

**PUNS ARE NOT TTT! Think about what you're posting!**

**Mhod Phost**

**Keep things on topic!**

- Puns/jokes are not TTT
  - Relatable memes are not automatically TTT
  - Murderedbywords twitter posts are not TTT
  - Dad jokes are not TTT
- ...unless there is a TTT statement in it!

**Posts that are simply the literal actual truth ARE OFF TOPIC.**

Can NLP be used to distinguish puns from truths?

By Chris Ratigan

# The Subreddits

- ▶ TechnicallyTheTruth (TTT)
  - ▶ Founded in 2017
  - ▶ Devoted to unexpected truths
  - ▶ Often reposts of images from other sites
- ▶ Dadjokes
  - ▶ Founded in 2011
  - ▶ When it becomes apparent
  - ▶ Primarily text-based posts.



**Why was 4 scared to ask out 5?**

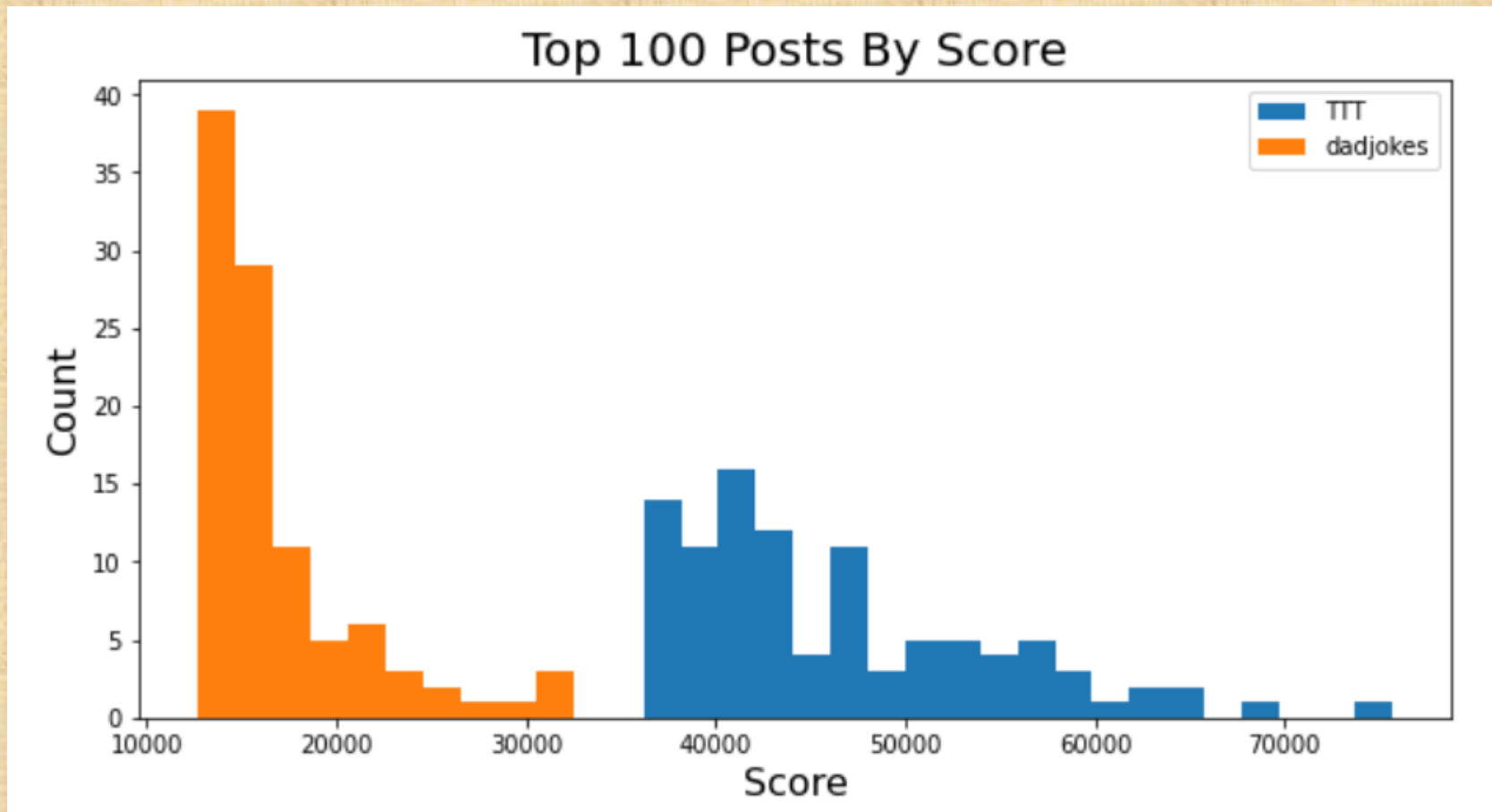
Because 4 was  $2^2$

# Getting the Data

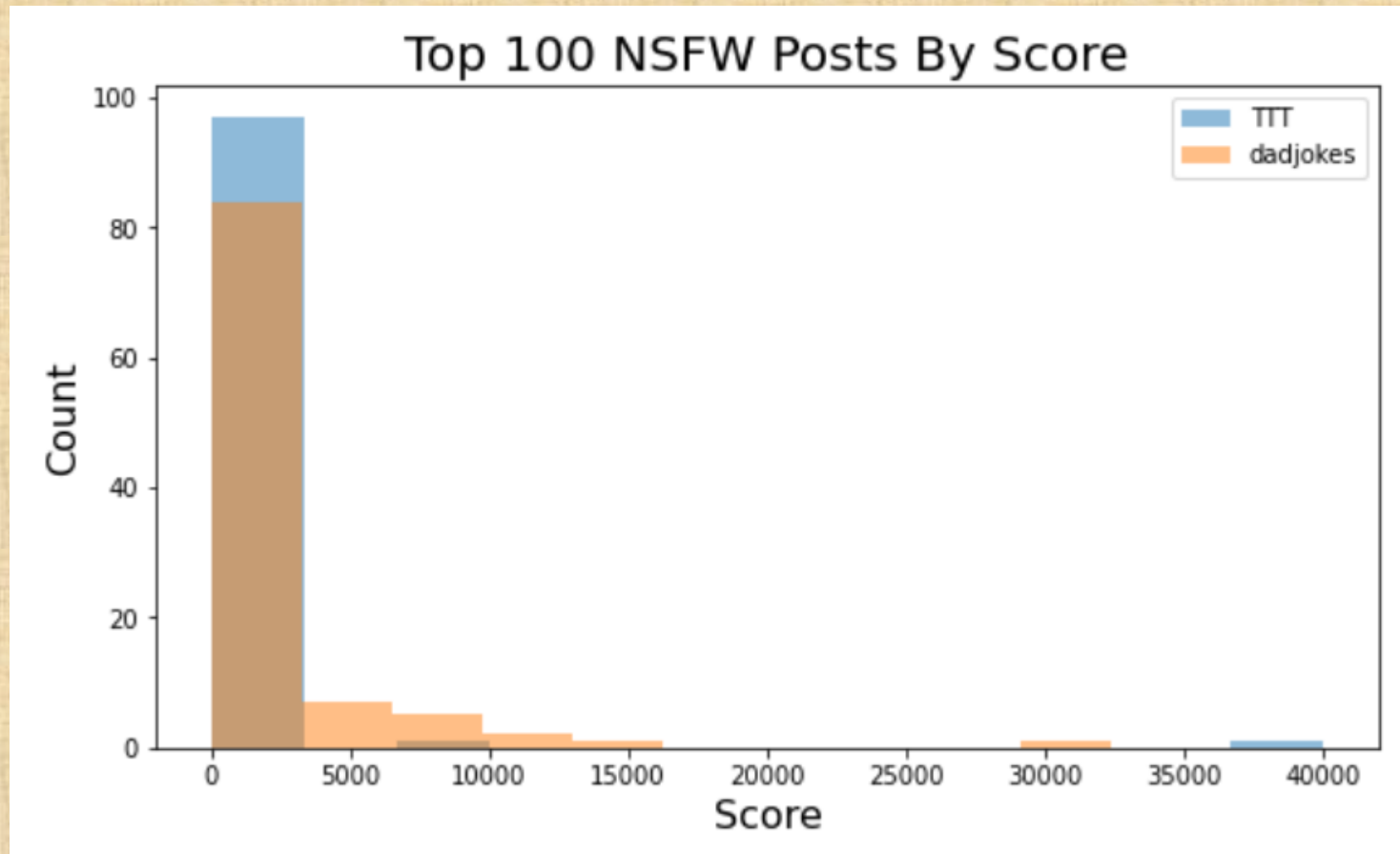
- ▶ Pushshift API
- ▶ 300 posts from TechnicallyTheTruth and dadjokes:
  - ▶ Top 100 posts by score
  - ▶ Top 100 NSFW posts by score
  - ▶ Top 100 Spoiler posts by score



# Some EDA



# NSFW



# Spoiler Outlier



15.1k

Don't open, it's a spoiler.

