

# Practical 1: Draft

Antonio Copete  
acopete@cfa.harvard.edu  
Camelot.ai username: Copete

February 3, 2018

## Notes

- 256 predictors. Linear Regression  $R^2 = 0.46$ , Random Forest  $R^2 = 0.55$ .
- 194 features are not significant in LR (coef = 0.0), and 225 in RF (importance = 0.0). All 194 unimportant LR predictors are also unimportant in RF. Drop them.
- Testing with only the 31 features that have RF importance  $> 0.0$  gives the same  $R^2$  in both cases, up to 5 significant digits. Keep working with these features only. Also, no molecules in the training set have ANY of the unimportant 225 features. Definitely drop them!

## Template begins here:

This is the template you should use for submitting your practical assignments. A full-credit assignment will go through each of these points and provide a thoughtful clear answer. Note that the limit for the report is 4 pages, please prioritize quality over quantity in your submission.

## 1 Technical Approach

How did you tackle the problem? Credit will be given for:

- Diving deeply into a method (rather than just trying off-the-shelf tools with default settings). This can mean providing mathematical descriptions or pseudo-code.
- Making tuning and configuration decisions using thoughtful experimentation. This can mean carefully describing features added or hyperparameters tuned.
- Exploring several methods. This can mean contrasting two approaches or perhaps going beyond those we discussed in class.

Thoughtfully iterating on approaches is key. If you used existing packages or referred to papers or blogs for ideas, you should cite these in your report.

Mention Features	
Feature	Value Set
Mention Head	$\mathcal{V}$
Mention First Word	$\mathcal{V}$
Mention Last Word	$\mathcal{V}$
Word Preceding Mention	$\mathcal{V}$
Word Following Mention	$\mathcal{V}$
# Words in Mention	$\{1, 2, \dots\}$
Mention Type	$\mathcal{T}$

Table 1: Feature lists are a good way of illustrating problem specific tuning.

Model	Acc.
BASELINE 1	0.45
BASELINE 2	2.59
MODEL 1	10.59
MODEL 2	13.42
MODEL 3	7.49

Table 2: Result tables can compactly illustrate absolute performance, but a plot may be more effective at illustrating a trend.

## 2 Results

This section should report on the following questions:

- Did you create and submit a set of predictions?
- Did your methods give reasonable performance?

You must have *at least one plot or table* that details the performances of different methods tried. Credit will be given for quantitatively reporting (with clearly labeled and captioned figures and/or tables) on the performance of the methods you tried compared to your baselines.

## 3 Discussion

End your report by discussing the thought process behind your analysis. This section does not need to be as technical as the others but should summarize why you took the approach that you did. Credit will be given for:

- Explaining the your reasoning for why you sequentially chose to try the approaches you did (i.e. what was it about your initial approach that made you try the next change?).
- Explaining the results. Did the adaptations you tried improve the results? Why or why not? Did you do additional tests to determine if your reasoning was correct?