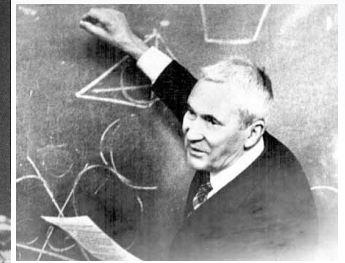
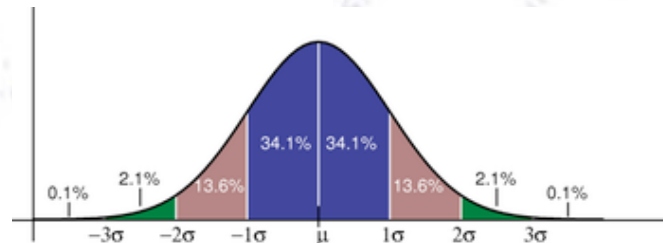


Applied ML

Results and Scores of Initial Project



Thomas Spieksma & Troels C. Petersen (NBI)



"Statistics is merely a quantisation of common sense - Machine Learning is a sharpening of it!"

Overall comments

Opening comment from 2021:

The name “Small Project” is misleading, and should have been “Initial project”, because it is **by no means small**.

Overall comments

Opening comment from 2021:

The name “Small Project” is misleading, and should have been “Initial project”, because it is **by no means small**.

You did very well, and so let me start by gently stating, that you have little/nothing to fear - in fact, you did really great!

Grading it was perhaps comparable to the project itself, but we have done our best to be as open as possible about the scoring. And to give you a maximum of feedback, we have produced a report for each of you.

The motivation

We wanted you to try the very **real challenge** of optimising models, without knowing their performance on the data it is applied to.

We also wanted you to **individually** run ML algorithms, so that you have the machinery in place after the course.

We insisted that you tried **both tree- and NN-based algorithms**, to get a feel for their differences and similarities.

We also wanted you to feel the “insecurity” about not knowing if you had gotten everything out of the data.

The description file was meant to trigger you to **think about your models**, and what you tried. Also, considerations of size and performance are in place.

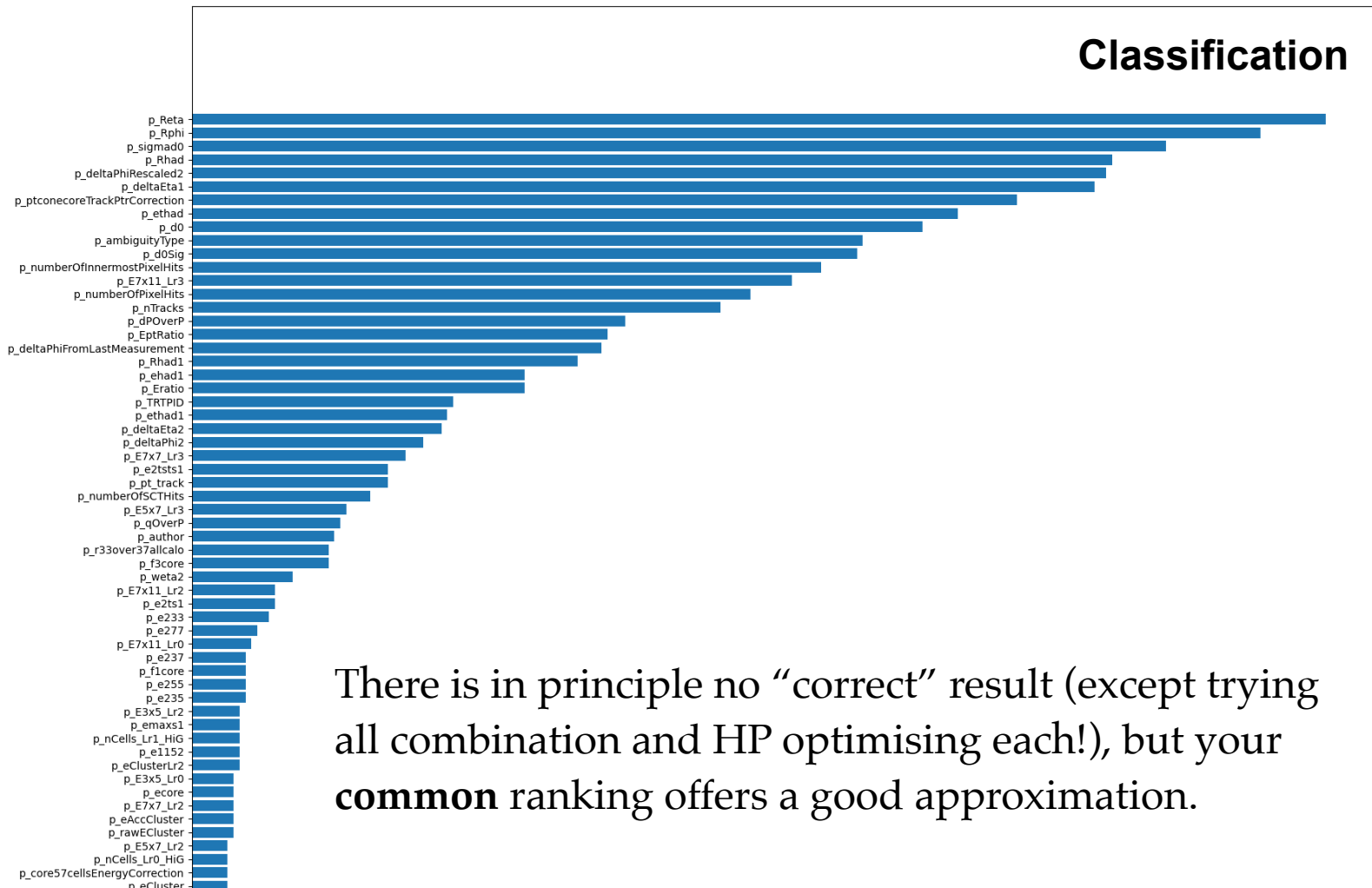
Finally, we wanted to **ensure** that you yourself tried all the work and things to consider, to put together ML models and apply them.



