# How to Get Your Questions Answered Quickly

● ● ●

Paul Lim
05/17/2017

# Objective

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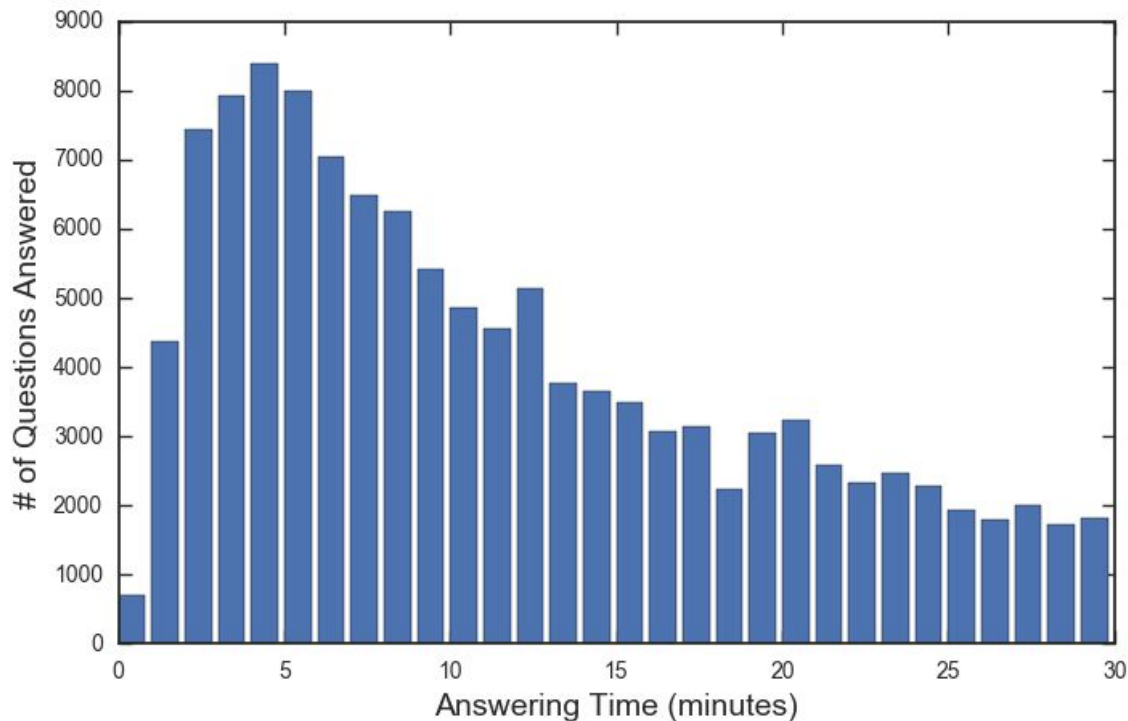
1. What features are important in getting quality answers?
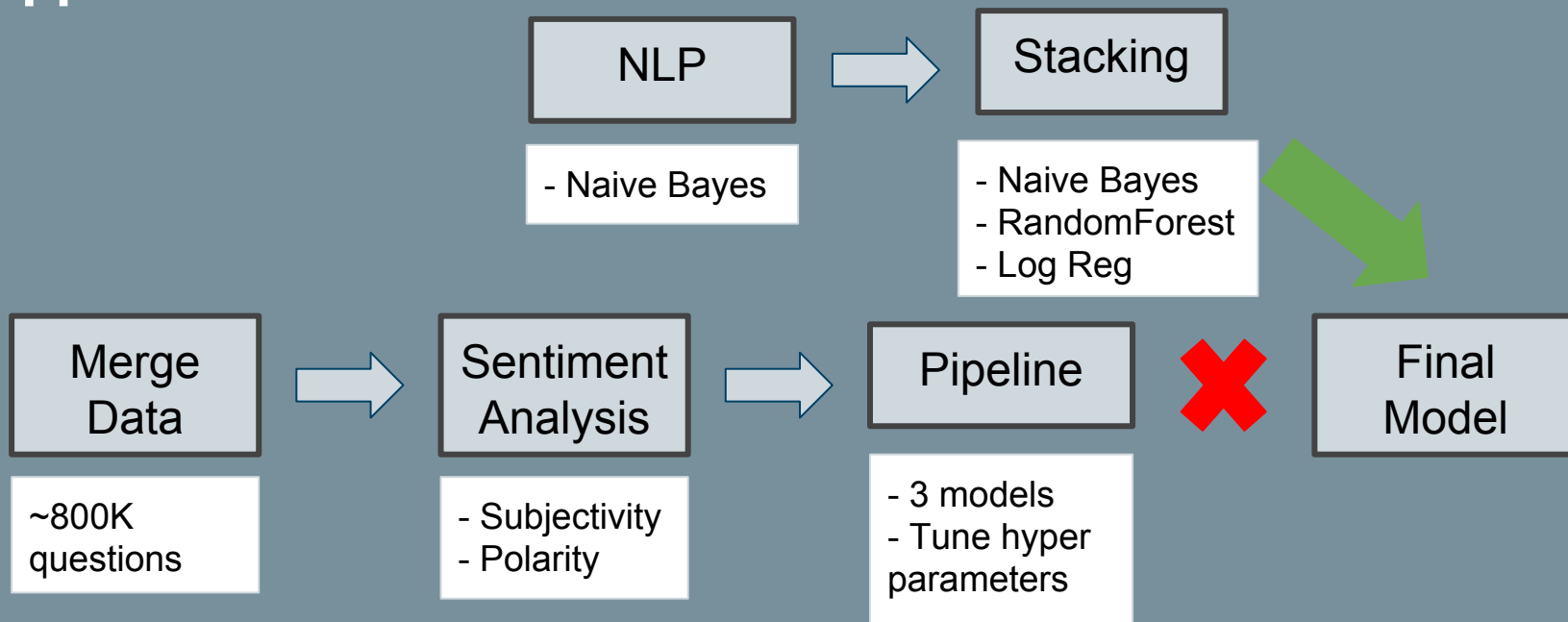2. Optimize the complexity of models and prediction time of new observations.

