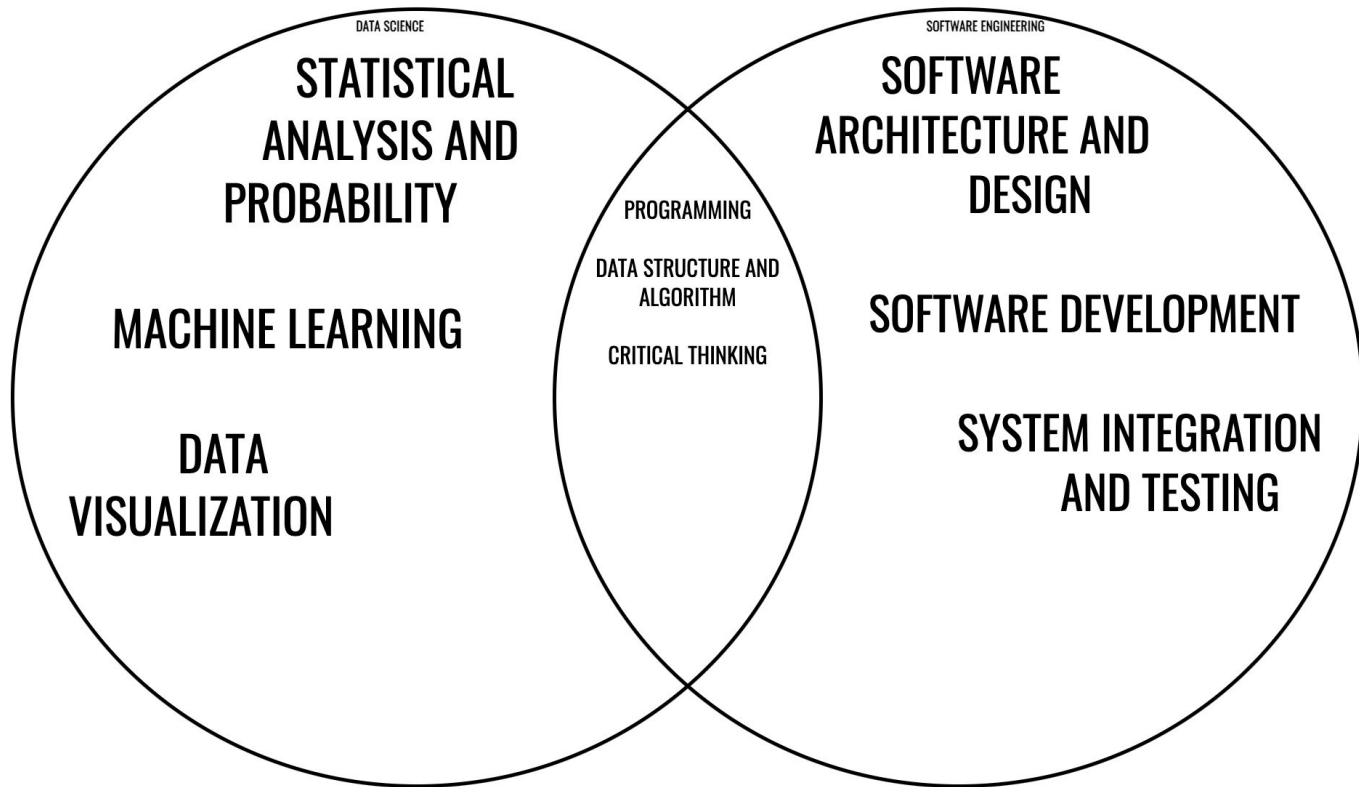


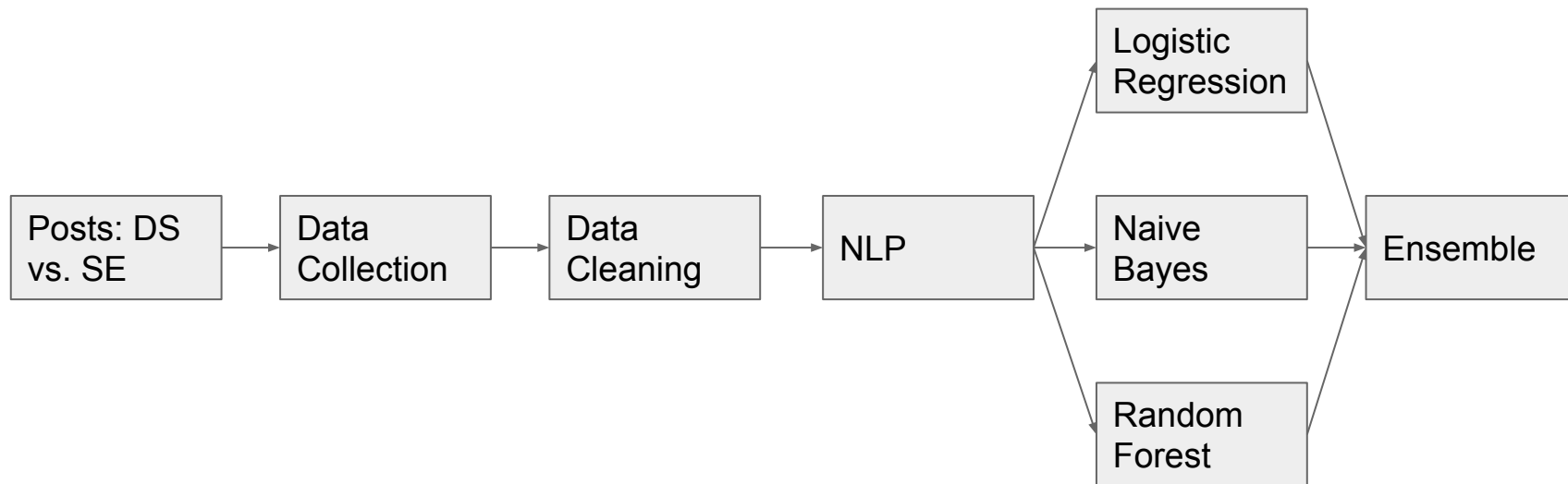
# DATA SCIENCE VS. SOFTWARE ENGINEERING: ANALYZING SUBREDDIT CONTENT WITH NLP AND CLASSIFICATION MODELS

