

Annotations, Evaluation, and Generalization

CS 485, Fall 2023

Applications of Natural Language Processing

https://people.cs.umass.edu/~brenocon/cs485_f23/

Brendan O'Connor

College of Information and Computer Sciences
University of Massachusetts Amherst

- Annotations
- • Chance-adjusted calculation
 - Practical ethics example in ChatGPT
- Evaluation
 - Held-out data and overfitting
 - Classification metrics
 - Statistical testing (J&M 4.9) - hold off until later in course

Cohen's Kappa for IAA

- If some classes predominate, raw agreement rate may be misleading
- Idea: normalize accuracy (agreement) rate such that answering randomly = 0.
 - From psychology / psychometrics / content analysis
- **Chance-adjusted agreement:**

p_o : **observed** agreement rate

p_e : **expected** (by chance) rate

$$K = \frac{p_o - p_e}{1 - p_e}$$

Other chance-adjusted metrics: Fleiss, Krippendorff... see reading

When is annotating ethical?

- Annotate stressful / traumatic content
- Confidentiality / privacy
- Transparency
- Pay!

Human labeling is key to ChatGPT

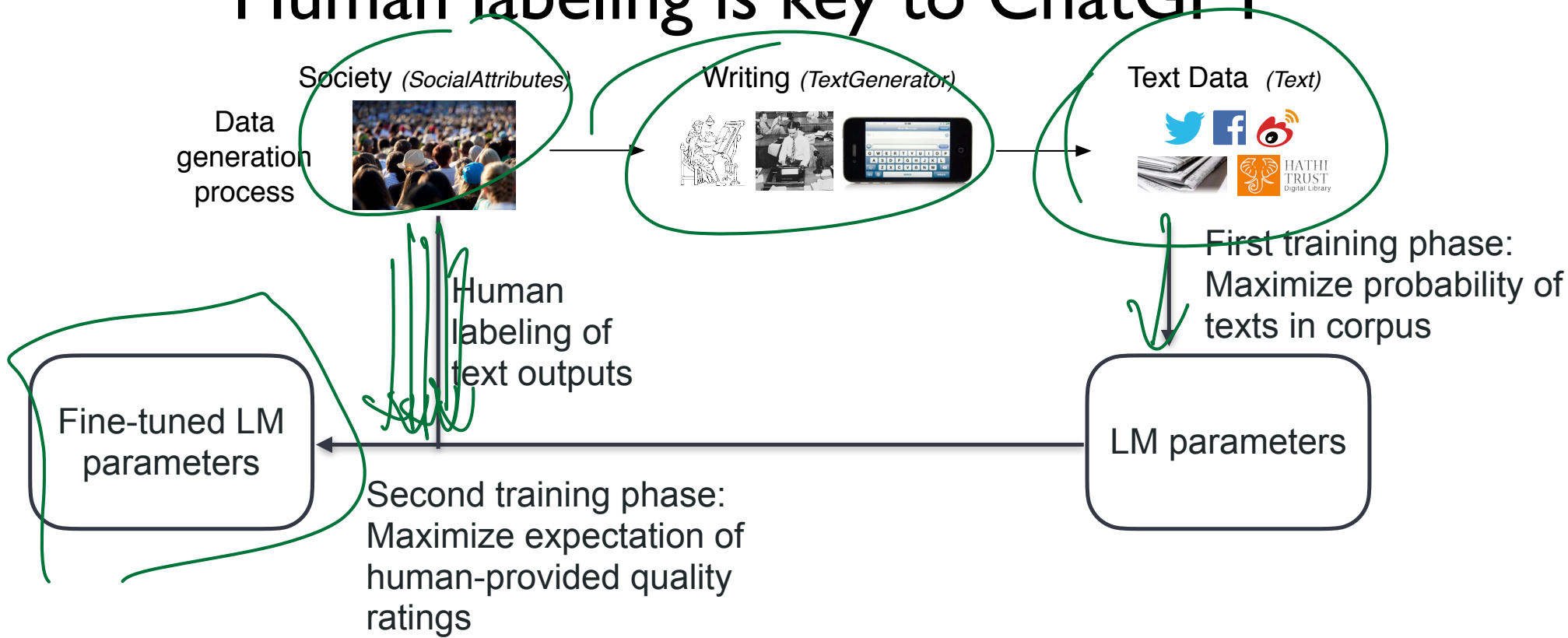


Table 3: Labeler-collected metadata on the API distribution.

