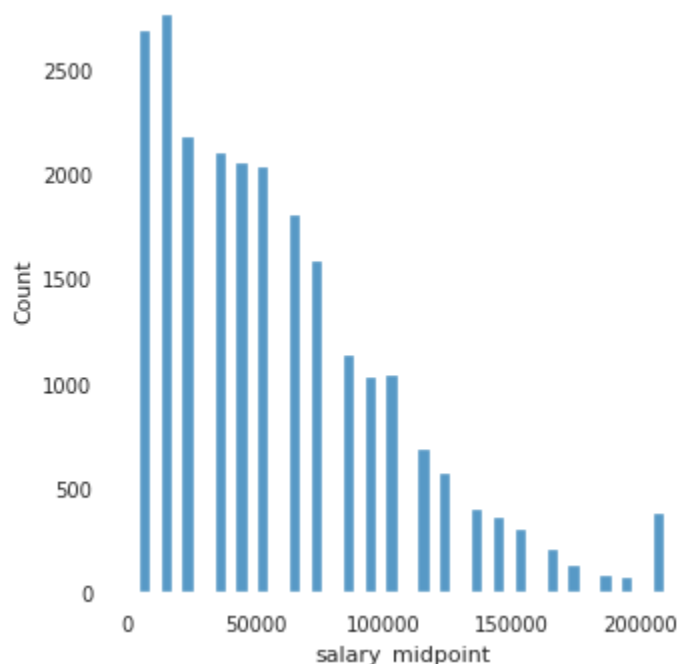


1/ Introduction and about the data

In this notebook, we will label a data scientist's salary based on information like location, gender, experience, etc. The original target value is the salary midpoint as shown in the below graph:



To simplify, we will narrow it down to groups: less than 20.000, from 20.000 to 85.000 and over 85.000 which are equivalent to label 1, 2, 3 respectively.

2/ Data Cleansing

2a/ We'll first remove absolutely irrelevant features like "aliens", "dogs_vs_cats".

2b/ Then convert the country names to its coordinates.

2c/ And convert the resting features like education, gender to binary values.

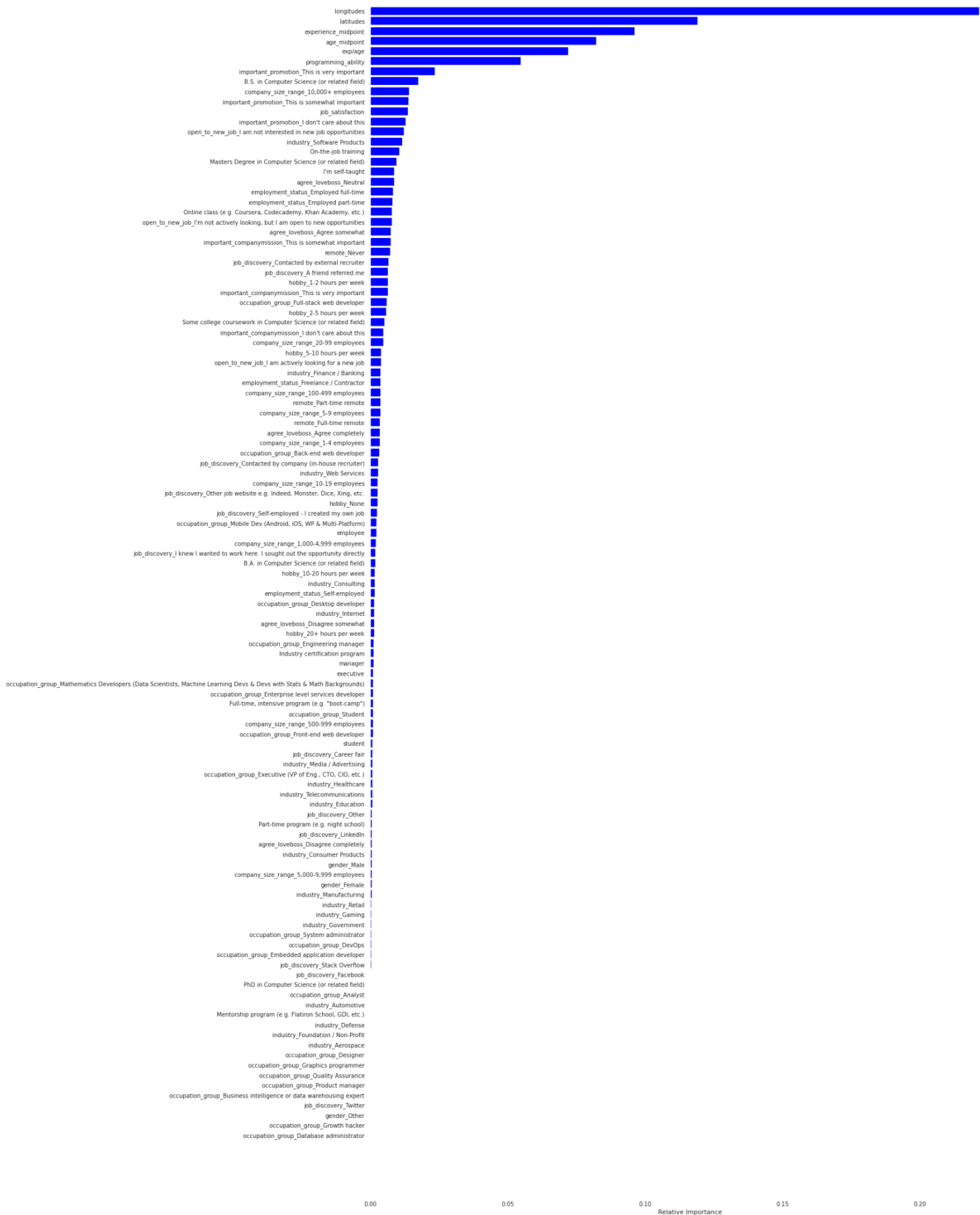
3/ Splitting and Sampling

After splitting the data, you will see that we don't have much data on the first and the third group. To address this imbalance, we'll use oversampling (or undersampling, if you try both, it will be pretty much the same).

4/ Randomly choose an algorithm and Feature Selection

So far we've had 113 features in total and we want to know how important each is. I am really sorry for the chart below, it is too small to read.

My idea is that we won't manually select features and have a check at each selection. Instead, we just need to drop the last features with random ratios. For example, in the notebook, I cut off 1/4 of all features



And here is the results comparison:

	set	label 1	label 2	label 3	f1_score
0	train_set	0.87	0.78	0.78	0.86
1	val_set	0.86	0.78	0.78	0.77
2	test_set	0.85	0.74	0.74	0.77

Using all features

	set	label 1	label 2	label 3	f1_score
0	train_set	0.8	0.78	0.79	0.8
1	val_set	0.8	0.77	0.78	0.77
2	test_set	0.78	0.74	0.76	0.78

Using 3/4 features (This is a bit more reasonable)

This result is acceptable in general. But as you can see, there is no difference between dropping or not dropping 1/3 of the features. Please let me know if you know why it is. I also tried dropping rates of 1/2 and 1/4 but the performance was worse.