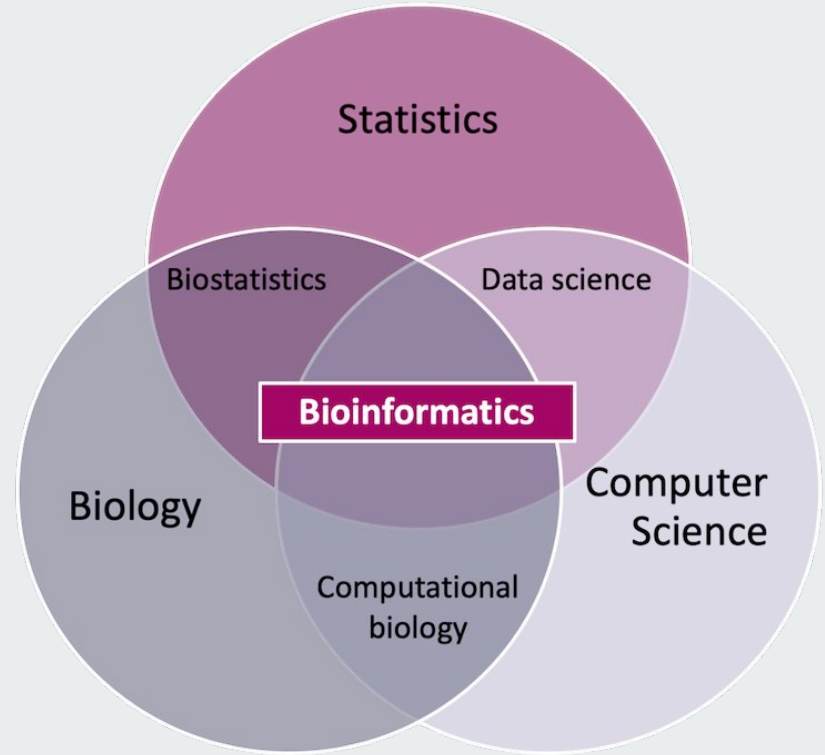


Project 3

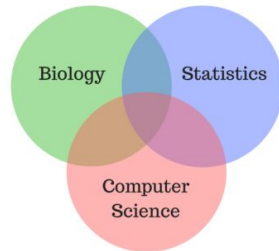
Bioinformatics and Data Science Reddits

by Tatiana Patrusheva
DSIR-1128

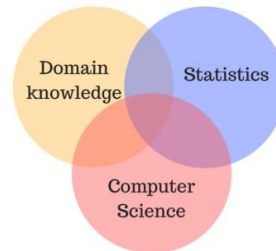


Find the model to differentiate Bioinformatics from Data Science Subreddits

