but are not check-worthy. The general public will not be interested in knowing whether these sentences are true or false.

- *Next Tuesday is Election Day.*
- *Two days ago we ate lunch at a restaurant*

- **Check-worthy Factual Sentence (CFS)**

They contain factual claims and the general public will be interested in knowing whether the claims are true.

- *He voted against the first Gulf War*
- *Over a million and a quarter Americans are HIV-positive*

# To what degree do annotators agree?

- For classification, could just measure % agreement
- Fails to capture “agreement by chance”
- Enter Cohen’s Kappa statistic:

Assume annotations  
are independent

$$\kappa \equiv \frac{p_o - p_e}{1 - p_e} = 1 - \frac{1 - p_o}{1 - p_e}$$

Accuracy

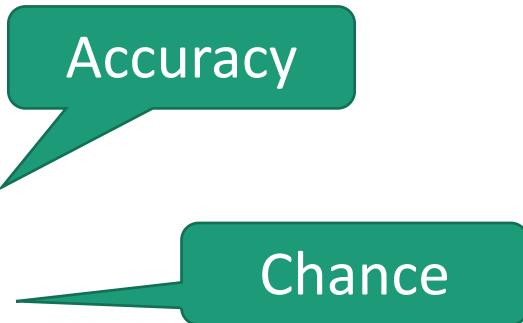
Chance

$$p_e = \sum_k \widehat{p_{k12}} = \sum_k \widehat{p_{k1}} \widehat{p_{k2}} = \sum_k \frac{n_{k1}}{N} \frac{n_{k2}}{N} = \frac{1}{N^2} \sum_k n_{k1} n_{k2}$$

$k$  is a category,  $i$  is rater,  $n_{ki}$  is # times rater  $i$  said category  $k$ ,  $N=\#$ items

# To what degree do annotators agree?

- For classification, could just measure % agreement
- Fails to capture “agreement by chance”
- Enter Cohen’s Kappa statistic:

$$\kappa \equiv \frac{p_o - p_e}{1 - p_e} = 1 - \frac{1 - p_o}{1 - p_e}$$


< 0.00 Poor agreement  
0.00 – 0.2 Slight agreement  
0.21 – 0.4 Fair agreement

0.41 – 0.6 Moderate agreement  
0.61 – 0.8 Substantial agreement  
0.81 – 1.0 Almost perfect agreement

[Landis & Koch, 1977]

# Usual annotation procedure

1. Look at the data
2. Develop initial annotation guidelines
3. Have two people go of and annotate a small number of examples
4. Compare results, update annotation guidelines
5. Throw out data from #3
6. Have multiple annotators annotate lots of data
7. Compute agreement
8. If too low, go back to #4
9. Adjudicate differences

# Case study: NUCLE

[Dahlmeier, Ng, Wu; 2013]



- English courses at CELC, all NUS undergraduate students
- Designed for students who need language support for academic studies
- Essays written as course assignments on wide range of topics

“Public spending on the aged should be limited so that money can be diverted to other areas of the country’s development.” Do you agree?

---

Surveillance technology such as RFID (radio-frequency identification) should not be used to track people (e.g., human implants and RFID tags on people or products). Do you agree? Support your argument with concrete examples.

---

Choose a concept or prototype currently in research and development and not widely available in the market. Present an argument on how the design can be improved to enhance safety. Remember to consider influential factors such as cost or performance when you summarize and rebut opposing views. You will need to include very recently published sources in your references.

Table 4: Example question prompts from the NUCLE corpus.

# Annotators

- 10 CELC English instructors
- 1414 essays, over 1.2m words
- Only doubly-annotated a subset of the corpus

# Case study: NUCLE

[Dahlmeier, Ng, Wu; 2013]

(8) Essay ID 38 ()

Your Annotation

Jump to: (8) Essay ID 38 () ▾

|<< <<

Bad Essay  Needs Editing

Assignment Prompt:  
EG1471 Assignment

- **Select** arbitrary, contiguous text spans using the cursor to identify grammatical errors.
- **Classify** errors by choosing an error tag from a drop-down menu.
- **Correct** errors by typing the correction into a text box.