Classification Results

Classification variable usage

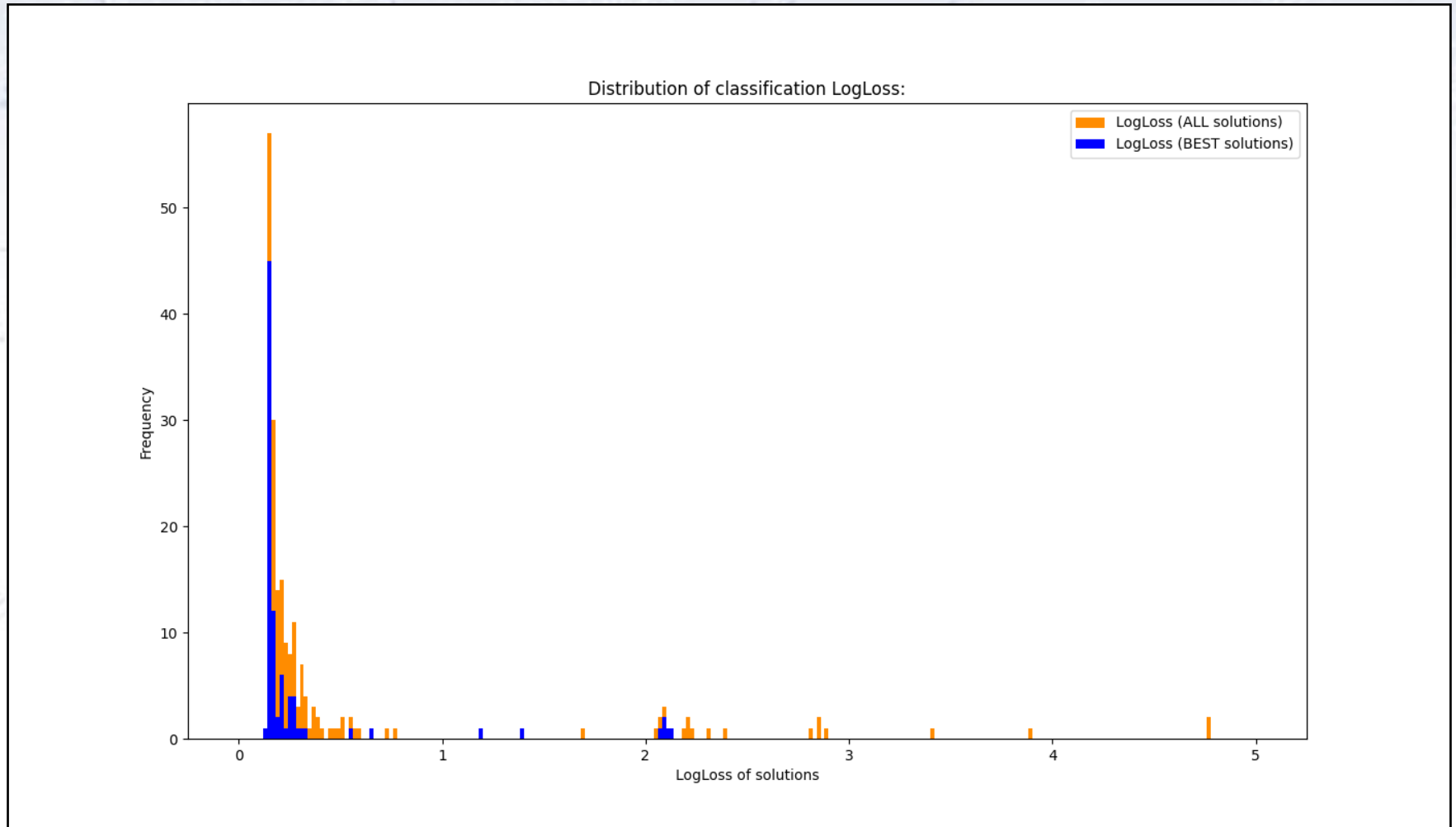
Many (most?) of you have made a good variable ranking. Below you find a variable usage frequency plot, showing how often a variable was used.



There is in principle no “correct” result (except trying all combination and HP optimising each!), but your **common** ranking offers a good approximation.