**Ran Ma**

# BACKGROUND

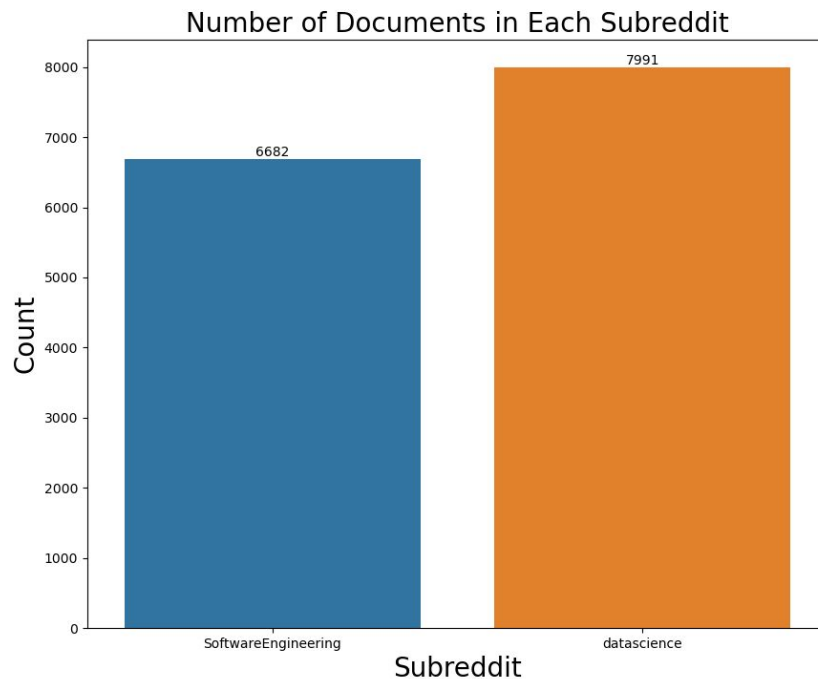


# GOAL AND WORKFLOW



# DATA COLLECTION AND CLEANING

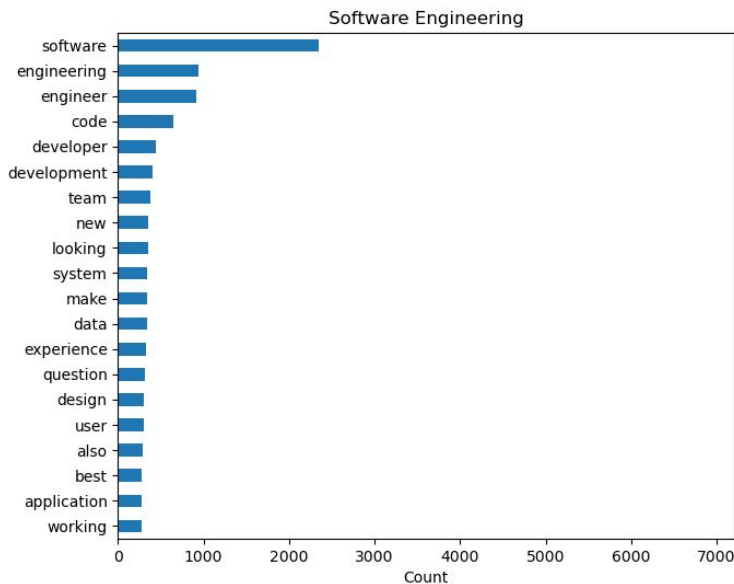
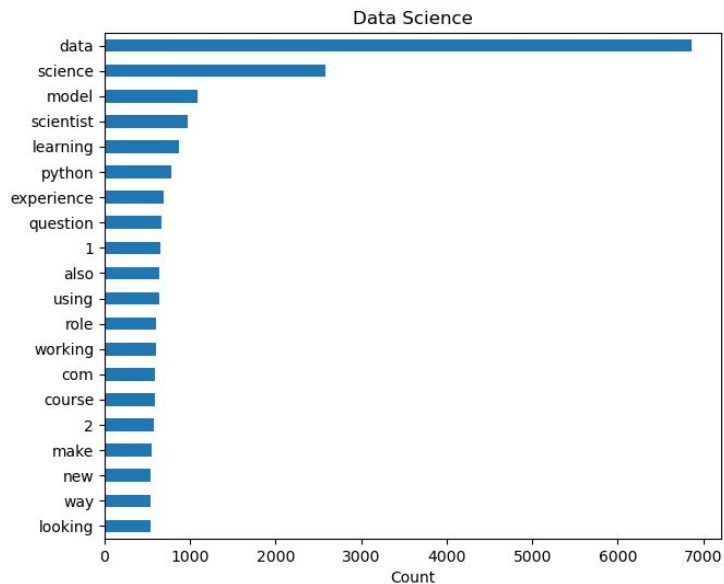
- Resources: **Pushshift API**
- Subreddit: Data Science vs. Software Engineering
- Documents: title or selftext
- Remove Missing values