Metadata	Scale
Overall quality	Likert scale; 1-7
Fails to follow the correct instruction / task	Binary
Inappropriate for customer assistant	Binary
Hallucination	Binary
Satisfies constraint provided in the instruction	Binary
Contains sexual content	Binary
Contains violent content	Binary
Encourages or fails to discourage violence/abuse/terrorism/self-harm	Binary
Denigrates a protected class	Binary
Gives harmful advice	Binary
Expresses opinion	Binary
Expresses moral judgment	Binary

'That Was Torture;' OpenAI Reportedly Relied on Low-Paid Kenyan Laborers to Sift Through Horrific Content to Make ChatGPT Palatable

The laborers reportedly looked through graphic accounts of child sexual abuse, murder, torture, suicide, and, incest.

By **Mack DeGeurin** Published January 18, 2023 | Comments (6) | Alerts

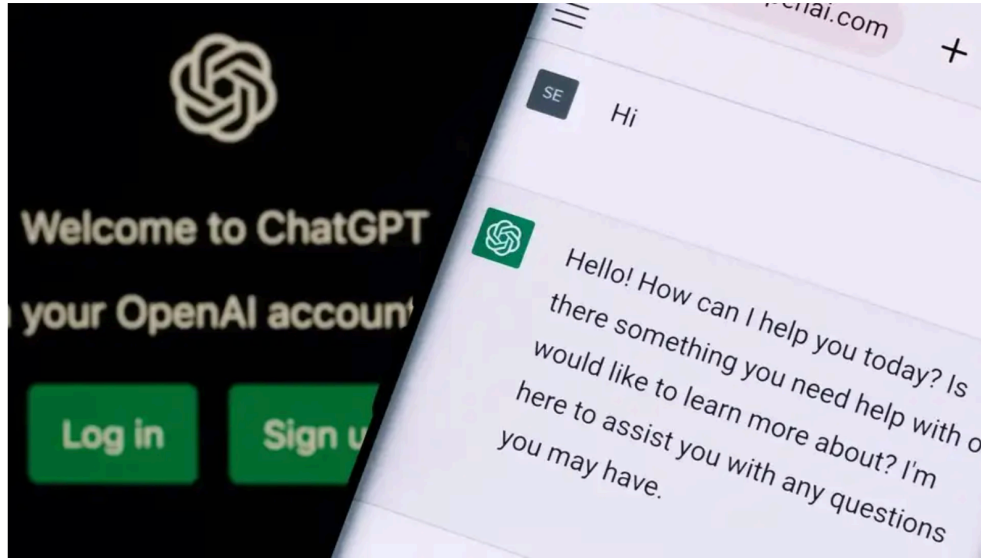


Image: Ascannio (Shutterstock)

You n



Held-out data for evaluation

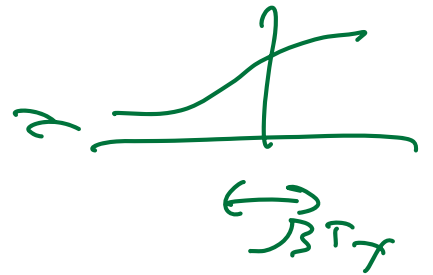
- How well will my classifier work in the future?
- Can we look at classifier accuracy on training data?