# Target Selection



Answering Time vs # of Questions Answered

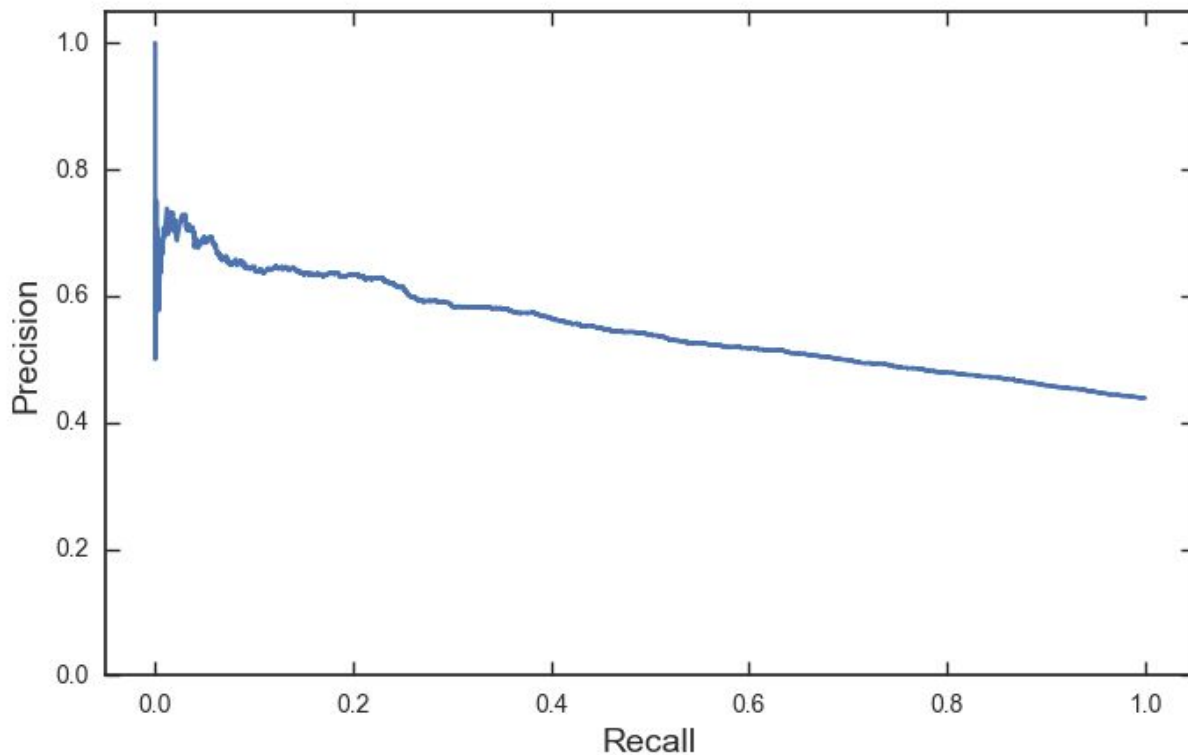- Pos. Label: < 30 minutes

- Neg. Label: >= 30 minutes

# Modeling
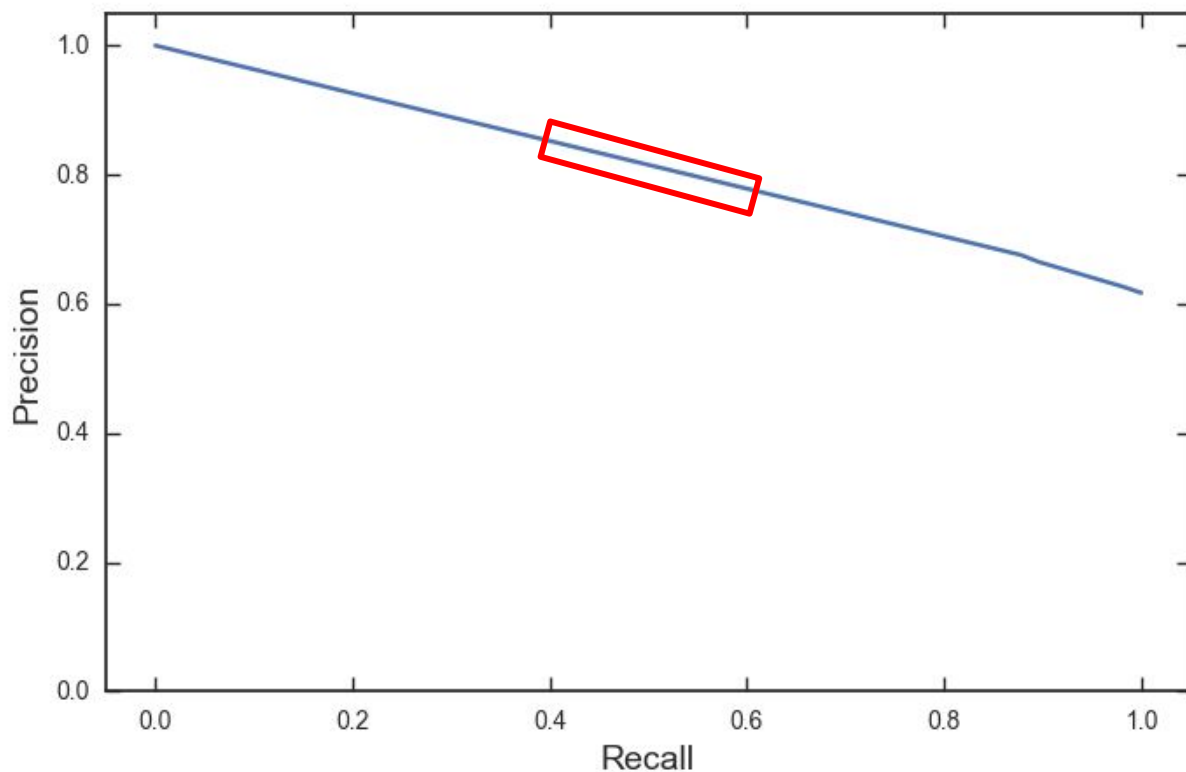
# Approach

# What went wrong?



Precision-Recall Curve

- Small precision increases

- Logistic Regression model was always predicting the positive label

# 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 + RandomForest → Logistic Regression

| Model: | FBeta: |
|---|---|
| Logistic Regression | 0.487 |
| Multinomial NB | 0.693 |
| Stacking | 0.698 |

