# How to Get Your Questions Answered Quickly

• • •

Paul Lim
05/17/2017

# Objective

# We all have questions...
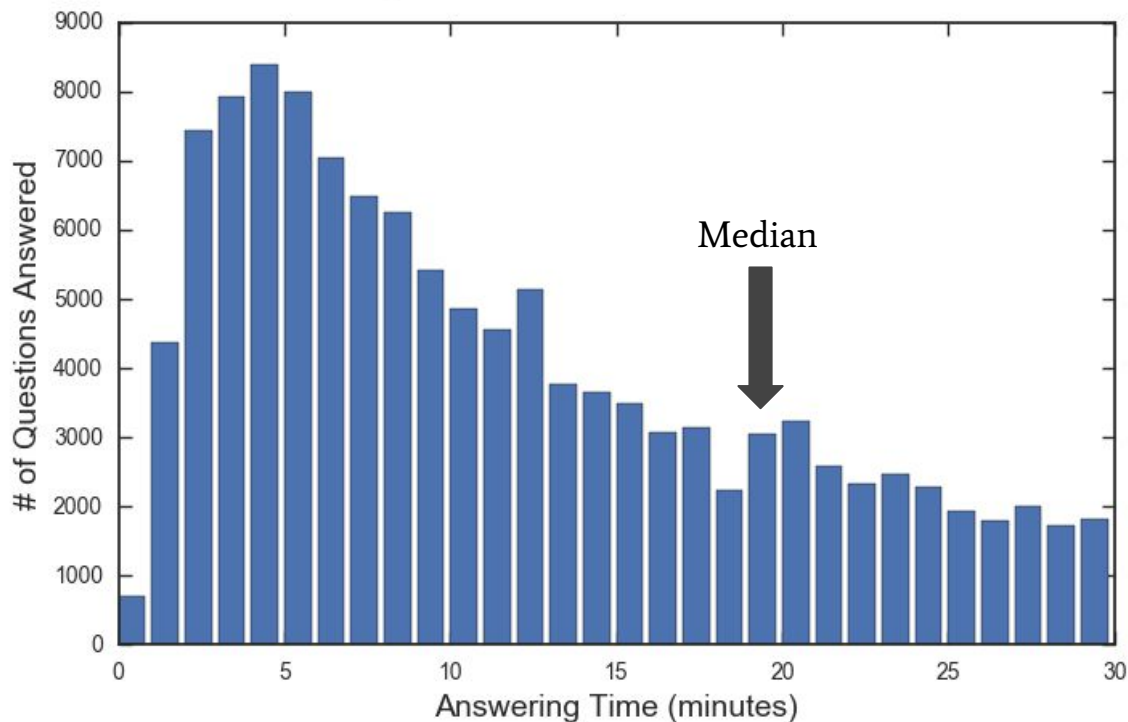
1. What features are important in getting answers?
2. Optimize the complexity of models and prediction time of new observations.