X	y	$P(y=1/x)$
I love dog	1	0.8
good dog	1	0.7
cat cat	0	0.3

"dog" occurs in $y=1$

never occurs in $y=0$

$$P(y=1/x; \beta) = \frac{1}{1 + e^{-\beta^T x}} = \sigma(\beta^T x)$$



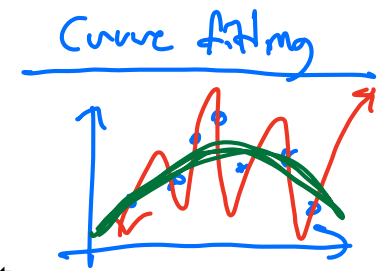
β_{dog} ↑

⇒

$$\beta^T x = \sum_i \beta_i x_i$$

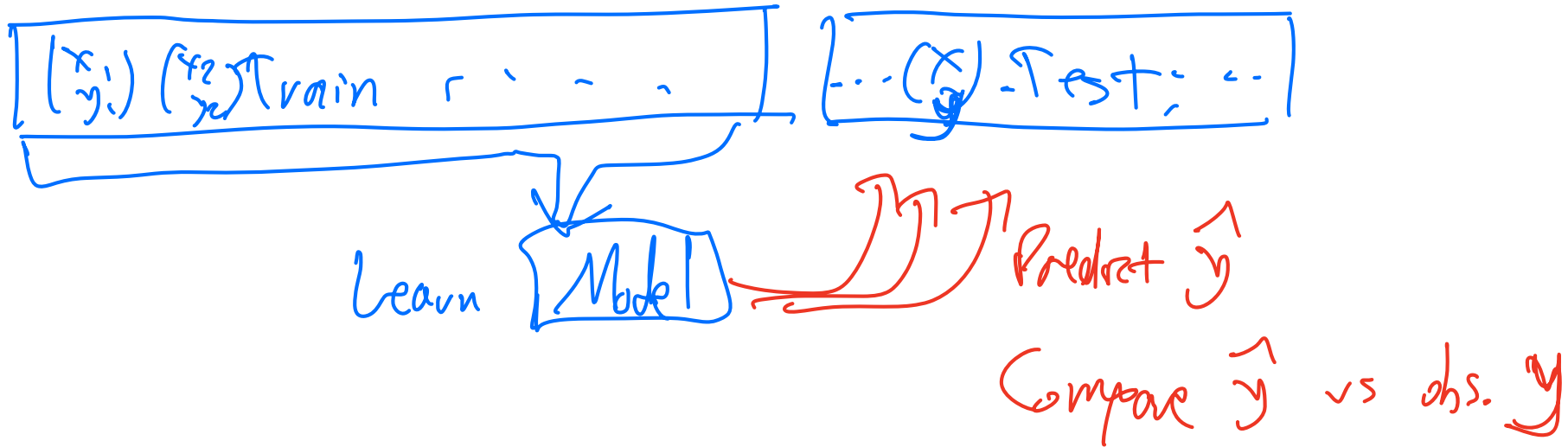
$$= \beta_{\text{dog}} x_{\text{dog}} + \beta_{|0|} x_{|0|} + \dots$$

Held-out data for evaluation



- Need to diagnose how much your model is **overfitting** the training set
- Data splits are key. Some ways to split:
 - Training set -vs- test set
 - Training set -vs- "validation"/"development" set -vs- test set
 - Cross-validation (within training set) -vs- test set

