

I had checked the consistency of the subject classification between the shrinkage model, $\alpha=0.3$ and the multi lm ols model.

In both models I focused on the outliers of the predictions in the normal weight class, so those subjects that were classified into overweight or obese, but are actually normal weight.

I first checked how many subject are persistently classified into overweight and obese, despite being normal weight in visit1 and visit2, in each model.

Then I checked how many subjects are consistently classified into overweight and obese, despite being normal weight in visit1 between the two models, same in visit2.

In the end I checked how many of the persistently misclassified as overweight or obese subjects are consistent between the two models.

All though multi lm ols model is giving us better classification scores and model criteria, the number of persistently misclassified as overweight subjects is smaller, while the number of misclassified as obese subjects is hard to compare as the number of classifications is very small.

As seen in the confusion matrices(I have added them again here), the distribution of prediction is quite different, with substantially more misclassified as overweight or obese subjects in multi lm ols model in visit 2.

But from the consistency checks, this does not mean there are more persistently misclassified as overweight or obese subjects in both visits.

Prediction tables for each visit, by taking regression coefficients from the independent dataset and constructing the diet score in that visit. **Green**=”right”, **Red**=”wrong” **Yellow** = “our interest”

Shrinkage model, alpha 0.3:

visit 1:

<div> <div></div> <div>true</div> </div> <div>predicted</div>	0	1	2
0	14501	7514	1910
1	2825	3656	774
2	2	0	1

visit 2:

<div> <div></div> <div>true</div> </div> <div>predicted</div>	0	1	2
0	8251	5477	2143
1	5002	7694	2501
2	37	44	34

OLS model, all significant variables:

visit 1:

<div> <div></div> <div>true</div> </div> <div>predicted</div>	0	1	2
0	14545	7560	1819
1	2761	3585	849
2	22	25	17

visit 2:

<div> <div></div> <div>true</div> </div> <div>predicted</div>	0	1	2
0	6347	3984	1406
1	6752	9005	3114
2	191	226	158

Elastic net shrinkage model, alpha=0.3, persistently misclassified as overweight in both visits:

1827 which is 0.365% of visit2(5002) and 0.646% of visit1(2825)

Elastic net shrinkage model, alpha=0.3, persistently misclassified as obese in both visits:

0 which is 0% of visit2(37) and 0% of visit1(2)*

*cant really work with these little numbers

Multi lm ols model, persistently misclassified as overweight in both visits:

1718 which is 0.254% of visit2(6752) and 0.622% of visit1(2761)

Multi lm ols model, persistently misclassified as obese in both visits:

2 which is 0.010% of visit2(191) and 0.090% of visit1(22)*

*this 2 is very little compare to the number of misclassifications

Consistency between the models:

visit1, consistently misclassified as overweight:

2285 which is 0.808% of visit1 shrinkage(2825) and 0.827% of visit1 multi lm ols(2761)

visit2, consistently misclassified as overweight:

4321 which is 0.863 % of visit1 shrinkage(5002) and 0.639% of visit1 multi lm ols(6752)

visit1, consistently misclassified as obese:

2 which is 100% of visit1 shrinkage(2) and 0.090% of visit1 multi lm ols(22)

visit2, consistently misclassified as obese:

36 which is 0.972 % of visit1 shrinkage(37) and 0.188% of visit1 multi lm ols(191)

Consistently and persistently misclassified as overweight:

1454 which is 0.795% of shrinkage persistently misclassified (1827) and 0.846% of multi ols persistently misclassified(1718)

Considering our aim of identifying the persistently lean, might be good to reconsider the modeling used, it will be more clear when we decide how exactly are we going to identify the persistently lean,

is it going to be just by applying a cutoff point of the total environment score and selecting those normal weight subjects that are above that cutoff point,

or ,

work with the predictions and take the persistently misclassified in to overweight or obese.

In the latter we are more blind in what environment were the misclassified really exposed to, but that could be said in any case in the first option as well and with the latter option, it might be better to do a more suitable model, like rf with oversampling of minority classes like obese class.

Maybe we can try both...