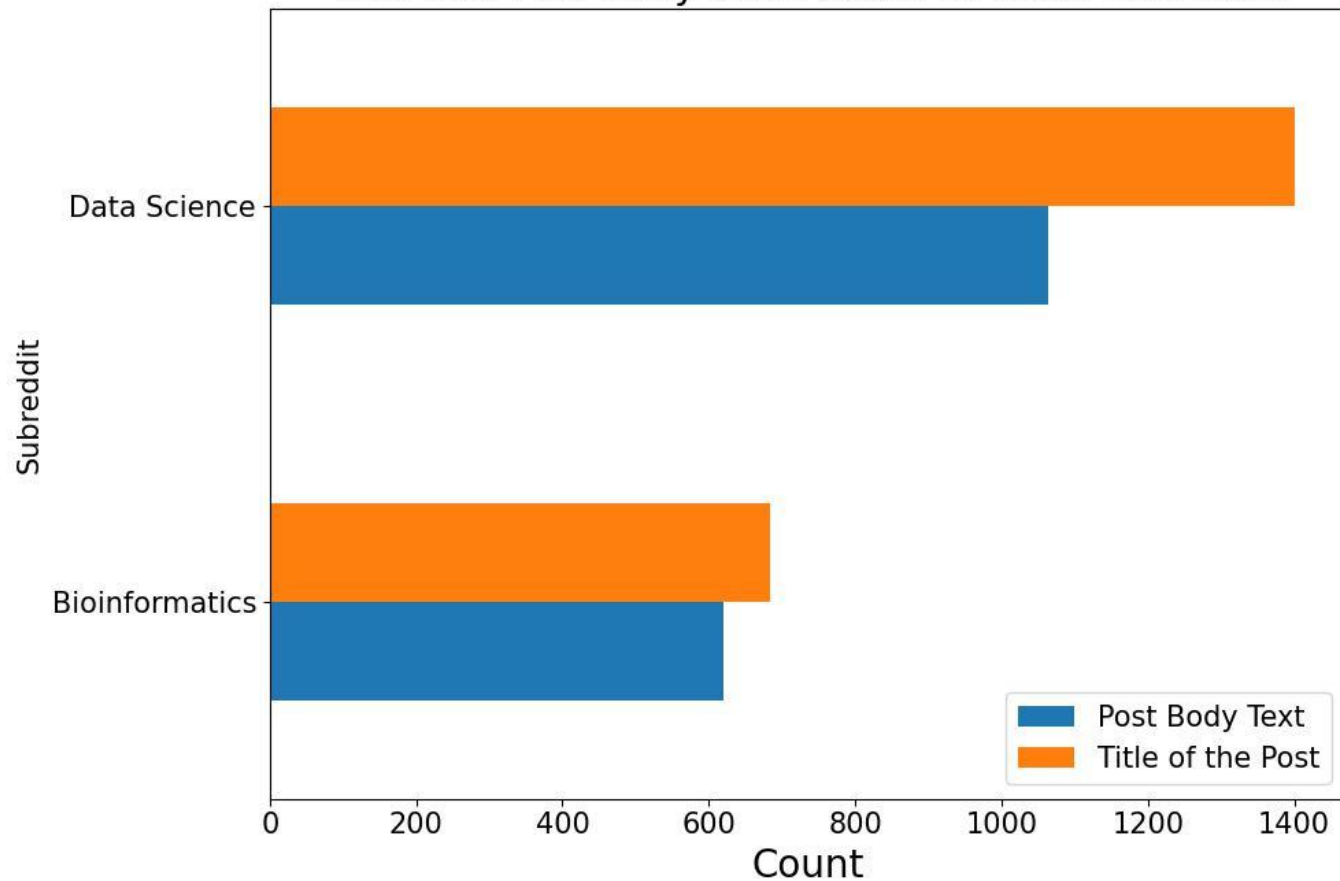
Bioinformatics



Data Science

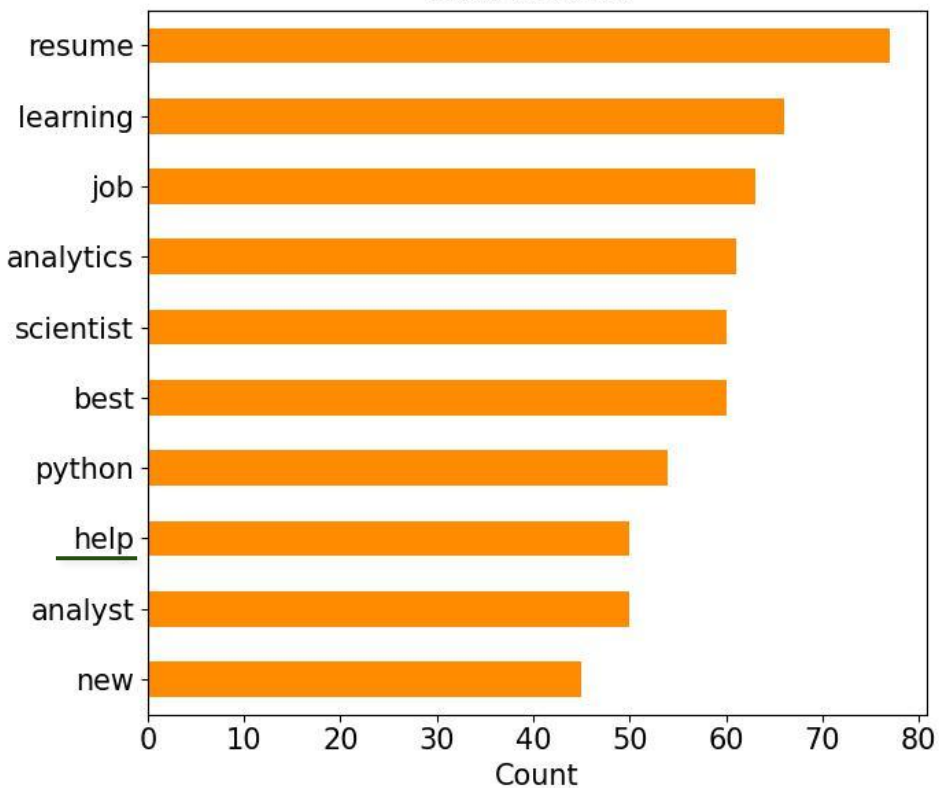


Title and Post Body texts count for each subreddit

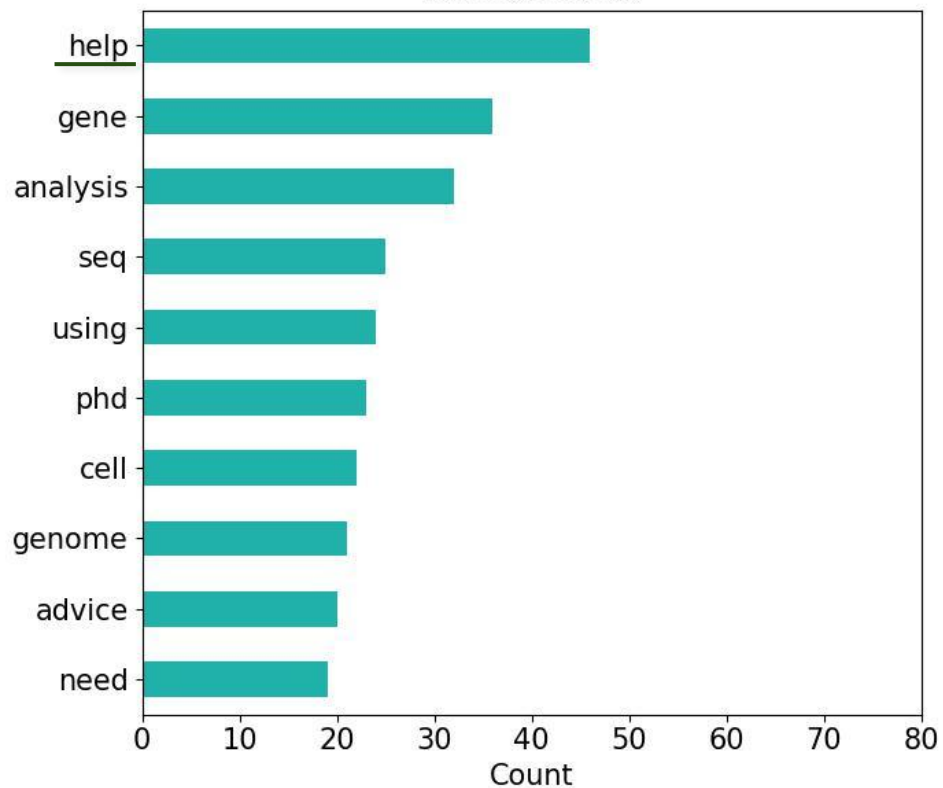


Titles Most Common Words (excluding data, science, bioinformatics, ds)