# Visualization

# Conclusions

# Takeaways

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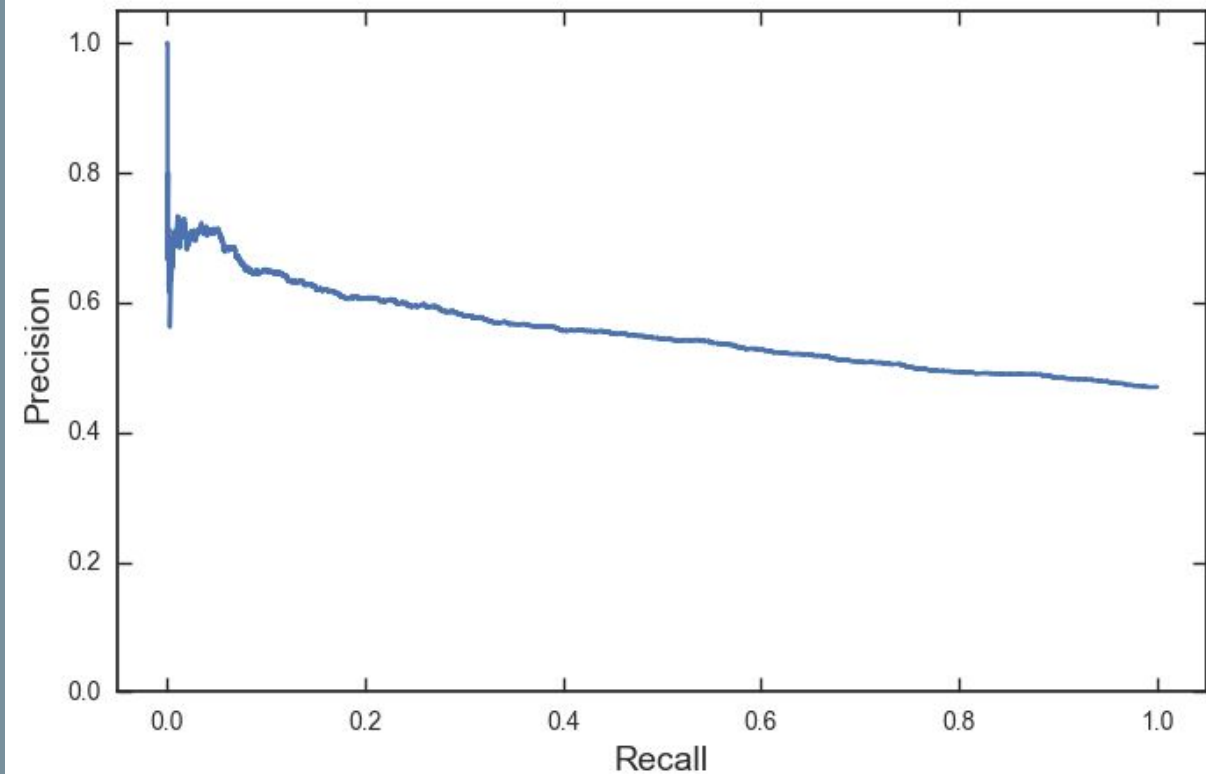- When something goes wrong.. try again!