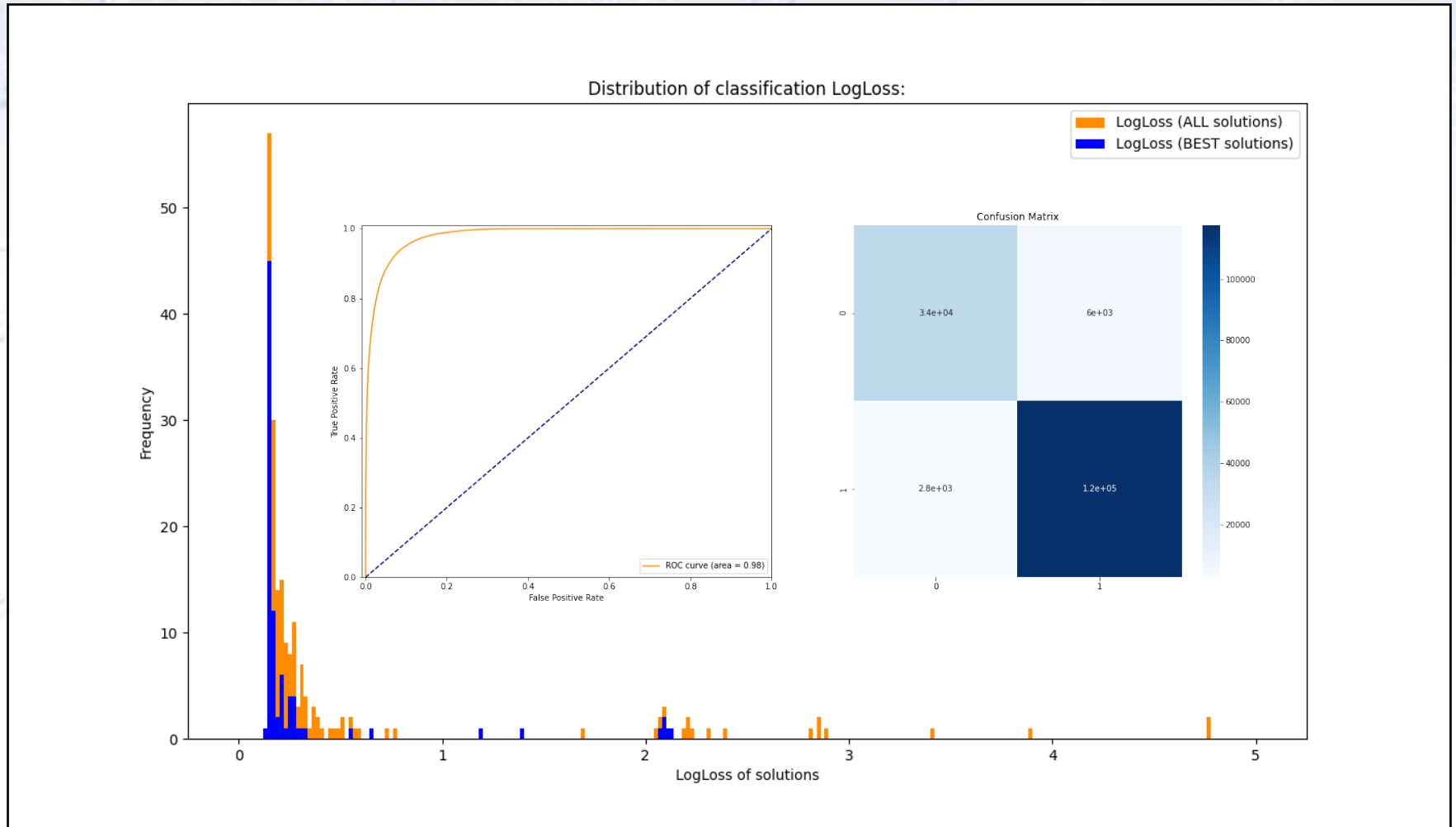
Classification score distribution

The distribution of the (Cross-Entropy) LogLoss values obtained was:



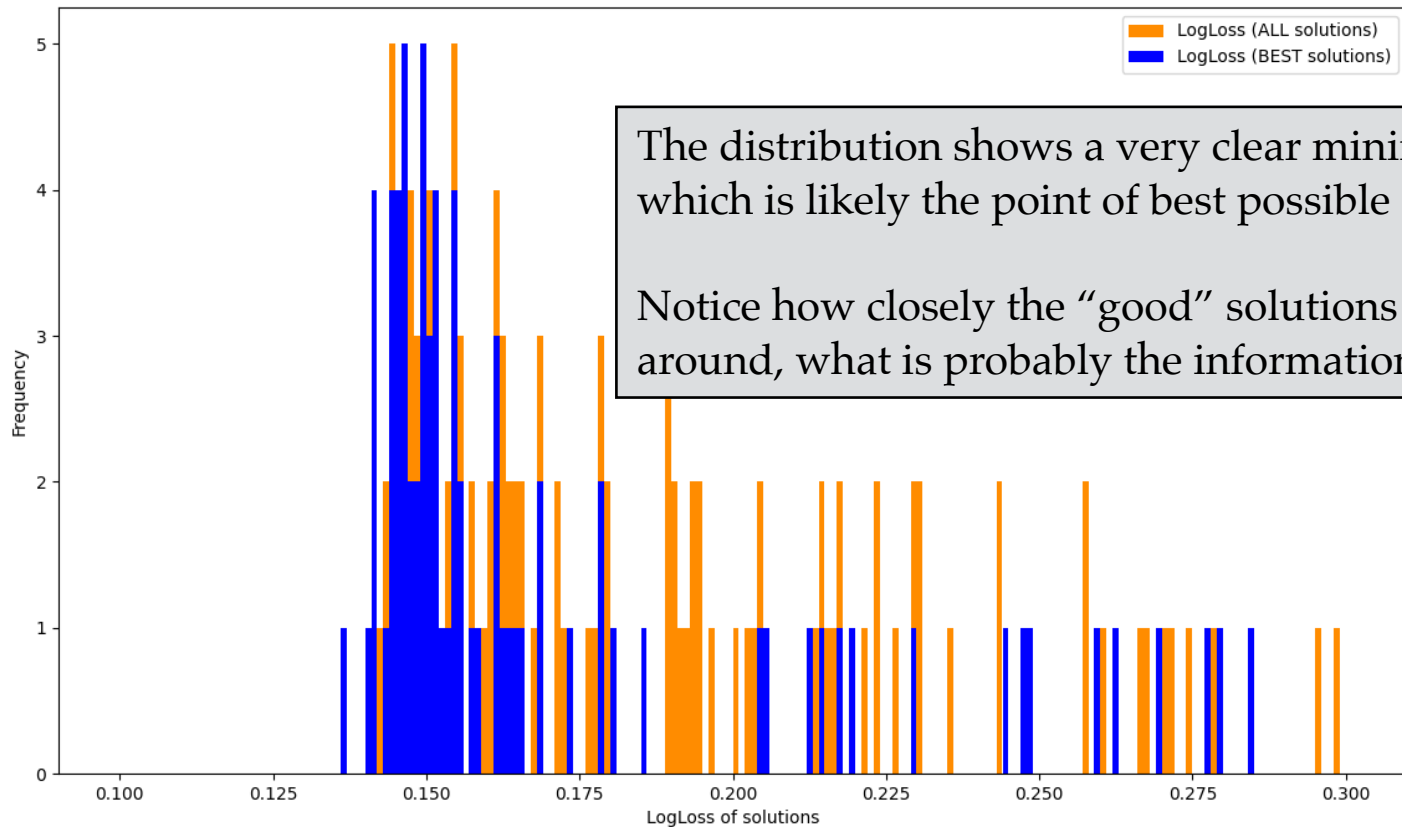
Classification score distribution

The distribution of the (Cross-Entropy) LogLoss values obtained was:



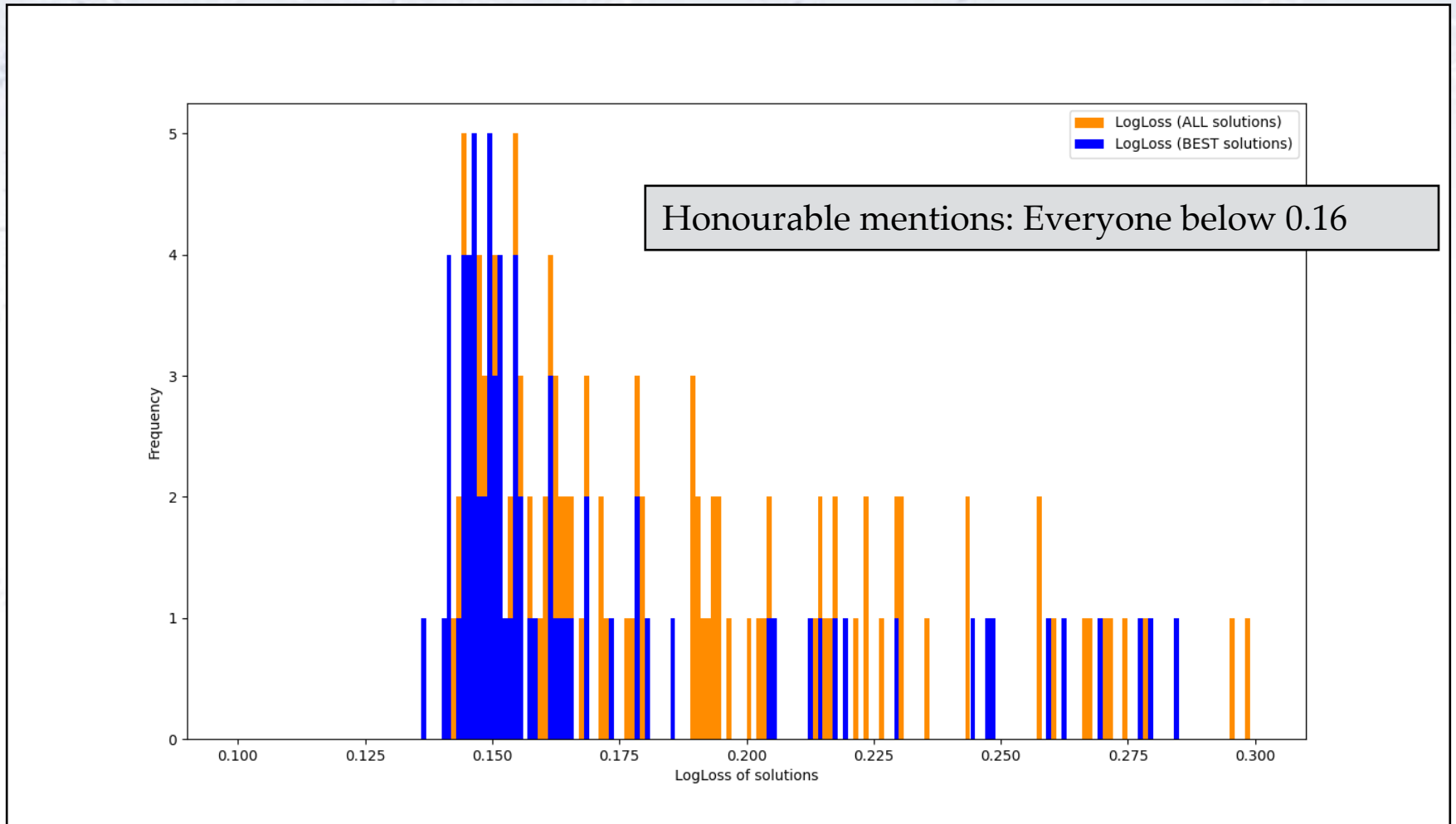
Classification score distribution

The distribution of the (Cross-Entropy) LogLoss values obtained was:



Classification score distribution

The distribution of the (Cross-Entropy) LogLoss values obtained was:

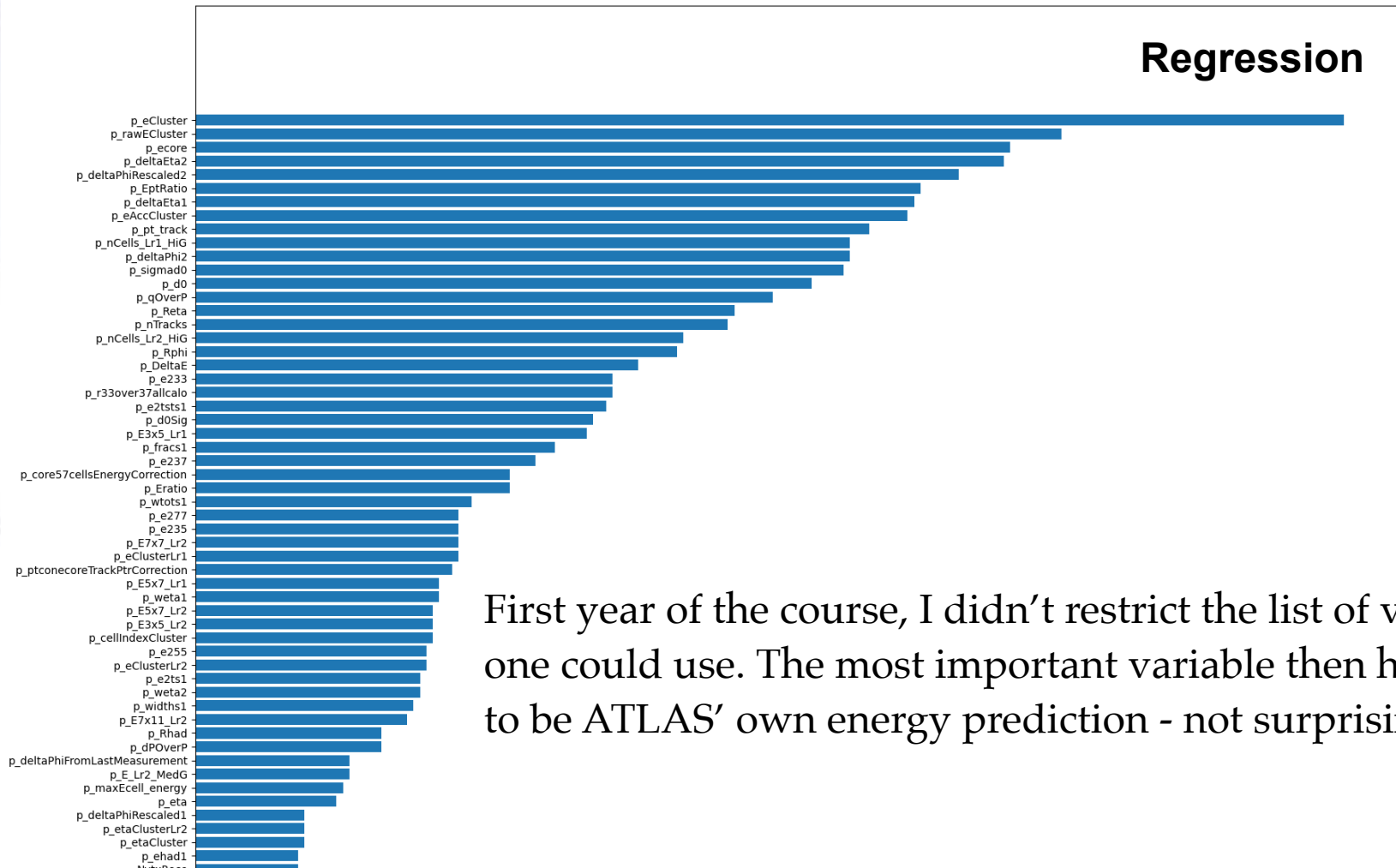




Regression Results

Regression variable usage

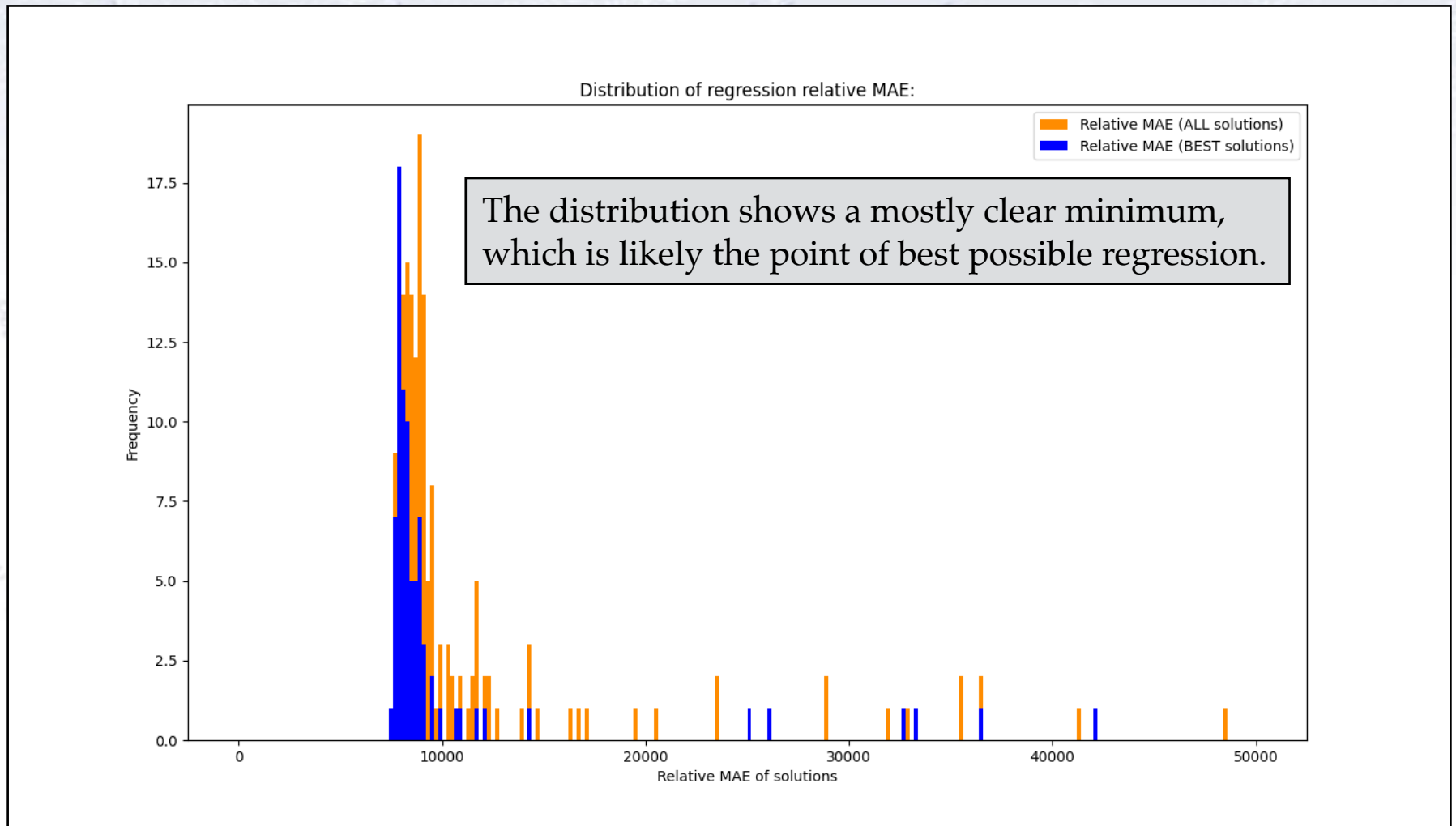
The variables have changed drastically from the classification case. There is NO overlap at all for the top 10-15 variables! Classification and Regression are in this case two very different tasks.



First year of the course, I didn't restrict the list of variables one could use. The most important variable then happens to be ATLAS' own energy prediction - not surprisingly!