Data Science

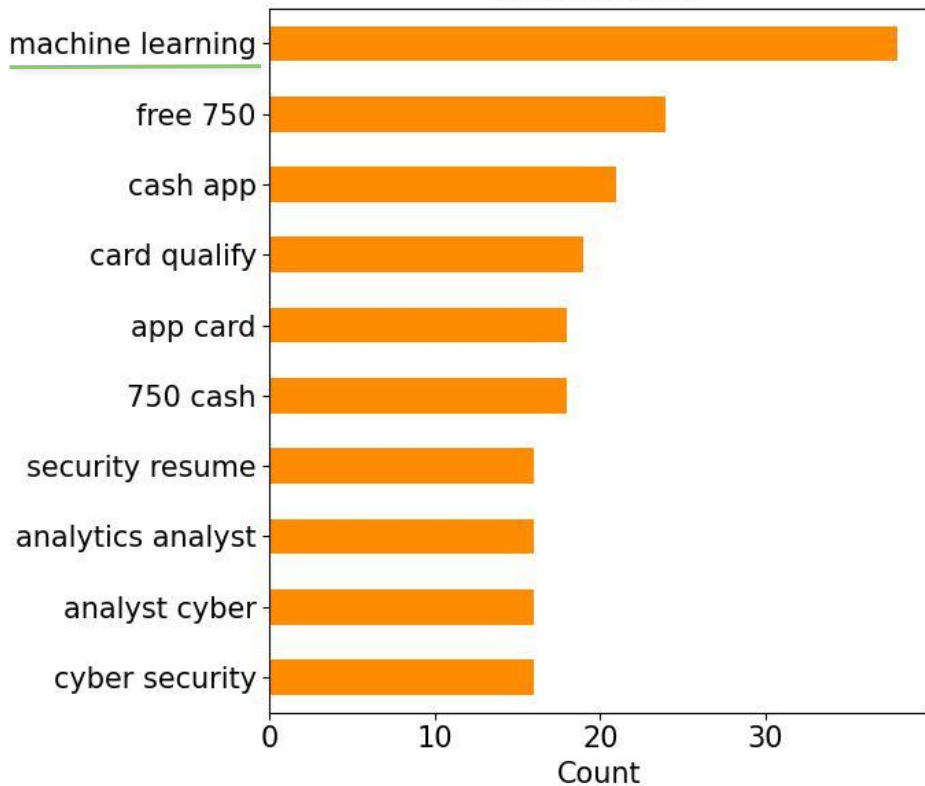


Bioinformatics

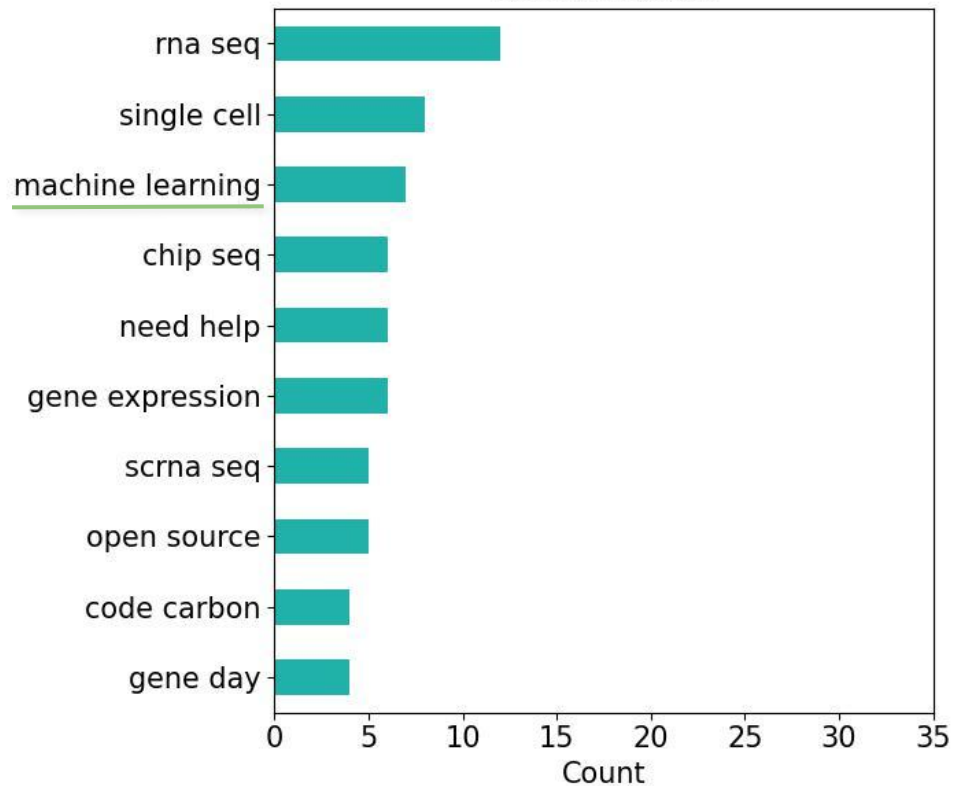


Title Most Common Bigrams (excluding data, science, bioinformatics)