# 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

# Logistic Regression



Precision-Recall Curve

# How to Get Your Questions Answered Quickly

● ● ●
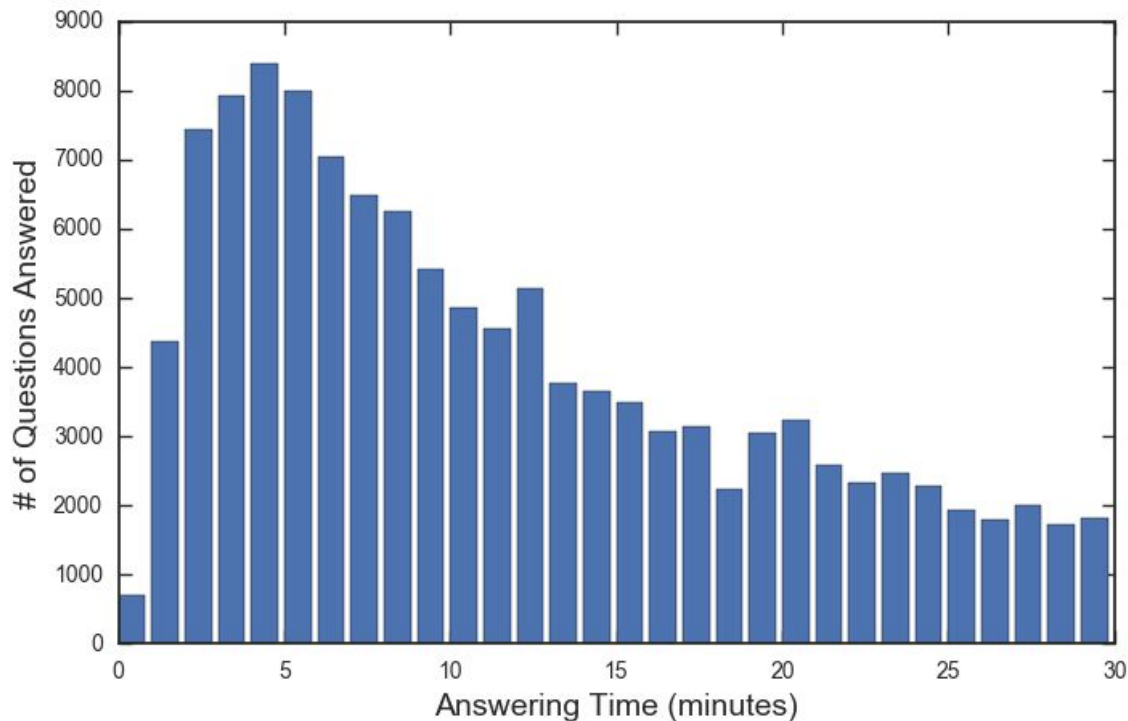
Paul Lim
05/17/2017

# Objective

# We all have questions...

1. What features are important in getting quality answers?
2. Optimize the complexity of models and prediction time of new observations.