Southeast Asia has the oldest and most consistent rainforests on the earth because it is in the equator zone. These forests are very necessary **for national** economies and for the living **of local** population in **the** Southeast Asia. And they are also globally essential requirements in terms of biodiversity and carbon stora **ArtOrDet (Article or Determiner)** early as a result of global demand and expanding economies. These direct causes of deforestation and forest **degrading** are mostly human **causes**.

One of the serious causes of rainforest destruction in **South East** Asia is commercial logging. Timber producing countries such as Myanmar and Indonesia log the trees for their countries' income. **For example, in** Myanmar, instead of cutting the trees in **sustainability level, it is determined based on the foreign currency earning goals**. So, this is just the short-term aim of the government rather than **long term** development to obtain foreign currency. Another thing is that the deforestation also becomes

# Case study: NUCLE

[Dahlmeier, Ng, Wu; 2013]

Error Tag	Error Category	Description / Example
<b>Word Order</b>		
WOinc	Incorrect sentence form	Why can [ <b>not we — we not</b> ] choose more intelligent and beautiful babies?
WOadv	Adverb/adjective position	It is similar to the murder of many valuable lives [ <b>only based — based only</b> ] on the couple's own wish.
<b>Transitions</b>		
Trans	Link words/phrases	In the process of selecting the gender of the child, ethical problems arise [ <b>where — because</b> ] many innocent lives of unborn fetuses are taken away.
<b>Mechanics</b>		
Mec	Punctuation, capitalization, spelling, typos	The [ <b>affect — effect</b> ] of that policy has yet to be felt.
<b>Redundancy</b>		
Rloc	Local redundancy	Currently, abortion is available to end a life only [ <b>because of — because</b> ] the fetus or embryo has the wrong sex.
<b>Citation</b>		
Cit	Citation	Poor citation practice.
<b>Others</b>		
Others	Other errors	Any error that does not fit into any other category, but can still be corrected.
Um	Unclear meaning	The quality of the passage is so poor that it cannot be corrected.

Error Tag	Error Category	Description / Example
<b>Verbs</b>		
Vt	Verb Tense	A university [ <b>had conducted — conducted</b> ] the survey last year.
Vm	Verb modal	No one [ <b>will — would</b> ] bother to consider a natural balance.
V0	Missing verb	This [ <b>may — may be</b> ] due to a traditional notion that boys would be the main labor force in a farm family.
Vform	Verb form	Will the child blame the parents after he [ <b>growing — grows</b> ] up?
<b>Subject-verb agreement</b>		
SVA	Subject-verb-agreement	The boy [ <b>play — plays</b> ] soccer.
<b>Articles/determiners</b>		
ArtOrDet	Article or Determiner	From the ethical aspect, sex selection technology should not be used in [ <b>non-medical — a non-medical</b> ] situation.
<b>Nouns</b>		
Nn	Noun Number	Sex selection should therefore be used for medical [ <b>reason — reasons</b> ] and nothing else.
Npos	Noun possessive	The education of [ <b>mother's — mothers</b> ] is a significant factor in reducing son preference.
<b>Pronouns</b>		
Pform	Pronoun form	90% of couples seek treatment for family balancing reasons and 80% of [ <b>those — them</b> ] want girls.
Pref	Pronoun reference	Moreover, children may find it hard to communicate with [ <b>his/her — their</b> ] parents.
<b>Word choice</b>		
Wcip	Wrong collocation/idiom/preposition	Singapore, for example, has invested heavily [ <b>on — in</b> ] the establishment of Biopolis
Wa	Acronyms	Using acronyms without explaining what they stand for.
Wform	Word form	Sex-selection may also result in [ <b>addition — additional</b> ] stress for the family.
Wtone	Tone	[ <b>Isn't it — Is it not</b> ] what you always dreamed for?
<b>Sentence Structure</b>		
Srun	Runons, comma splice	[ <b>Do spare some thought and time, we can make a difference! — Do spare some thought and time. We can make a difference!</b> ] (Should be split into two sentences)
Smod	Dangling modifier	[ <b>Faced — When we are faced</b> ] with the unprecedented energy crisis, finding an alternative energy resource has naturally become the top priority issue.
Spar	Parallelism	The use of sex selection would prevent rather than [ <b>contributing — contribute</b> ] to a distorted sex ratio.
Sfrag	Fragment	Although he is a student from the Arts faculty.
Ssub	Subordinate clause	It is the wrong mindset of people that boys are more superior than girls [ <b>should — that should</b> ] be corrected.