Data Science

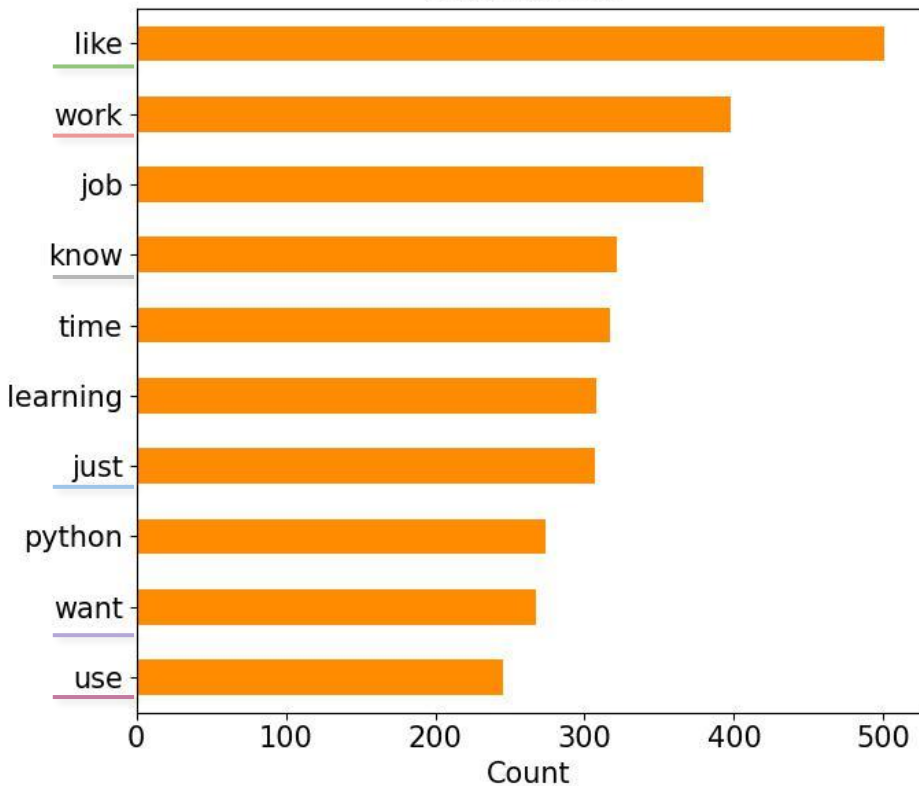


Bioinformatics

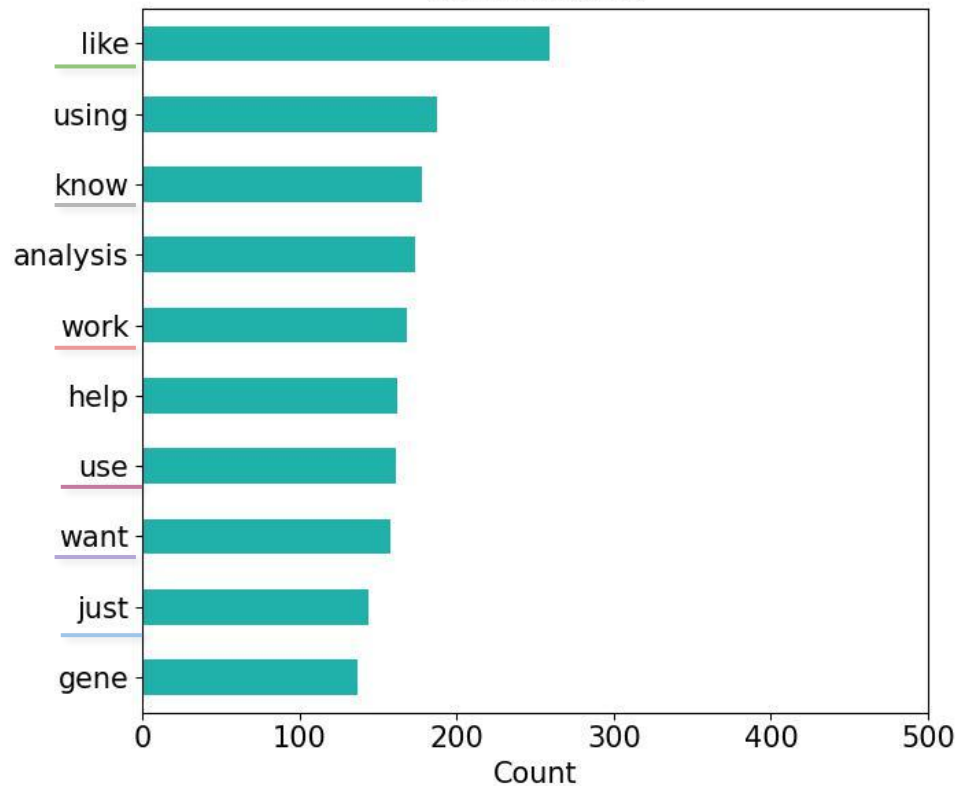


All Text Most Common Words (excluding data, science, bioinformatics)