# Target Selection



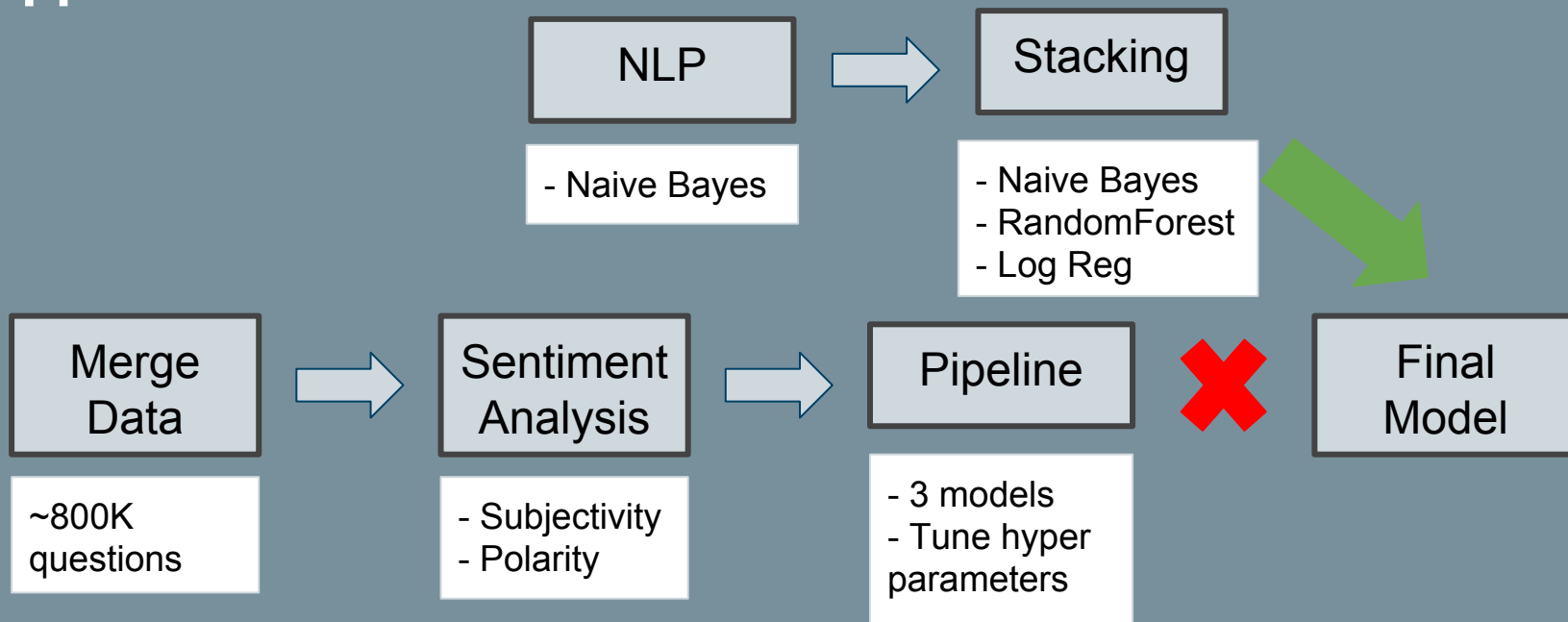Answering Time vs # of Questions Answered
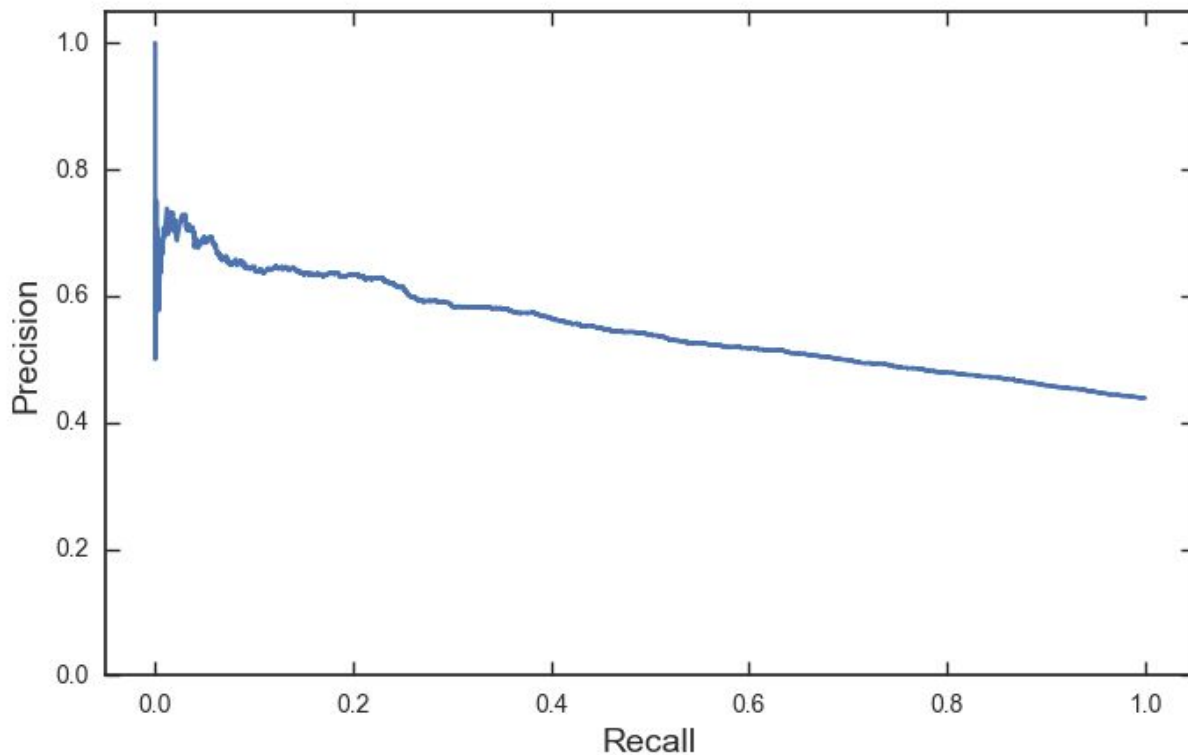
- Pos. Label: < 30 minutes