# Target Selection



Answering Time vs # of Questions Answered

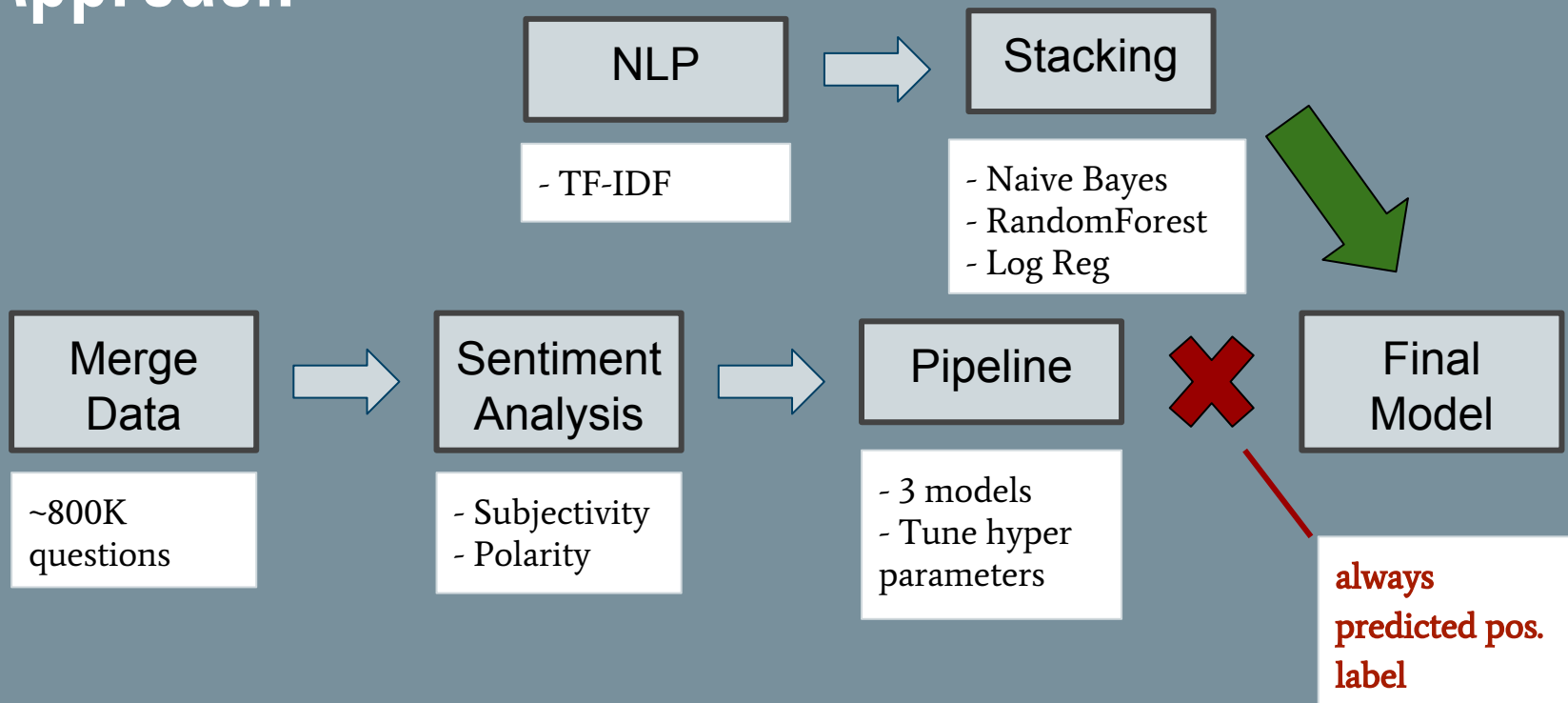- Data limited to questions answered within 24 hours

