



Human-in-the-Loop Synthetic Text Data Inspection with Provenance Tracking



Hong Jin* Kang



Fabrice* Harel-Canada



Muhammad Ali Gulzar



Nanyun Peng



Miryung Kim

So what's the problem?

Data augmentation (DA) is commonly used to expand datasets and improve models but often compromises quality



QA to the rescue?







Human effort is expensive

Individual inspection is time-consuming

Repetitive tasks undermine focus

A new tool for scaling text quality analysis



Provenance Tracking

- 1. Transforms applied during augmentation
- 2. Text features that might cause context-specific failures (e.g. negations)

Assistive Labeling

- **1. Quality Metrics**
 - Alignment
 - Fluency
 - Grammaticality

2. LLM Guidance

Transformation Provenance

- Tracking **sequence of transforms** applied during augmentation
- Helps identify specific transforms that frequently result in low-quality or mislabeled data
- Implemented as a simple wrapper compatible with many augmentation frameworks
 - NL-Augmenter, Sibyl,
 EDA, TextAttack and more



Input	Transform	Args	Output
l cannot	Τ ₁	{prob: 0.2}	I can't be
I can't be	Τ ₂	{idx: 3}	l can't un
l can't un	T _n	{rand: 0.4}	l can't un

Feature Provenance

- Grouping on shared linguistic features derived from Abstract Meaning Representation graphs
- Helps identify how specific features might misbehave in the presence of certain transforms

"We should probably focus on every case where **WordDeletion** and :negation co-occur!"



Example Features

- :negation
- :polarity
- :purpose
- :range
- :condition
- :expressive
- :choice

Quality Metrics

• Alignment

- Does the label match the new text?
- Scored via confident learning, which predicts label errors

• Fluency

- Does the new text sound realistic?
- Scored by change in LLM perplexity
- Grammaticality
 - Does the new text violate grammar?
 - Scored by change in grammar errors



LLM Guidance

- How do humans respond to LLM assistance while inspecting texts?
- Using an LLM (gpt-3.5-turbo) to provide pre-loaded:
 - Prediction
 - Justification
- New flag indicating when LLM disagrees with current label





Computer Science

User Study Design



- A simplified version of INSPECTOR
- Lacks provence and assistive labeling



- Recruited 15 Computer Science students 11 PhD, 2 MS, 2 Undergrad
- Users had varied levels of familiarity with data inspection + ML



- Sentiment Analysis: SST-2
- Hate Speech: Tweet Eval (Hate)



- Identify texts with correct labels as efficiently and accurately as possible
- Inspect one dataset w/ INSPECTOR and the other w/ Annotator

Research Questions



Reduction in Human Effort

Does INSPECTOR increase efficiency in identifying texts with correct labels?



Technique Utility

How useful is each effort reduction technique offered by INSPECTOR?



Model Robustness

Are models more robust when trained using data identified with INSPECTOR?

Reduction in Human Effort



INSPECTOR helped analyze **3 * - 4 *** more data in the same amount of time

Technique Utility



Model Robustness

Attack Success Rate Before + After INSPECTOR



Takeaways



https://github.com/UCLA-SEAL/ProvenanceInspector



Thank You!

Questions?

UCLA Computer Science

Outline

- Motivation
 - DA expands datasets by applying intuitive transforms to existing text
 - While DA generally improves model performance, the new data is often noisy and of lower quality
 - Inspecting data quality at scale is challenging because of costs and human psychology (boredom with the task).
- We propose Inspector, a tool for scaling text quality analysis which features:
 - Provenance Grouping:
 - Transformations
 - AMR Features
 - Assistive Labeling: alignment, grammaticality, fluency, and LLM best-guess
- Demo interface
 - Show main panels and describe the intended workflow

Outline

- Study Design
 - Participants
 - Datasets
 - o Task
- Results
 - RQ1: reduction in human effort? 3-4X
 - RQ2: how helpful is each of inspector's features?
 - RQ3: models more robust? Up to a 32% improvement
- Takeaways
 - Transform provenance regarded as more useful than LLM best guess (features not at all)
 - Assistive labeling is helpful for trust and users felt more confident in their data quality
- Final slide with links

Grouping by Provenance

Provenance Grouping

Transformations

Tracking all transforms used to arrive at potentially noisy augmented text



Features

Extract abstract meaning representation (AMR) features from texts







Provenance Grouping

Transformations

Tracking all transforms used to arrive at potentially noisy augmented text



Features

Extract abstract meaning representation (AMR) features from texts





- A simplified version of INSPECTOR
- Lacks provence and assistive labeling



- Sentiment Analysis: SST-2
- Hate Speech: Tweet Eval (Hate)



- Recruited 15 Computer Science students 11 PhD, 2 MS, 2 Undergrad
- Users had varied levels of familiarity with data inspection + ML



- Identify texts with correct labels as efficiently and accurately as possible
- Inspect one dataset w/ INSPECTOR and the other w/ Annotator



Number of Inspected Texts

UCLA Computer Science

Takeaway

INSPECTOR helped analyze

3* - **4*** more data

in the same amount of time



Rating vs. Technique



Attack Success Rate Before + After INSPECTOR

UCLA **Computer Science**

Takeaways

Main Takeaway

"[I could] reason about whether a specific transform can lead to a reduction of data quality."

Effective tool for **scaling** inspection

Empowers humans to build diverse, **provenance-guided** strategies

"sentences with a grammar score < 0.92 are almost always low-quality."

UCLA Computer Science

Transform	Mark all as high quality	Mark all as inspected	Examples of transformed and labels	texts
RandomSwapQwerty (substitues random characters with adjacent keys on a keyboard) (35 total matching instances, select all 2 selected instances (28,47	Batch Labeling Inspection Statistic Transform Example	ch ling ction ctics	Niggas be getting on here crying bgout not getting no play everyday lol there are women that like weird niggas,funny niggas, athletic niggas,thug niggas etc. Find a bitch in vo arena and stfu	not hate
), 7 already inspected; 5 high quality)		med les	In a country sca <mark>ff</mark> red by violence, LGBTI activism is not always welcome. This is Sofia and Daniel's story.	not hate

Feature	Mark all as Mark all as high quality inspected	Examples of original texts and labels		
Has a description of a location, e.g., "in Los Angeles"	Batch	it takes a strange kind neg of laziness to waste the talents of robert forster, anne meara, eugene levy, and reginald		
(70 total matching instances, select all	Labeling Inspection Statistics	veljohnson all in the same movie. audrey tatou has a pos knack for picking roles that magnify her		
1 selected instances (9), 24 already inspected; 11 high quality)	Examples w/ Features	outrageous charm, and in this literate french comedy, she's as morning-glory exuberant as she was		

sentence	Align	Fluency	Grammar	LLM predictions	LLM explanation	label-LLM label consistency	label	👍 high quality	√ Inspected
runs on the pure adrenalin of pacino'c performance.	1	0.5	1	pos	The sentiment of the text is positive. The phrase 'pure adrenalin' suggests excitement and intensity, which is	consistent	pos		
also reachn't embarrassed to make you is for the tissues		0.33 <mark>uality Me</mark>	1 trics	neutral	The sentiment of the text is neutral. The text is difficult to interpret as it appears to be a jumbled collection of words and does	inconsistent	pos		