- Neg. Label: >= 30 minutes

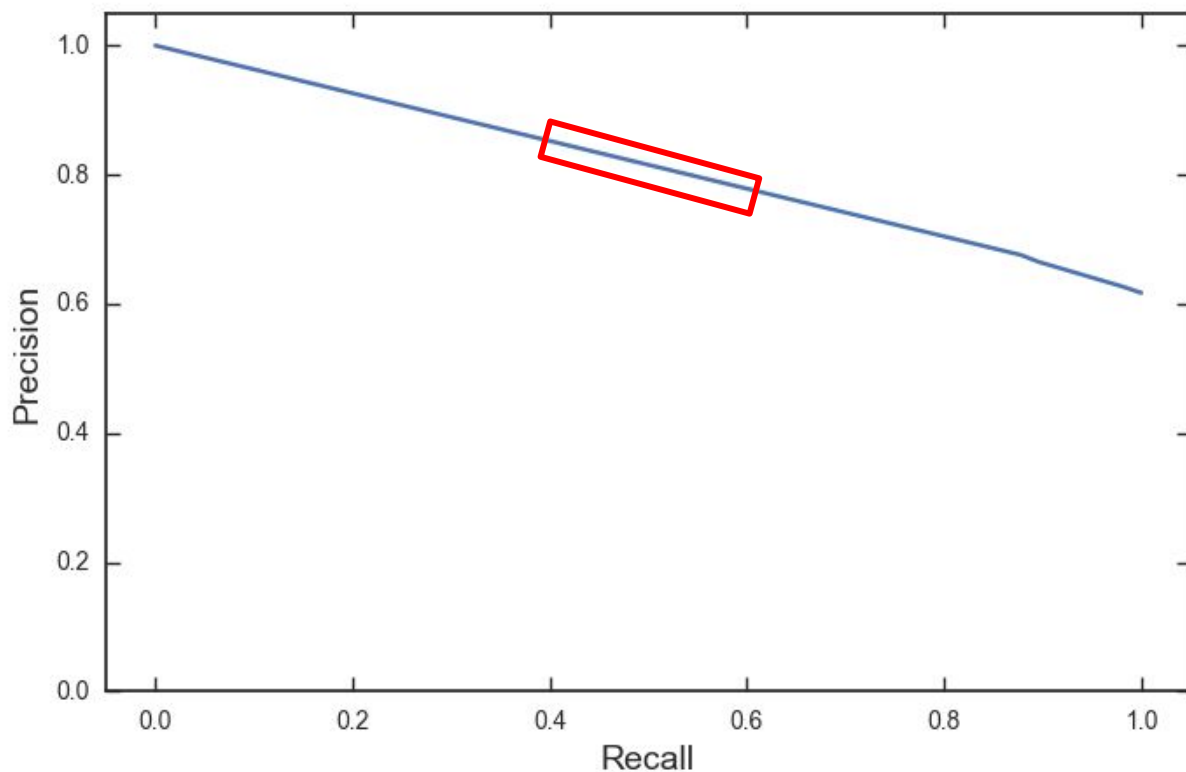# Modeling

# What went wrong?



Precision-Recall Curve

- Small precision increases

- Logistic Regression model was always predicting the positive label

# 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 + RandomForest → Logistic Regression

| Model: | FBeta: |