# Case study: NUCLE

[Dahlmeier, Ng, Wu; 2013]

Source	: This phenomenon opposes the real .
Annotator A	: This phenomenon opposes (the → ε (ArtOrDet)) (real → reality (Wform)) .
Annotator B	: This phenomenon opposes the (real → reality (Wform)) .

- **Identification** Agreement of tagged tokens regardless of error category or correction.
- **Classification** Agreement of error category, given identification.
- **Exact** Agreement of error category and correction, given identification.

# Case study: NUCLE

[Dahlmeier, Ng, Wu; 2013]

- **Identification** Agreement of tagged tokens regardless of error category or correction.
- **Classification** Agreement of error category, given identification.
- **Exact** Agreement of error category and correction, given identification.

Annotators	Kappa-iden	Kappa-class	Kappa-exact
A – B	0.4775	0.6206	0.5313
A – C	0.3627	0.5352	0.4956
B – C	0.3230	0.4894	0.4246
<b>Average</b>	0.3877	0.5484	0.4838

# Re-annotation protocols

- Annotate a small amount of data
- Train a simple system
- Use that system to annotate more data
- Have humans correct the annotations
- Beware the Brill Tagger Effect!  
(aka confirmation bias or automation bias)  
(also note that these annotations are no longer independent)
- Can also do active learning, eg [aclweb.org/anthology/W07-1516](https://aclweb.org/anthology/W07-1516)

# What about non-tagging problems?

- Generally no agreed-upon measure of agreement
- Standard practice:
  - Treat one annotator as “gold standard”
  - Compute (insert\_your\_metric\_here) of another annotator against that
  - Repeat for all annotators
- Does *not* take into account “chance agreement”
- But provides useful *skyline* for systems
  - *Baseline-skyline* gap tells you how much to work on this task

# When others collected data....

[Mieskes, 2017]



- Do you make your data available? HOW?

Venue	# papers	# data published	Ratio
NAACL	182	57	31.3%
ACL	231	63	27.3%
EMLNP	264	81	30.7%
Coling	337	89	26.4%
LREC	744	414	55.6%
total	1758	704	40.0%

Category	Percentage
Link available	65.2%
Link does not work	15.7%
No Link	31.4%
On Request	1.8%
Proprietary data	< 1%

- How have you protected participants of data collection process?
  - 32.8% (of 704) collected data from social media or other sensitive data
  - 96.5% did not specify anything done to anonymize
  - 7% of LREC data included recordings of elderly, children, FLLs; 90% non-anon
- Is there enough information to *re-collect* the data?

# When you collect data....

[Mieskes, 2017; Gebru, Morgenstern, et al., 2018; Bender & Friedman, 2018]



- Has data been collected?
- How was this data collected and processed?
- Was previously available data used/extended – which one?
- Is a link or a contact given?
- Where does it point (private page, research institute, official repository)?

Important to think about these issues *before* you begin collecting and annotating data!

- Data statements for NLP
  - Curation rationale
  - Language variety
  - Speaker demographic
  - Annotator demographic
  - Speech situation
  - Text characteristics
  - Recording quality
  - Other
  - Provenance Appendix
- Datasheets for Datasets
  - Motivation/curators
  - Composition
  - Collection process
  - Preprocessing/cleaning
  - Uses
  - Distribution
  - Maintenance

# Other things possibly worth knowing....

- There are various extensions to Cohen's kappa to deal with
  - Lots of annotators
  - Missing data (not every annotator annotated everything)
  - Scalar responses (vs categorical)
  - Significance tests on kappa
  - Decent catchall: Krippendorff's alpha

Krippendorff's *alpha* brings several known statistics under a common umbrella, each of them has its own limitations but no additional virtues.

[see Wikipedia page]

# Just because you *can* collect a dataset...

- Many examples over past few years of “easily broken datasets”  
(I won’t point fingers in the slides...)
- The Riloff Rule (as remembered by Hal):

If you collect a dataset and distribute it,  
if it’s not high quality and worthwhile,  
you will waste tons of other people’s time.



- Hal’s heuristic:

If tomorrow someone got 98% accuracy on this dataset,  
would I be *happy* or *sad*. If sad, don’t do it.

# Today

- Where do you get your data from?
- How do you annotate it?
- How do you measure annotator agreement?
- How can you make sure you've produced high quality data