early eval for model tuning



Cross-validation

- Cross-validation (within training set) -vs- test set

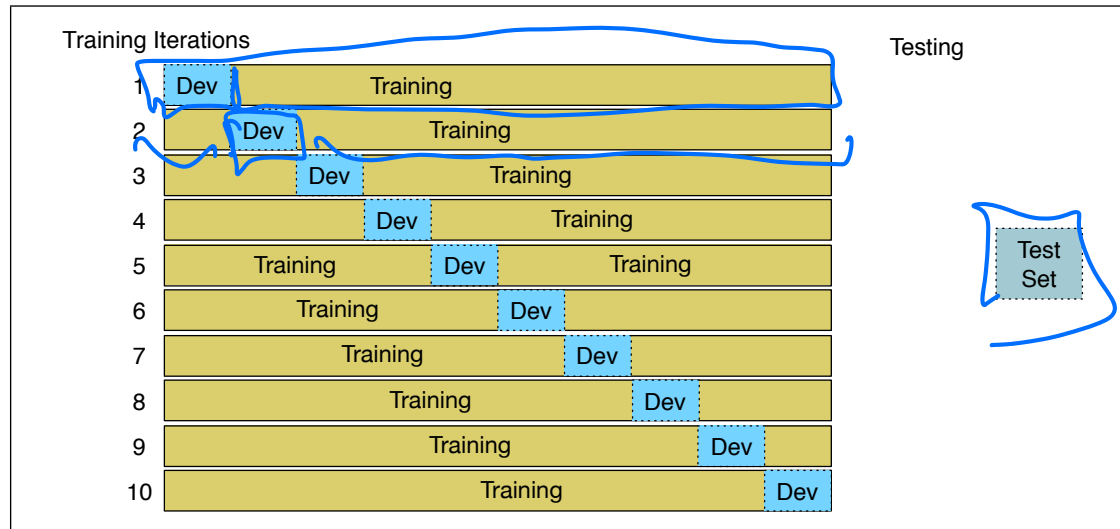


Figure 4.7 10-fold cross-validation

Regularization in Naive Bayes

$$\frac{p(y=1|\vec{w})}{p(y=0|\vec{w})} = \frac{p(y=1) \prod_i p(w_i|y=1)}{p(y=0) \prod_i p(w_i|y=0)}$$

$$= \frac{p(y=1)}{p(y=0)} \prod_i \frac{p(w_i|y=1)}{p(w_i|y=0)}$$

↑ 2 ↗
 higher regul.
 more biased model

LR > 1: favors
 y=1

LR < 1: favors
 y=0

LR(w)

$$LR(dog) = \frac{[c(dog, y=1) + \alpha] / [c(y=1) + \alpha |V|]}{[c(dog, y=0) + \alpha] / [c(y=0) + \alpha |V|]}$$

ignore

$\alpha = 0$ " "
 $\alpha = \text{tiny}$ → LR ↑
 $\alpha \rightarrow \infty$ → LR = 1

Regularization in logistic regression

- If "dog" only occurs for class **k**, what weight will it get? $\Rightarrow \beta_{\text{dog}, k} = \infty$
- Consider MLE training:

$$\text{arg max}_{\beta} \log P(y_1, \dots, y_n / x_1, \dots, x_n; \beta)$$

See L2 norm
Quadr. penalty

- Solution: **regularized** training for logistic regression

$$\text{arg max}_{\beta} \left[\log P(y_1, \dots, y_n / x_1, \dots, x_n) - \lambda \sum_{j=1}^{M_{\text{feat}}} (\beta_j)^2 \right]$$