|---|---|
| Logistic Regression | 0.487 |
| Multinomial NB | 0.693 |
| Stacking | 0.698 |

# Conclusions

# Takeaways

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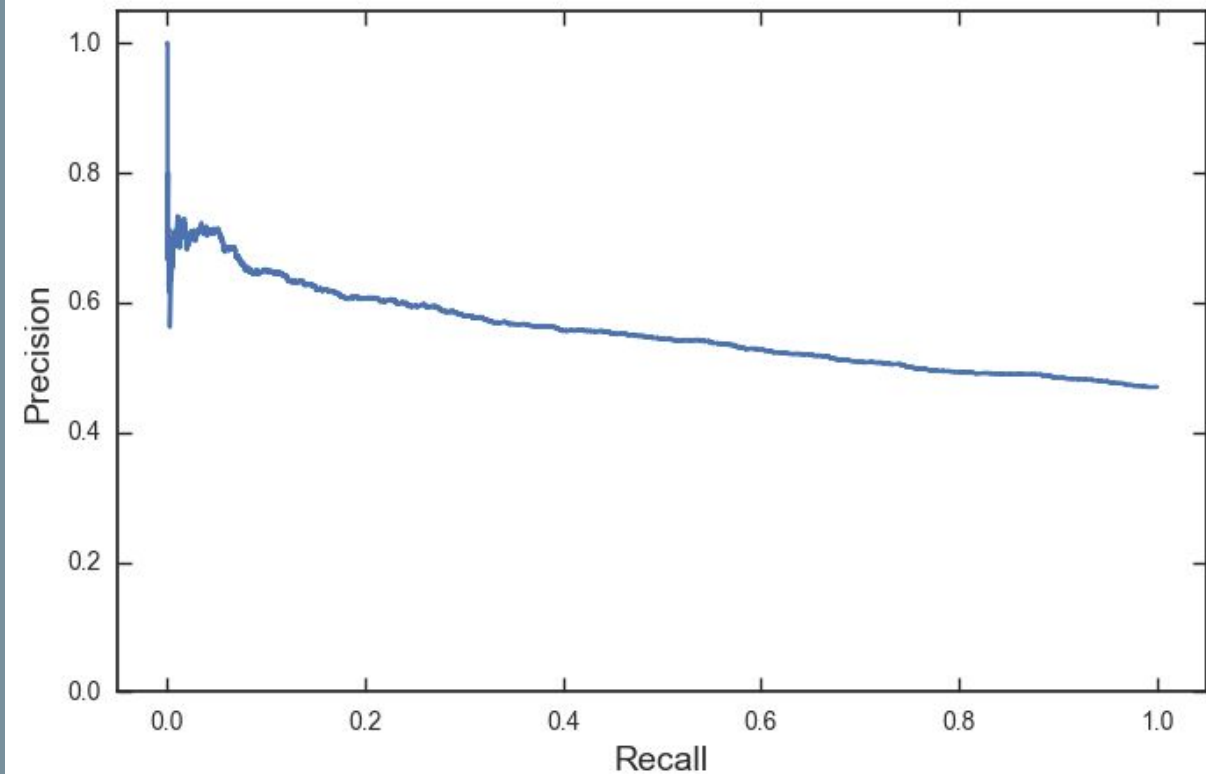- When something goes wrong.. try again!