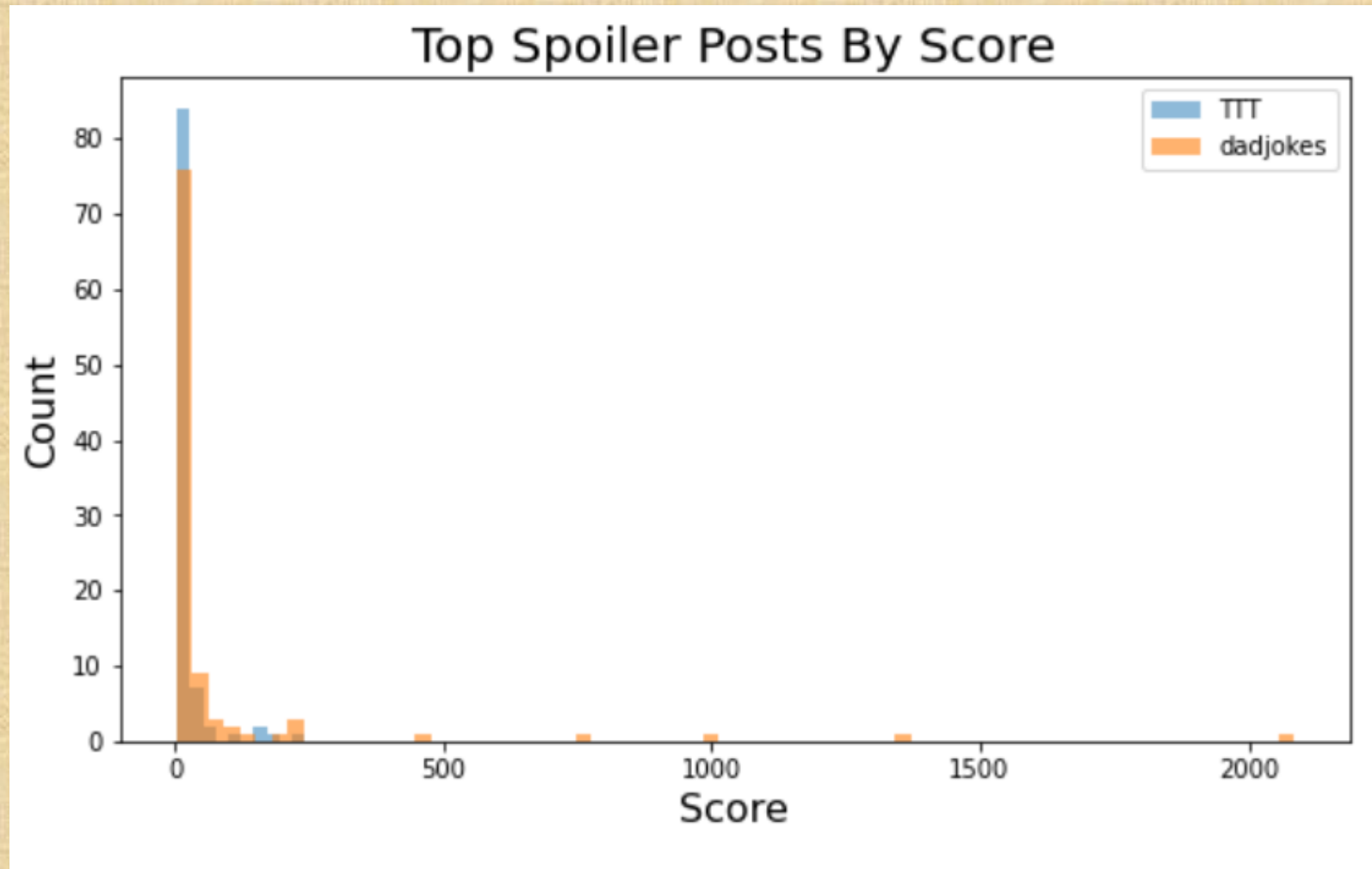


spoiler





# Spoiler



# Is\_self

This feature alone could be used to build a very accurate classifier.