Data Science



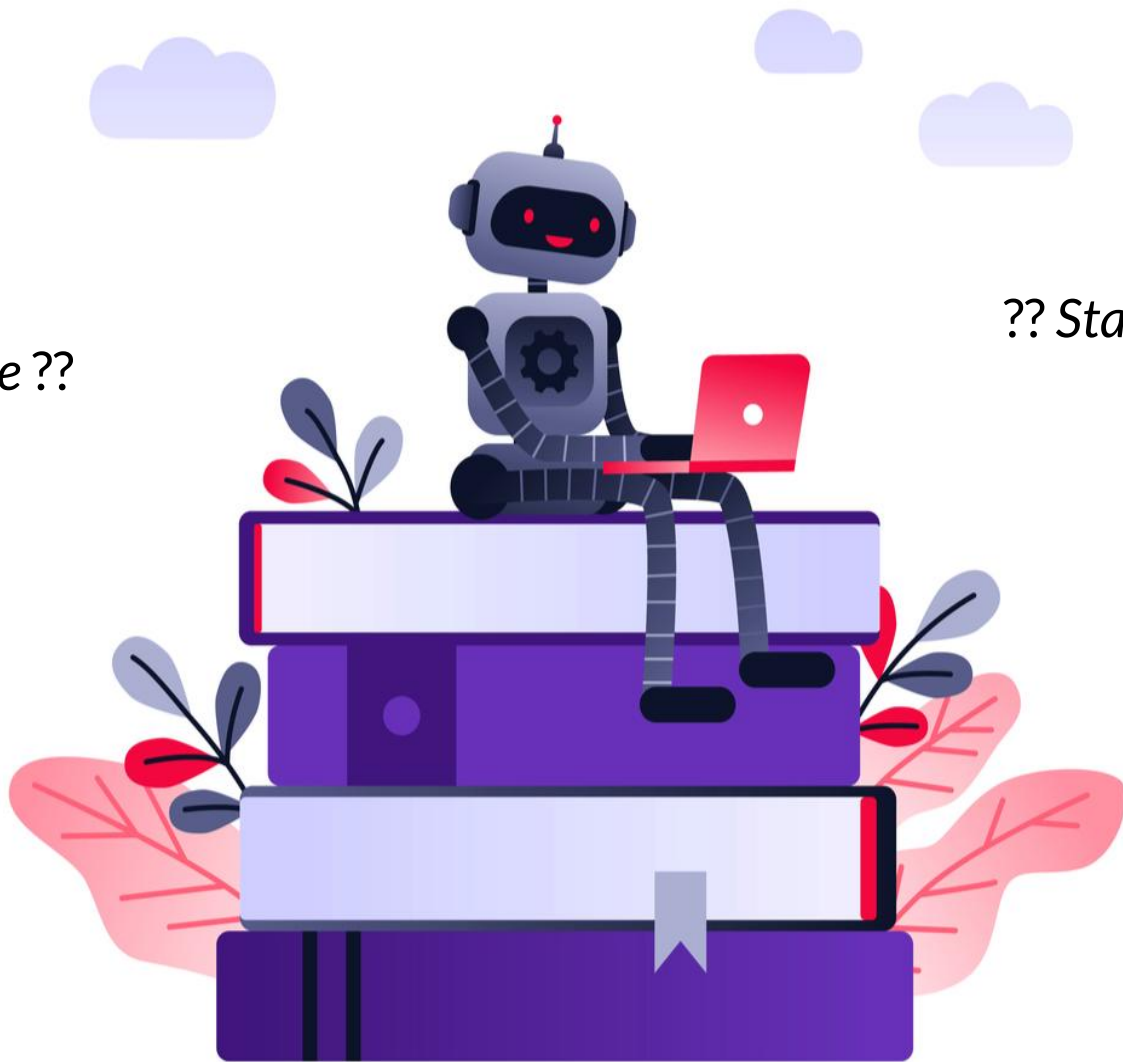
Bioinformatics



?? XGBoost ??

?? Naive Base ??

?? Stacking ??

