

# Inter-Annotator Agreement

## Adjudication & Modeling Introduction

Jin Zhao

Brandeis University

March 30, 2025

# Today's Agenda

- ① Resolving annotator disagreements
- ② Adjudication strategies
- ③ Creating gold standard datasets
- ④ LLM-assisted adjudication
- ⑤ Introduction to modeling with annotated data
- ⑥ From annotations to trained models

*Note: No class April 1 (Passover), No class April 6-8 (Passover Break)*

**Assignment:** HW 3 due

# Why Resolve Disagreements?

**Disagreements are inevitable**

**But ML models need single labels:**

- Supervised learning requires ground truth
- Multiple labels per instance is problematic
- Evaluation needs clear correct answers

**Creating gold standard:**

- Authoritative, adjudicated labels
- Used for training and evaluation
- Represents “correct” annotation

# Adjudication Strategies

## Common approaches:

### ① Majority voting:

- Most common label wins
- Simple, scalable
- Requires odd number of annotators

### ② Expert adjudication:

- Expert reviews disagreements
- Higher quality, more expensive

### ③ Discussion:

- Annotators discuss until consensus
- Time-intensive but educational

### ④ Probabilistic:

- Weight by annotator reliability
- More sophisticated

# Majority Voting

## Simple and effective

### Process:

- ① Collect labels from all annotators
- ② For each item, select most frequent label
- ③ Handle ties (random, expert review, or skip)

### Advantages:

- Easy to implement
- Scalable to large datasets
- No additional annotation needed

### Disadvantages:

- Ignores annotator quality differences
- Majority can be wrong
- Doesn't improve guidelines

# Expert Adjudication

**For higher quality gold standard**

**Process:**

- ① Identify items with disagreement
- ② Expert reviews each case
- ③ Expert makes final decision
- ④ Optionally: update guidelines based on patterns

**When to use:**

- High-stakes tasks
- Creating evaluation benchmarks
- Complex annotation requiring domain expertise

**Cost:** More expensive and slower

# LLM-Assisted Adjudication

## New approach for 2025

### Options:

#### ① LLM as tie-breaker:

- When annotators disagree, ask LLM
- Use as third “vote”

#### ② LLM reasoning:

- Ask LLM to explain which label is correct
- Human reviews LLM reasoning

#### ③ LLM confidence:

- Accept human majority if LLM confidence low
- Review if LLM disagrees with high confidence

**Caution:** LLM shouldn't be sole adjudicator for evaluation data

# Creating Gold Standard Dataset

## Complete workflow:

- ① Multiple annotators label data
- ② Calculate IAA to verify quality
- ③ Identify disagreements
- ④ Apply adjudication strategy
- ⑤ Verify adjudicated labels
- ⑥ Document process

## Best practices:

- Keep original annotations (for analysis)
- Document adjudication decisions
- Track problematic patterns
- Update guidelines for future annotation

# From Annotations to Models

## The ML pipeline:



**Annotation quality directly affects model quality**

# Training on Annotated Data

## Classification tasks:

- Input: text features (BoW, embeddings)
- Output: predicted label
- Algorithms: Logistic Regression, SVM, Neural Networks

## Sequence labeling tasks:

- Input: token sequences
- Output: BIO label sequences
- Algorithms: CRF, BiLSTM-CRF, BERT-NER

**Key insight:** Model learns patterns from annotations

## Classification:

- Accuracy, Precision, Recall, F1
- Confusion matrix

## Sequence labeling:

- Token-level accuracy
- Entity-level Precision/Recall/F1
- Exact match vs. partial match

## Important:

- Evaluate on held-out test set
- Report multiple metrics
- Compare to baseline

# Annotation Quality and Model Performance

## The relationship:

- Higher IAA → cleaner training signal
- Cleaner signal → better model
- Better model → higher accuracy

## Upper bound:

- Human agreement is ceiling for model
- If humans agree 80%, model likely  $\leq 80\%$
- (Though models can sometimes beat individuals)

**Rule of thumb:** Expect model F1  $\approx$  Human IAA

# Common Modeling Pitfalls

## Data issues:

- Training on unadjudicated data
- Test set contamination
- Class imbalance

## Evaluation issues:

- Overfitting to dev set
- Cherry-picking metrics
- Not comparing to baseline

## Process issues:

- Not documenting preprocessing
- Inconsistent tokenization
- Not releasing data/code

# Semester Project: Gold Standard Due

**After Passover break (Apr 15):**

**Deliverables:**

- ① Adjudicated gold standard dataset
- ② Documentation of adjudication process
- ③ Final IAA report
- ④ Updated guidelines (if changed)

**Next steps:**

- Train baseline model
- Evaluate on gold standard
- Analyze errors
- Prepare final report

# Next Class: Modeling and Evaluation I

**Lecture 21 (Apr 13): Annotation-Informed Modeling I**

*Note: No class April 6-8 (Passover Break)*

## Topics:

- Training models on annotated data
- Handling annotation uncertainty
- Multi-annotator learning
- Soft labels vs. hard labels
- Human vs. LLM annotation comparison

**Project:** Gold standard dataset due

**Assignment:** HW 4 assigned

# Key Takeaways

- ① **Adjudication** resolves disagreements to create gold standard
- ② **Majority voting** is simple but may ignore annotator quality
- ③ **Expert adjudication** produces higher quality but costs more
- ④ **LLMs can assist** adjudication but shouldn't be sole judge
- ⑤ **Model quality** is bounded by annotation quality
- ⑥ **Document** your adjudication process for reproducibility

# Questions?

## Questions?

Office Hours: Wednesdays 1-3pm, Volen 109

 [jinzhaob@brandeis.edu](mailto:jinzhaob@brandeis.edu)