# Future Works

# Next steps

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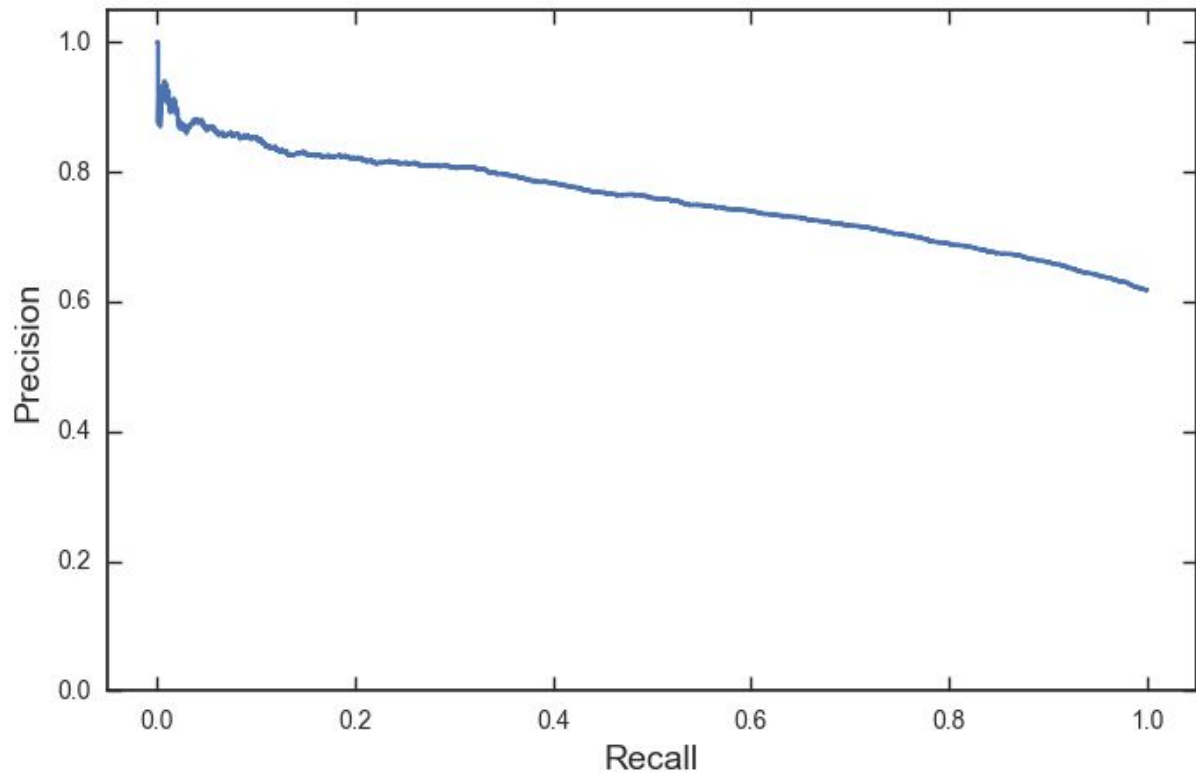- Find a way to use all of the available data rather than a subset of it.

# Appendix

# Logistic Regression

# Multinomial NB

# Multinomial NB



Precision-Recall Curve