

11791 Homework 4 Report

I started evaluating the performance among different models on two metrics by doing basic statistical analysis. Results are shown in the following table. As can be observed from the table, the baseline model has the highest mean, median and mode value for both metrics.

metric	model	mean	median	mode	min	max
ROUGE-2	Baseline	0.4165538	0.33333	1 (count=12)	0	1
ROUGE-2	Baseline+Fusion	0.3029308	0.218255	0 (count=17)	0	1
ROUGE-2	Baseline+Ordering	0.4124338	0.318015	0 (count=19)	0	1
ROUGE-2	Baseline+Ordering+Fusion	0.2929442	0.207615	0 (count=17)	0	1
ROUGE-SU4	Baseline	0.4165538	0.33333	1 (count=12)	0	1
ROUGE-SU4	Baseline+Fusion	0.3029308	0.218255	0 (count=17)	0	1
ROUGE-SU4	Baseline+Ordering	0.4124338	0.318015	0 (count=19)	0	1
ROUGE-SU4	Baseline+Ordering+Fusion	0.2929442	0.207615	0 (count=17)	0	1

I then calculated p-values with the T-Test, Wilcoxon test and KS test on each model and metric. The following table displays the results using red to indicate p-values ≥ 0.05 and green to indicate p-values < 0.05 . The results show that improved models are not necessarily better than the baseline model.

Specifically, the three models (fusion, ordering, ordering+fusion) are not significantly better than the baseline model in all three tests. In both metrics, the improvement of the ordering model is not significant as measured by all three tests. The fusion and ordering+fusion model are significantly better than the baseline model in T-test and Wilcoxon test but not significant in the KS test. When comparing the fusion and the ordering model, the ordering model is significantly better than the baseline fusion model in the T-test and Wilcoxon test. The ordering+fusion model is not significantly better than the fusion model in all metrics and significance measures. The ordering+fusion model is significantly better in all three significance measures but is not marked as better in the mean diff metrics.

metric	model	mean diff	P(T test)	P(wilcoxon test)	P(ks test)
ROUGE-2	Baseline & Fusion	-0.1136230000000009	0.010449188538310899	0.0065922694080604635	0.13997518016385846
ROUGE-2	Baseline & Ordering	-0.004119999999999957	0.93451403694589597	0.62148994888870157	0.19304165192469011

ROUGE-2	Baseline & Ordering+Fusion	-0.1236096000 000001	0.00510505601 23188041	0.0044566094455 675952	0.069092434889 398091
ROUGE-2	Fusion & Ordering	0.1095030000 0000013	0.01870174963 7786994	0.0056470445023 683548	0.069092434889 398091
ROUGE-2	Fusion & Ordering+Fusion	-0.0099866000 000000121	0.79877547262 841786	0.7533456839567 4884	0.999996898897 31391
ROUGE-2	Ordering & Ordering+Fusion	-0.1194896000 0000014	0.00993021465 04130131	0.0009836990976 1899317	0.031376652153 072448
ROUGE-SU4	Baseline & Fusion	-0.1136230000 0000009	0.01044918853 8310899	0.0065922694080 604635	0.139975180163 85846
ROUGE-SU4	Baseline & Ordering	-0.0041199999 99999957	0.93451403694 589597	0.6214899488887 0157	0.193041651924 69011
ROUGE-SU4	Baseline & Ordering+Fusion	-0.1236096000 000001	0.00510505601 23188041	0.0044566094455 675952	0.069092434889 398091
ROUGE-SU4	Fusion & Ordering	0.1095030000 0000013	0.01870174963 7786994	0.0056470445023 683548	0.069092434889 398091
ROUGE-SU4	Fusion & Ordering+Fusion	-0.0099866000 000000121	0.79877547262 841786	0.7533456839567 4884	0.999996898897 31391
ROUGE-SU4	Ordering & Ordering+Fusion	-0.1194896000 0000014	0.00993021465 04130131	0.0009836990976 1899317	0.031376652153 072448

For the implementation, I basically followed the starter code but refactored the way to fill in the table. I used the starter code in writeOutput function as a new function called `init_table` that fills in basic headers, row names, etc. I added a function called `stat_analysis` that calculates the basic statistical metrics and stores the values as class members. The function `fill_results` evaluates on the four metrics (mean diff, ttest, wilcoxon test and ks test) and fills in the table. To ease the process of filling the table, I created another function called `cross_evaluate` that does cross evaluations on each pair of the models and fills the table with respective results. The final `write_output` function just writes the whole table into the file. I simply used the existing packages in Scipy to do the p-value tests.