- Remove Duplicate values
- Remove “[removed]” and “[deleted]”



# NLP

- Vectorization: Count Vectorize/Tf-idf Vectorize
- Applying NLTK: Stop words Lemmatize
- Retains: Letters and Digits

Top 20 Popular words in Each Subreddit



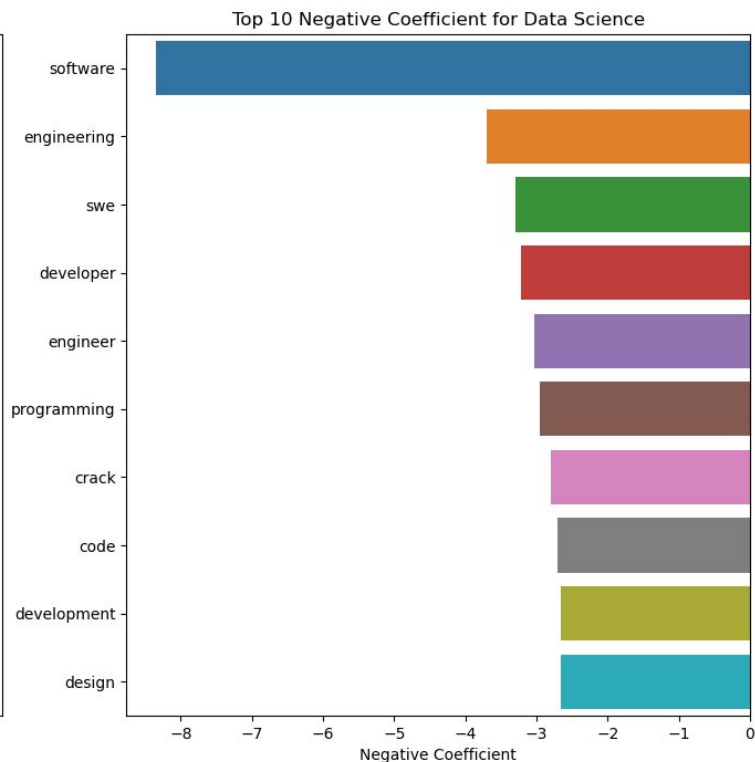
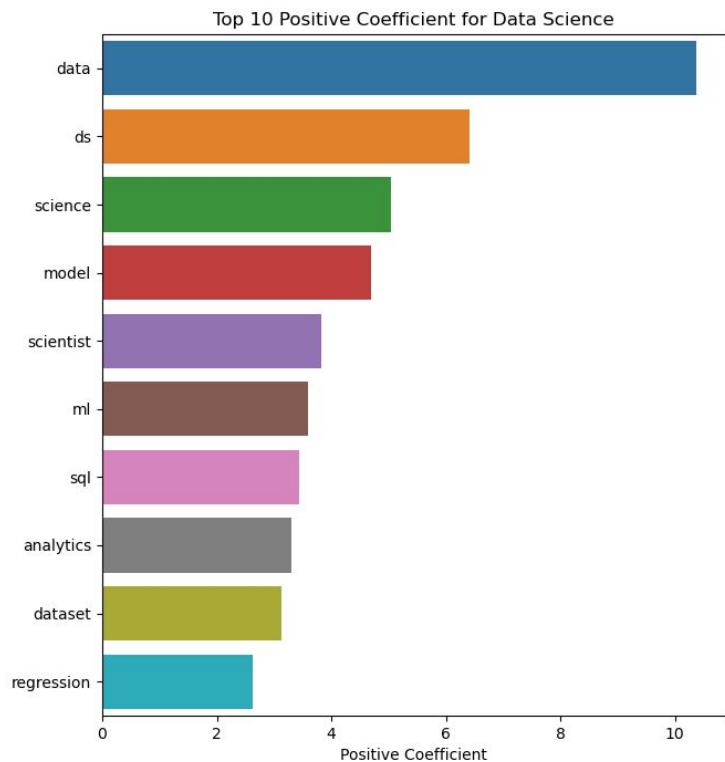
# MODEL BUILDING