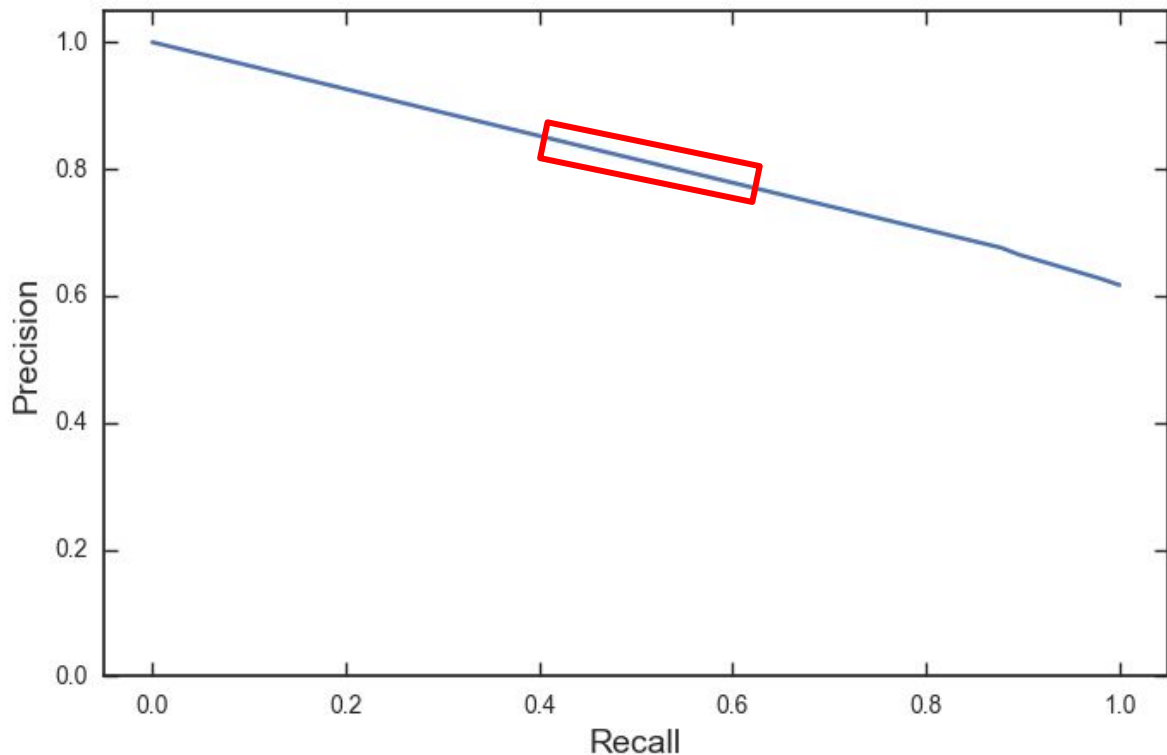| Label: | Criteria: |
|---|---|
| Quick Answer | < 30 minutes |
| Slow Answer | >= 30 minutes |

# Modeling

# Approach

# NLP → Multinomial NB → Stacking


Precision-Recall Curve

- Find the balance between precision and recall

- No more guessing only the positive label

- Limited to ~10,000 obs.

# Scores

- FBeta with a beta of 0.5 places a higher weight on precision.

- **Stacking**:

Multinomial NB + Random Forest → Logistic Regression

| Model: | FBeta: |
|--------|--------|
| Multinomial NB | 0.68 |
| Random Forest | 0.70 |
| Logistic Regression | 0.487 |
| Stacking | 0.70 |

← **Fastest model**

← Small gains for higher complexity

# Visualization

# Conclusions

# Takeaways

- According to the model, ~70% FBeta is possible.
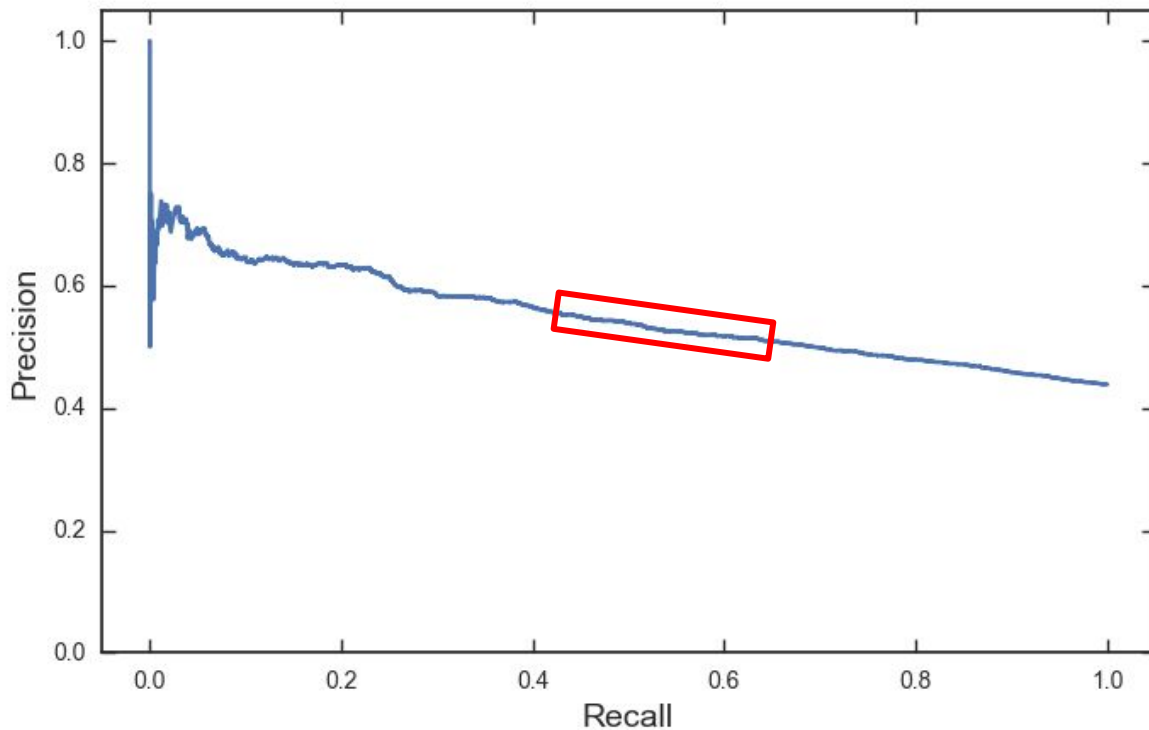- When something goes wrong.. try again!