λ control overfitting

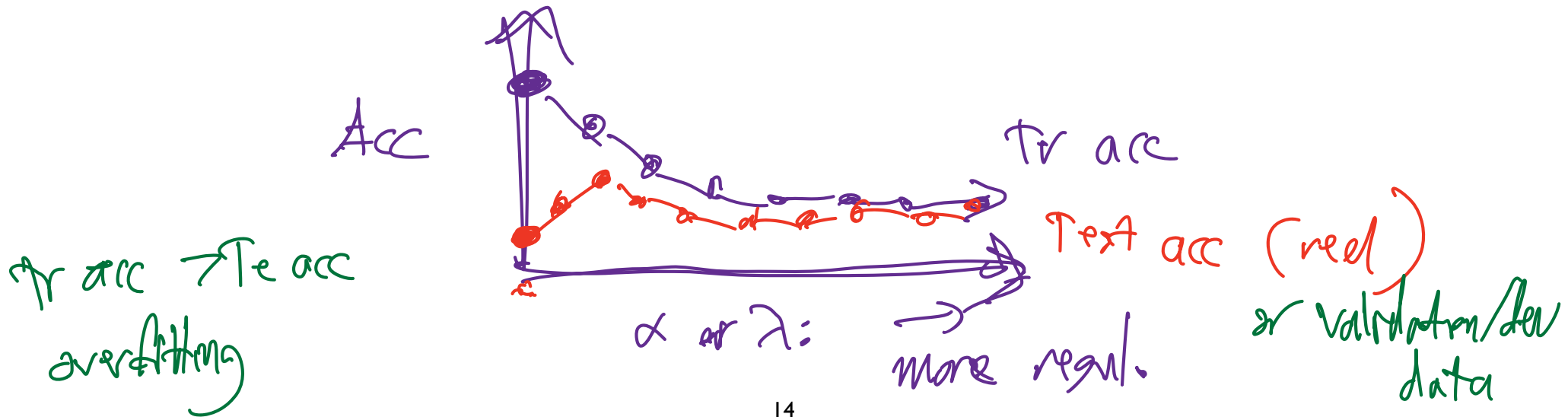
"penalty hyperparameter"
e.g. $\lambda = 1$

Regularization tradeoffs

- No regularization <-----> Very strong regularization

Bad generalization
Great training acc!

Similar preds, on all examples
Similar test & tr. acc



Do I have enough labels?

- For training, hundreds to thousands of annotations may be needed for reasonable performance
- Current work: how to usefully make NLP models with <10 or <100 training examples. "Few-shot learning"
- Exact amounts are difficult to know in advance. Can do a **learning curve** to estimate if more annotations will be useful.

Evaluation metrics

- Accuracy =

$$\sum (y = y^{(k)})$$

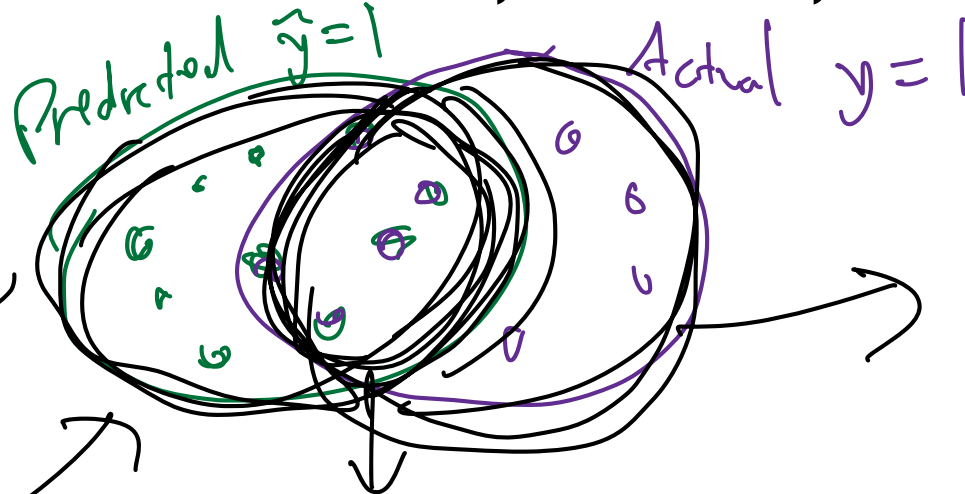
		gold standard labels		
		gold positive	gold negative	
system output labels	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

Figure 4.4 A confusion matrix for visualizing how well a binary classification system performs against gold standard labels.

- But do we care about false positives and negatives equally?
- What about rare classes?

Precision, recall, F1

DWS
 average



False Pos (FP) True Pos (TP)

		gold standard labels		
		gold positive	gold negative	
system output labels	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$	accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$	

- macro vs. micro F1

Figure 4.4 A confusion matrix for visualizing how well a binary classification system performs against gold standard labels.