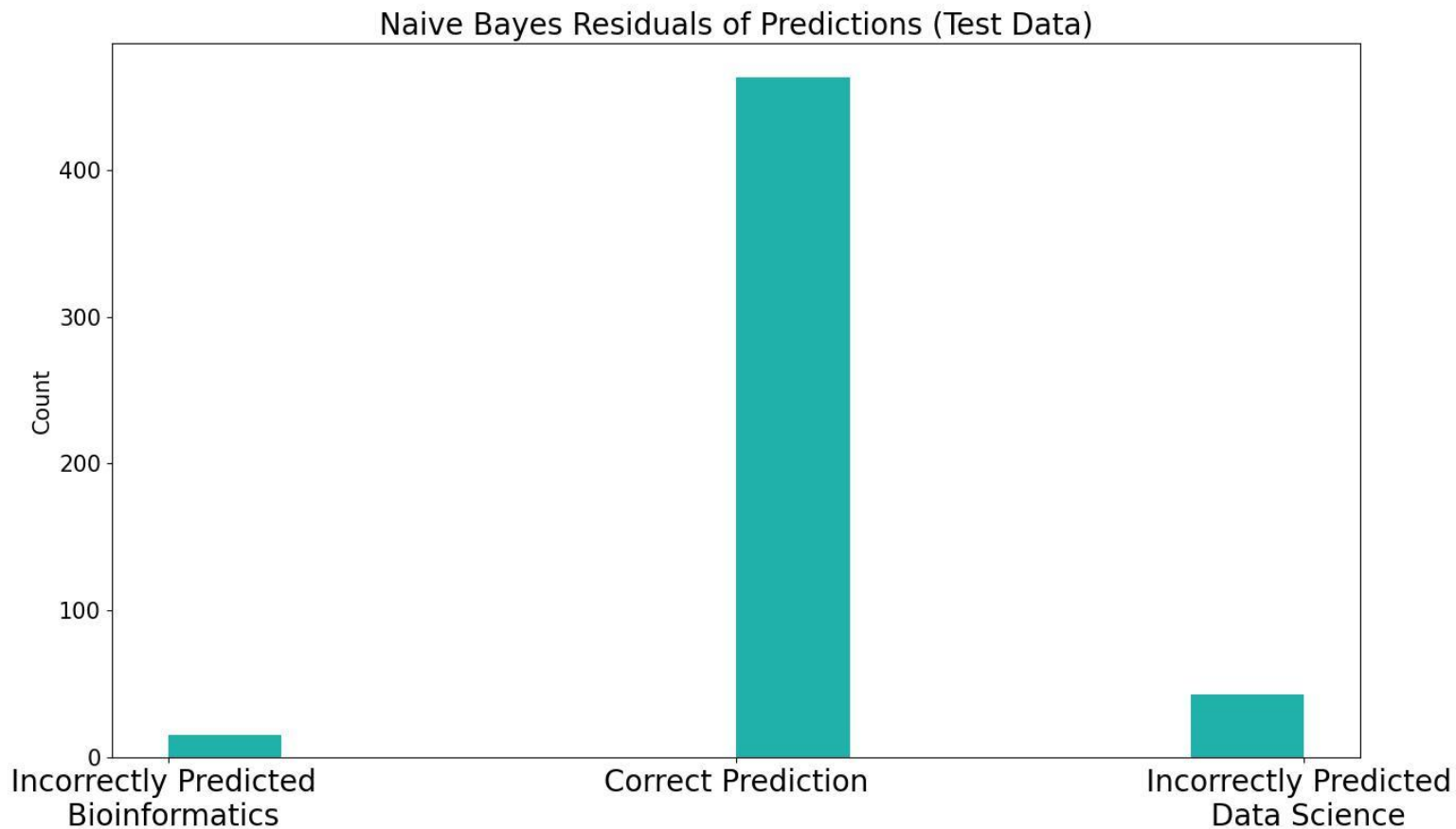


My baseline
is 67%

Naive Bayes best model

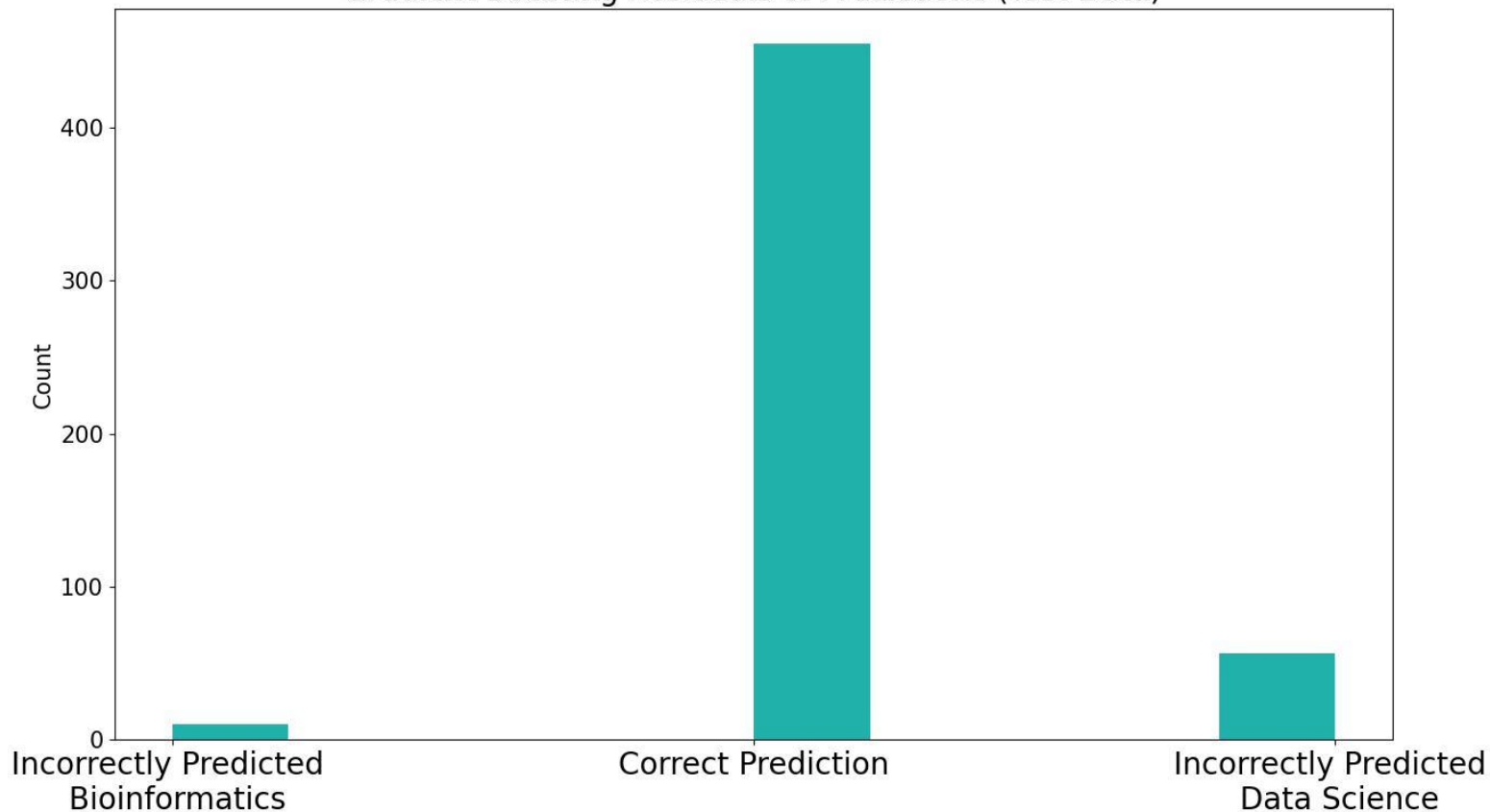
F1 score:
0.82

Test
accuracy:
89%



Gradient Boosting best model

Gradient Boosting Residuals of Predictions (Test Data)