# Future Works

# Next steps

- Different combinations of **stacking** or *boosting* for better scores
- Find a way to use all of the available data rather than a subset of it.

# Appendix
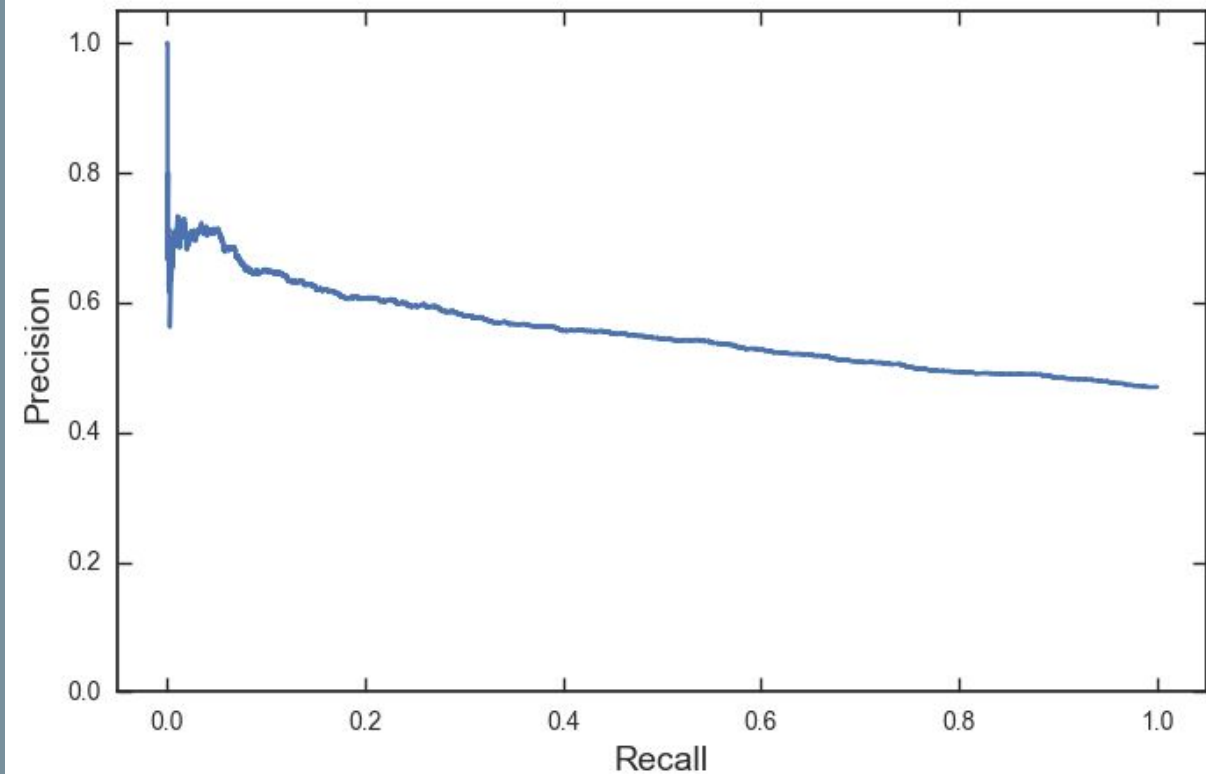
# What went wrong?


Precision-Recall Curve

- Small precision increases

- Logistic Regression model was always predicting the positive label
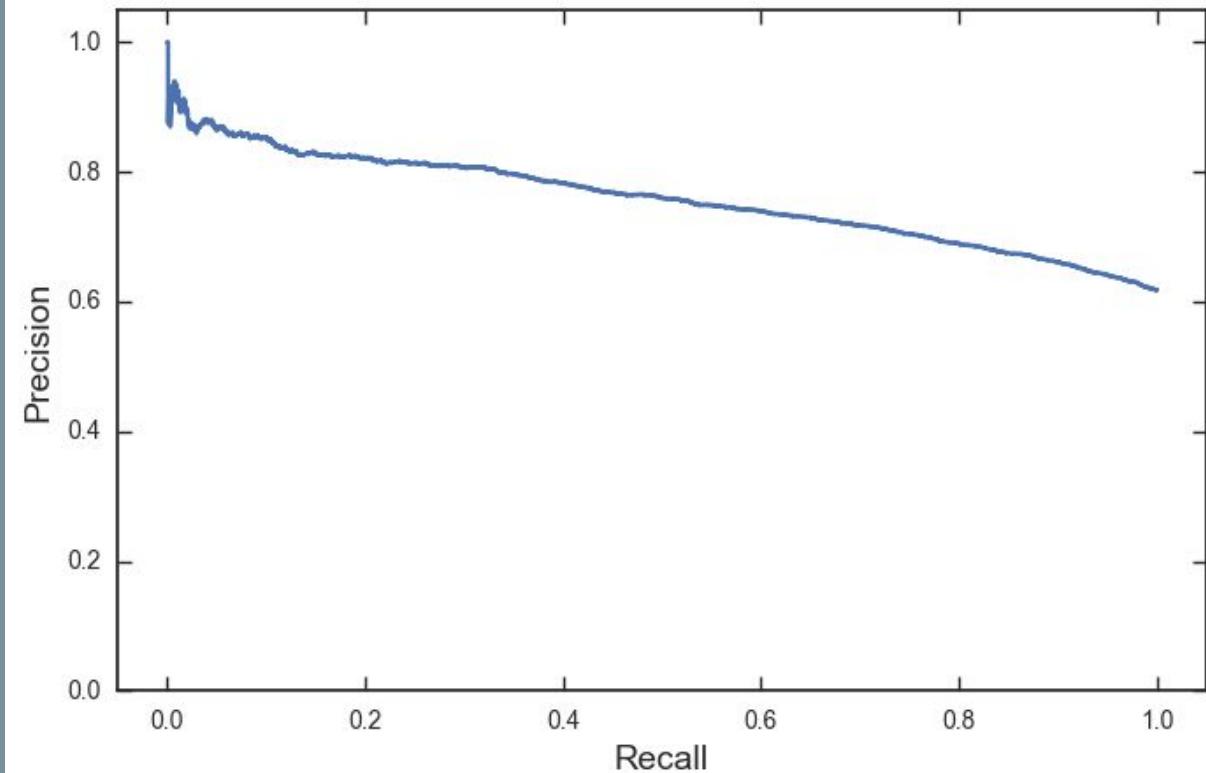
# Logistic Regression



Precision-Recall Curve

# Multinomial NB

# Multinomial NB