Dataset	True Is_self
TTT top 100	0%
TTT NSFW	6%
TTT Spoilers	5%
Dadjoke top 100	97%
Dadjoke NSFW	100%
Dadjoke Spoilers	99%

But, that's not in the spirit of the problem



# Data Cleaning/Processing

- ▶ For posts with images of text, we used pytesseract to convert the images to strings. E.g.

I am bullet proof until proven  
otherwise

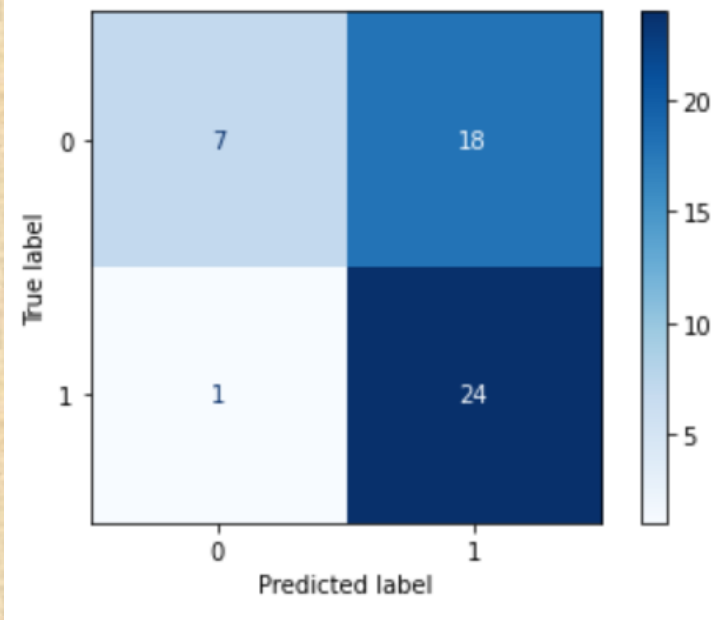
Becomes “| am bullet proof until proven \n otherwise”

- ▶ CountVectorize text for the resulting posts.
- ▶ Train/Test Split

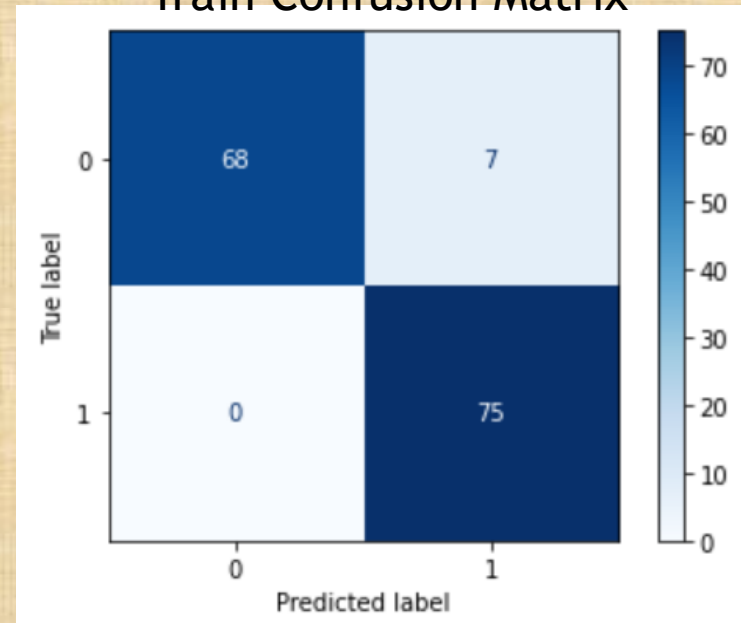
# Modeling

- ▶ Binary Classification (is\_dadjoke)
  - ▶ Logistic Regression
  - ▶ 62% accurate on test data But, 95% accurate on the training data
- ▶ Baseline model: 50% accuracy score

Test Confusion Matrix



Train Confusion Matrix



# Conclusion

- ▶ Posts on r/dadjokes and r/technicallythetruth are formatted differently
- ▶ The language used on these subreddits is different
  - ▶ Though similar
- ▶ Future research could
  - ▶ Use different data (e.g. the average, not just the hits)
  - ▶ Further process the data.
  - ▶ Use a combination of other models.



# Thank you

► Any Questions or Comments?