- Logistic Regression
- Multiple Naive Bayes
- Random Forest
- Weighted Averaging
- Stacking

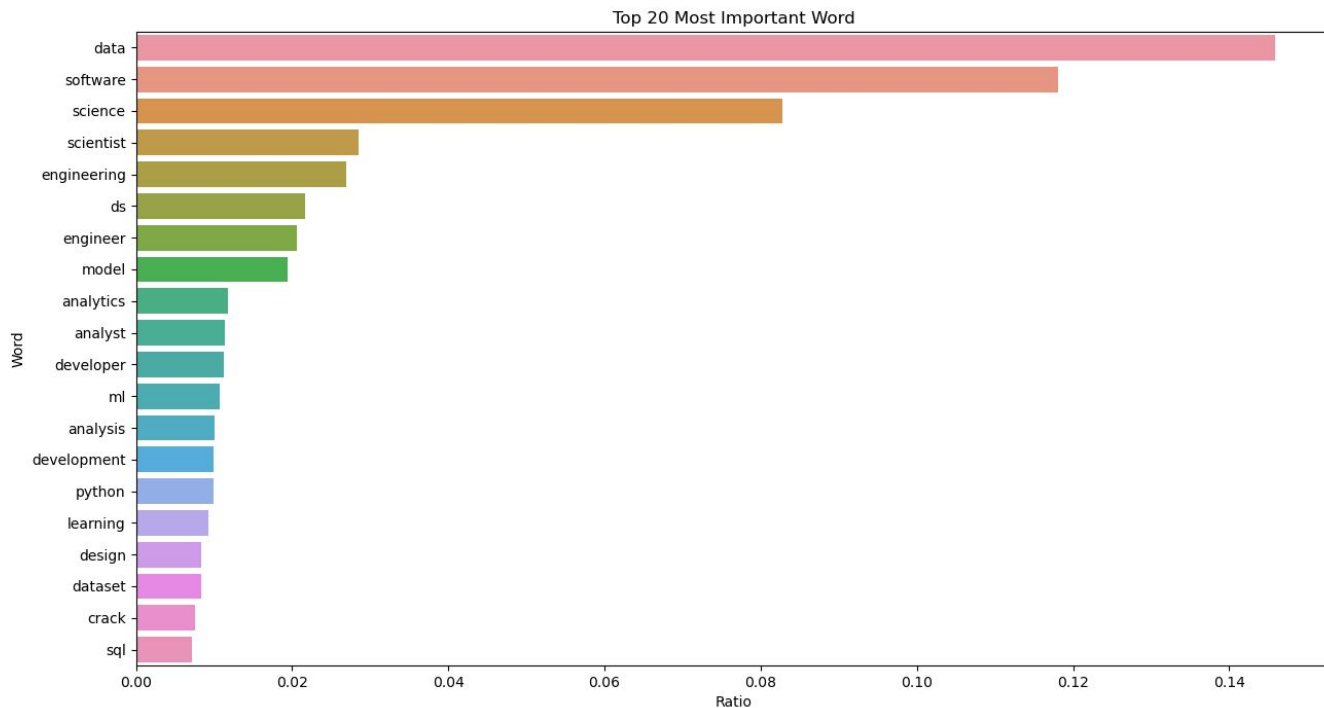
	accuracy
Logistic Regression	0.843
Multiple Naive Bayes	0.834
Random Forest	0.822
Weighted Averaging	0.846
Stacking	0.847

# LOGISTIC REGRESSION: COEFFICIENT

Coefficient in Logistic Regression



# RANDOM FOREST: IMPORTANT WORD





# CONCLUSION

1. Data Science: Model, Machine Learning, SQL, Analytics
2. Software Engineering: Software, Code, Development, Design
3. Best Accuracy: Ensemble Method

THANK YOU