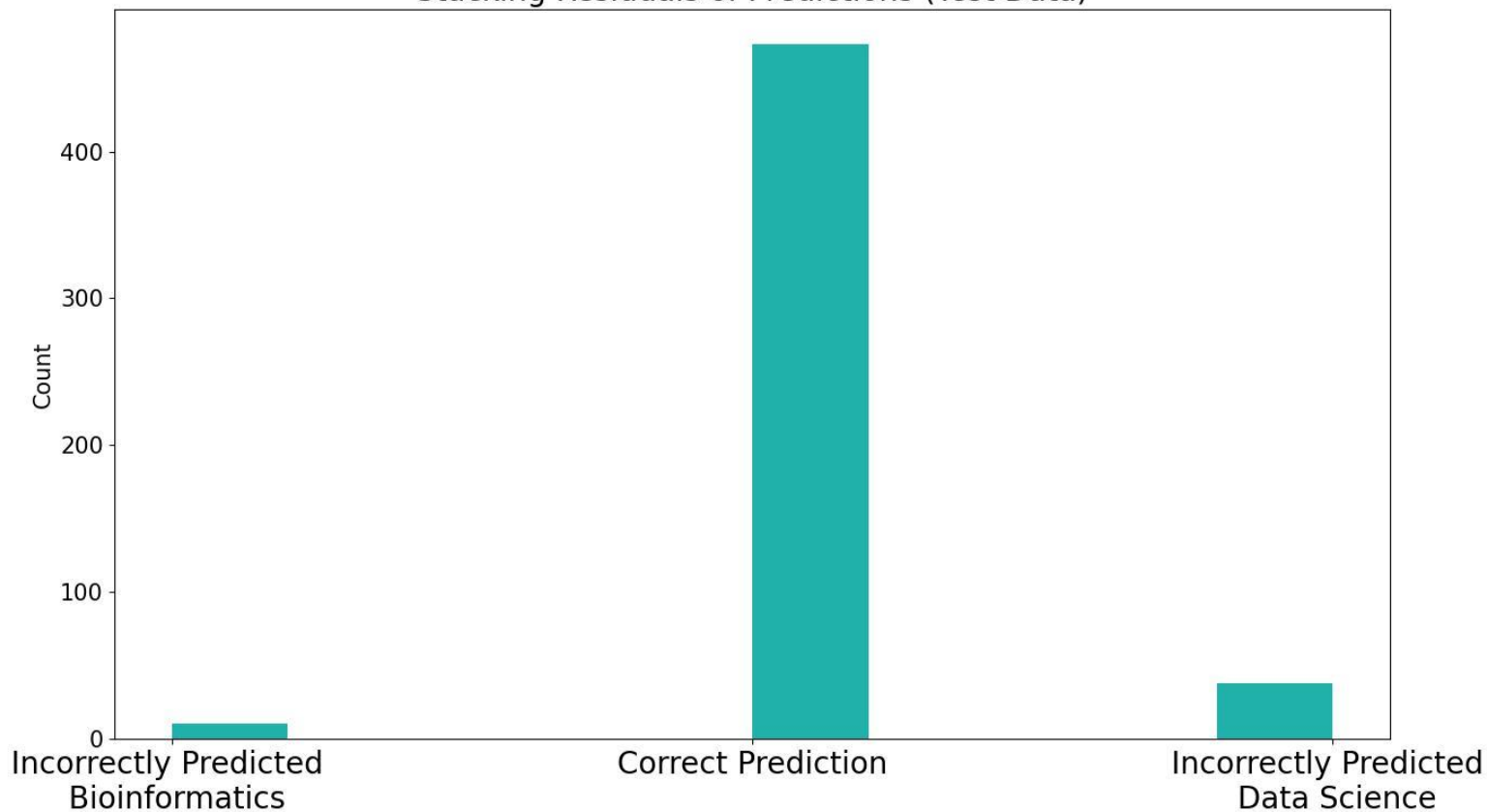


F1 score:
0.78

Test
accuracy:
87%

Stacking Best Model

Stacking Residuals of Predictions (Test Data)



F1 score:
0.85

Test
accuracy:
91%



Conclusions

Best Model - Stacking with Tokenization



Recommendations

- Data from different websites
- Balance between amount of time to fit and performance



Thank you!