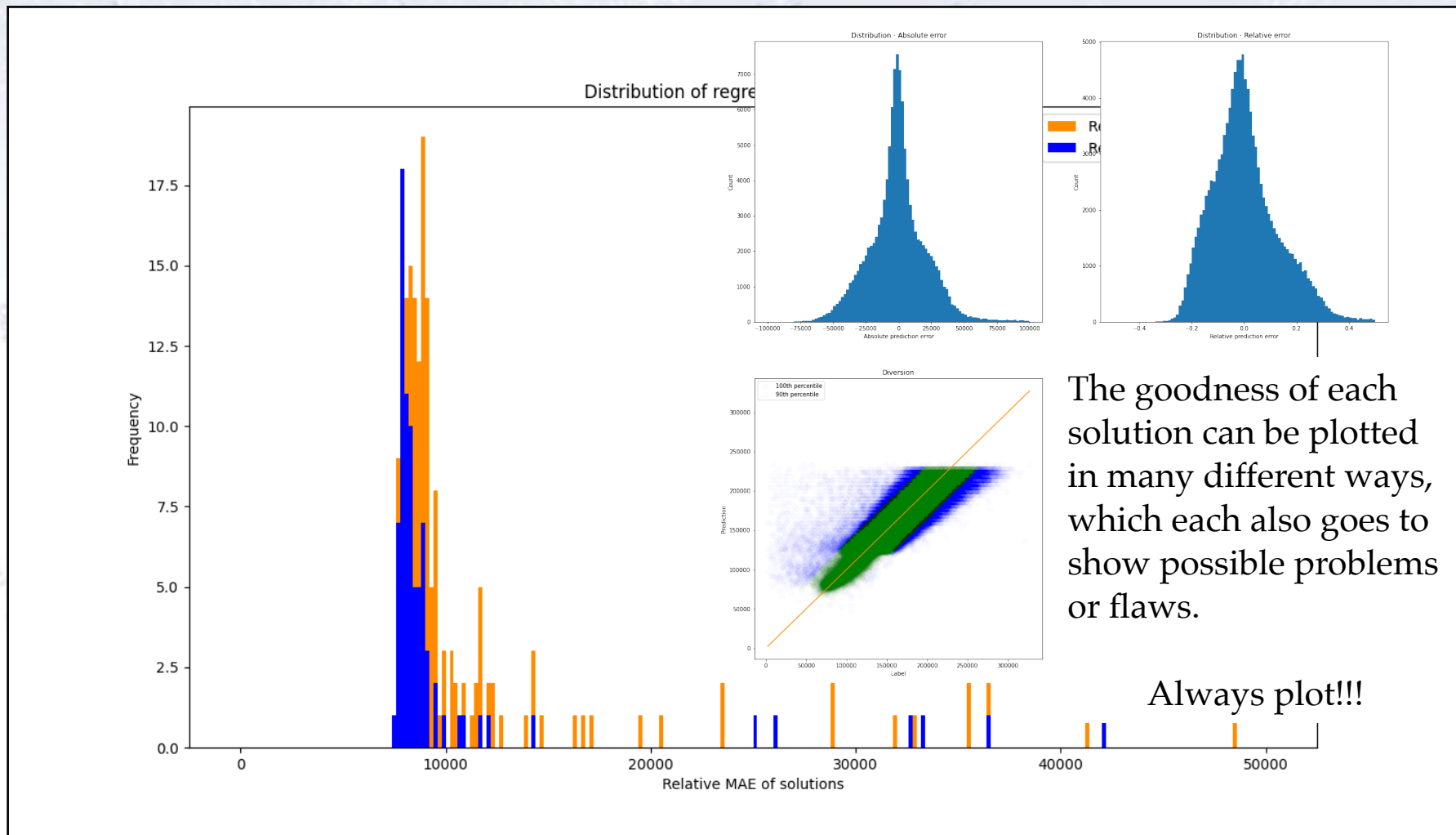
Regression score distribution

The distribution of the relative MAE (i.e. $MAE((E-T)/T)$) values obtained was:



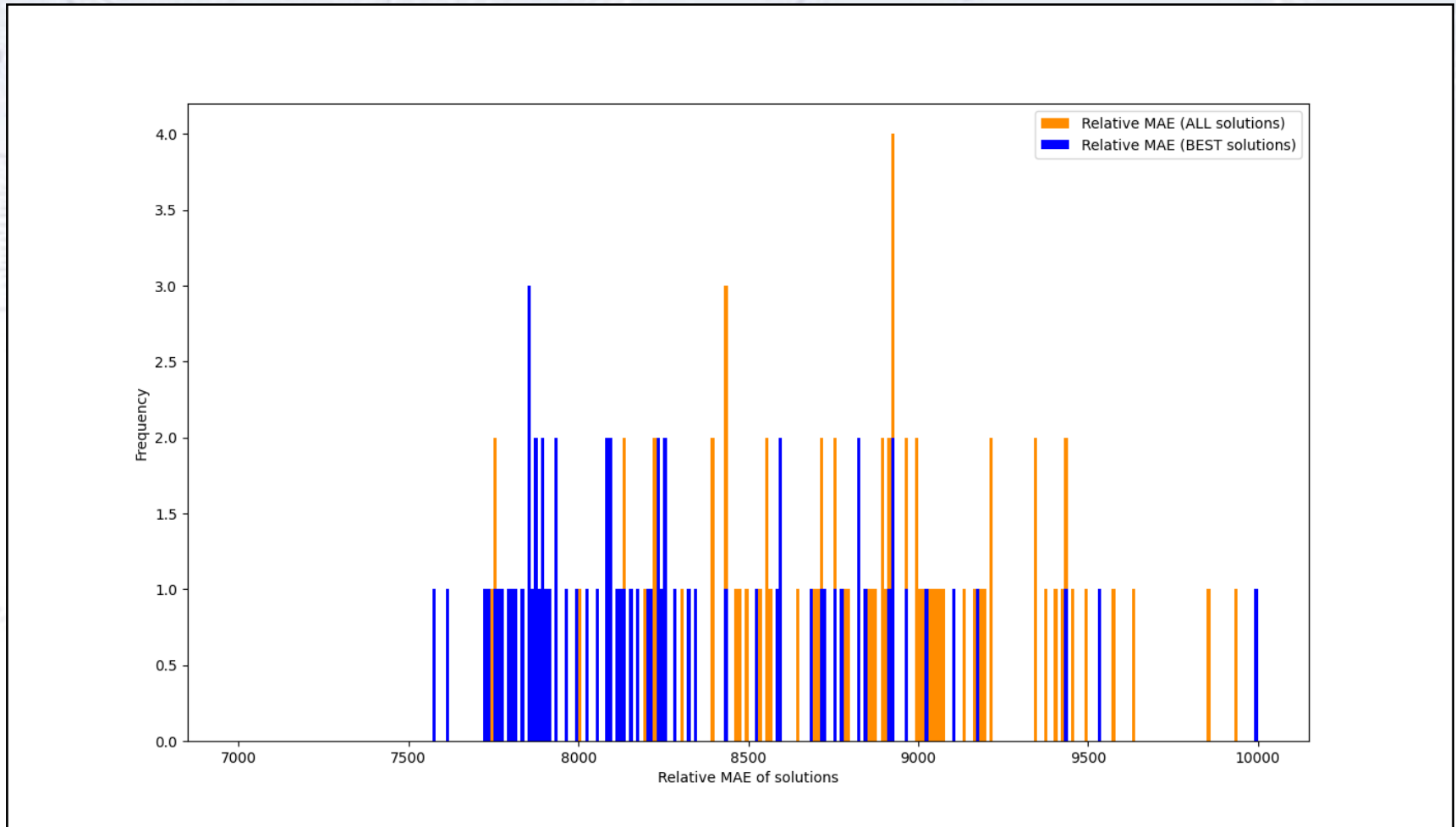
Regression score distribution

The distribution of the relative MAE (i.e. $MAE((E-T)/T)$) values obtained was:



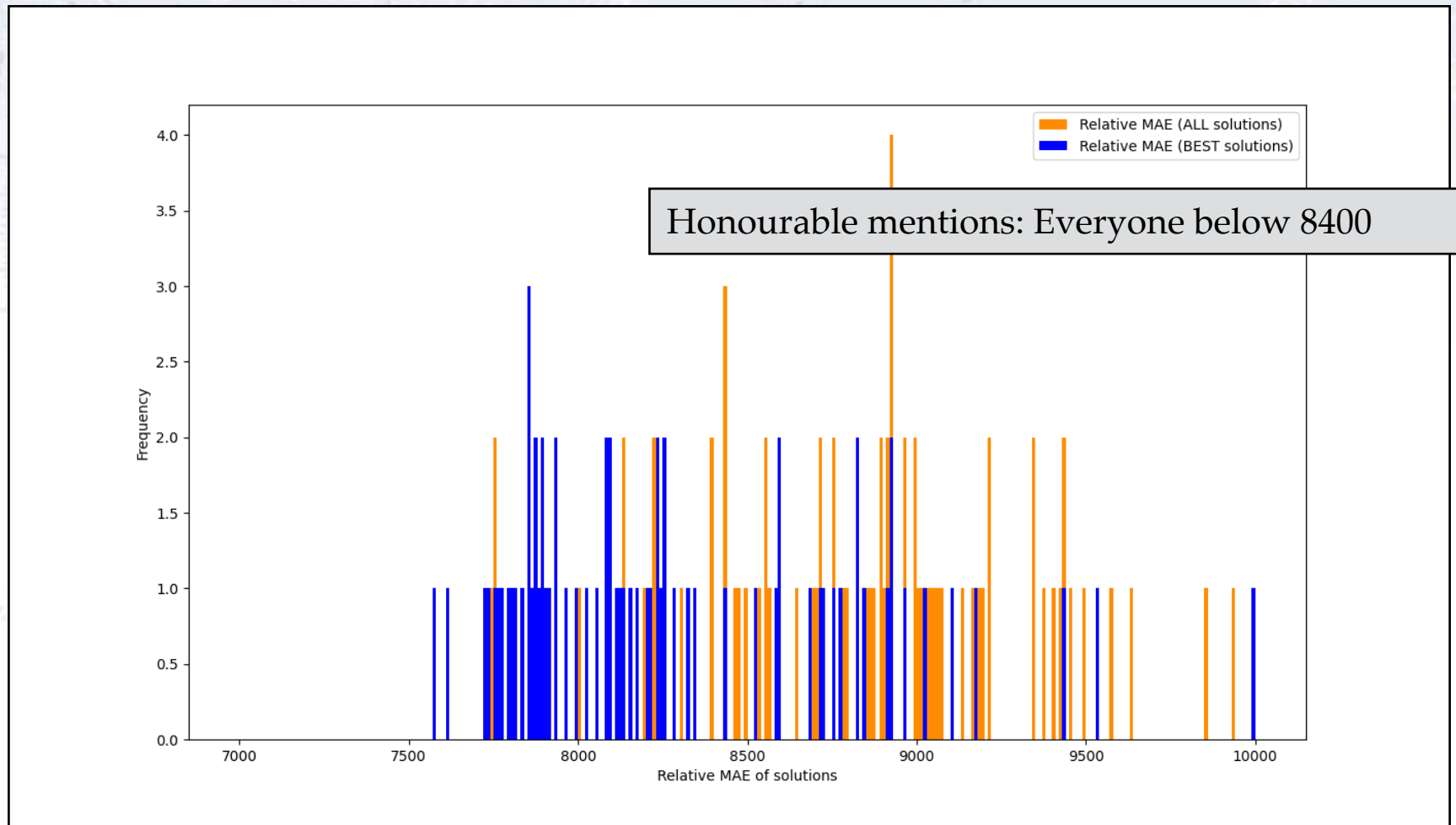
Regression score distribution

The distribution of the relative MAE (i.e. $\text{MAE}((E-T)/T)$) values obtained was:



Regression score distribution

The distribution of the relative MAE (i.e. $\text{MAE}((E-T)/T)$) values obtained was:

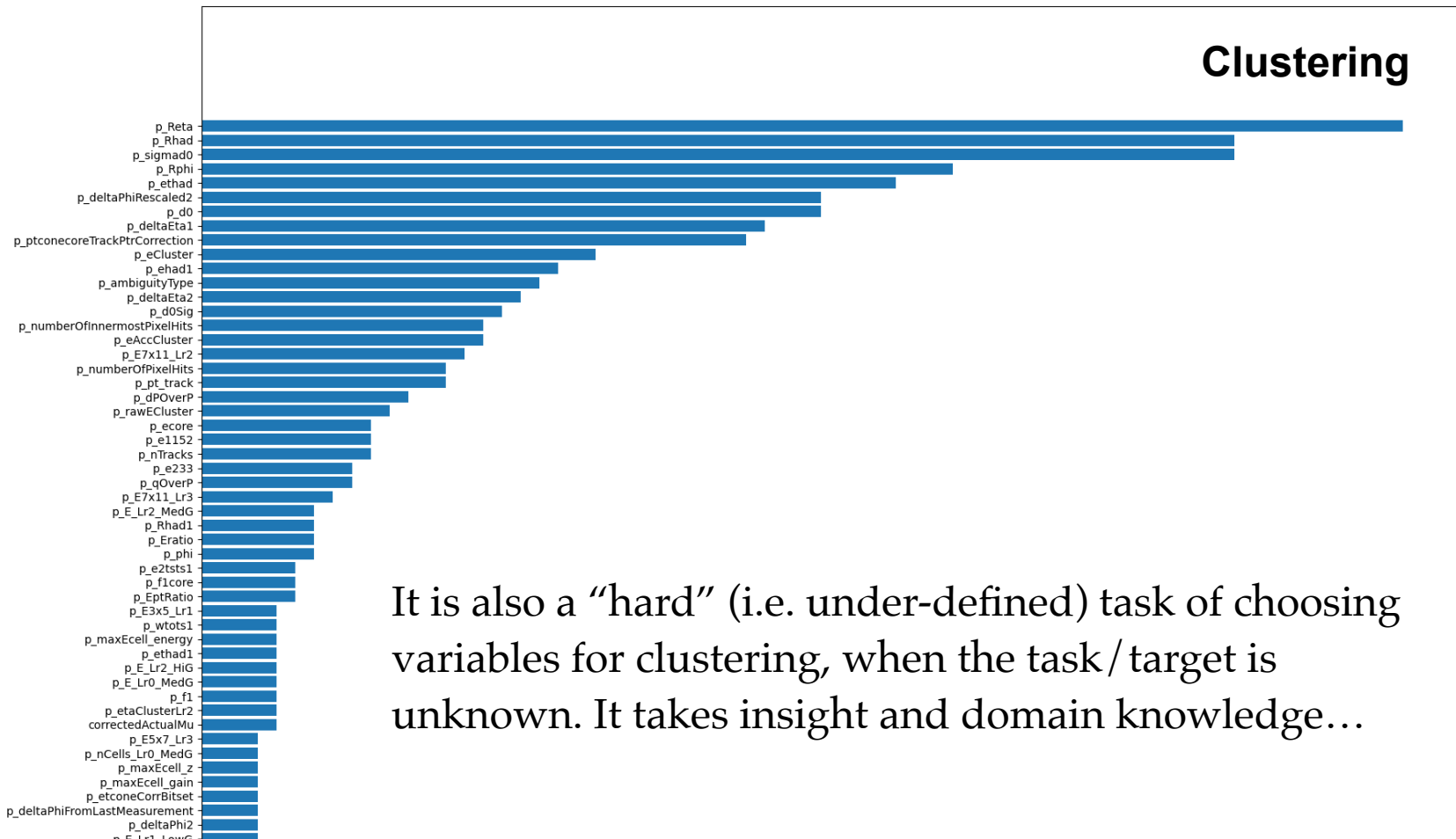




Clustering Results

Clustering variable usage

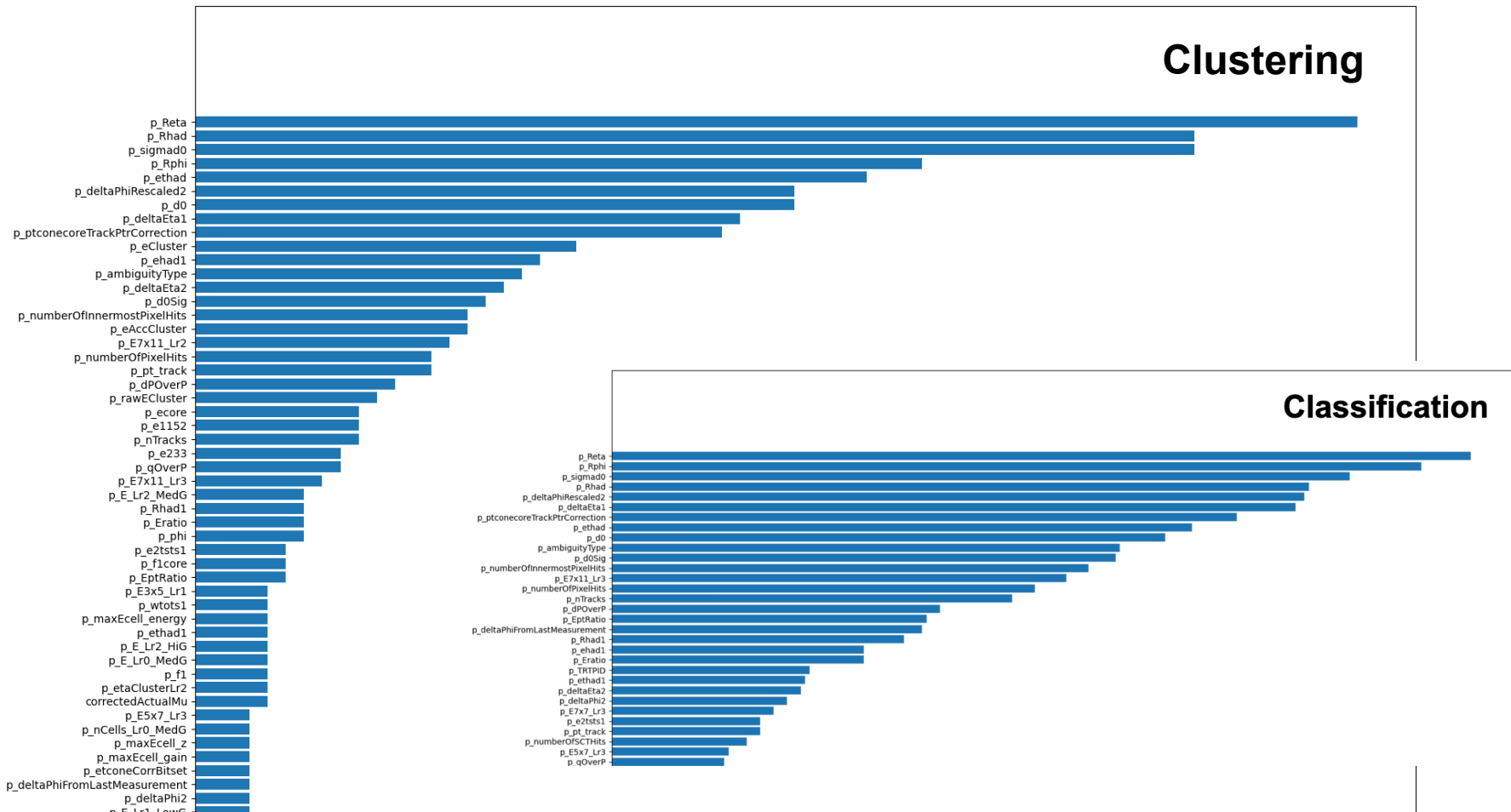
I would have thought, that the clustering variable usage would be near-identical to that of the (supervised) classification task. However, it is not entirely...



It is also a “hard” (i.e. under-defined) task of choosing variables for clustering, when the task / target is unknown. It takes insight and domain knowledge...

Clustering variable usage

I would have thought, that the clustering variable usage would be near-identical to that of the (supervised) classification task. However, it is not entirely...

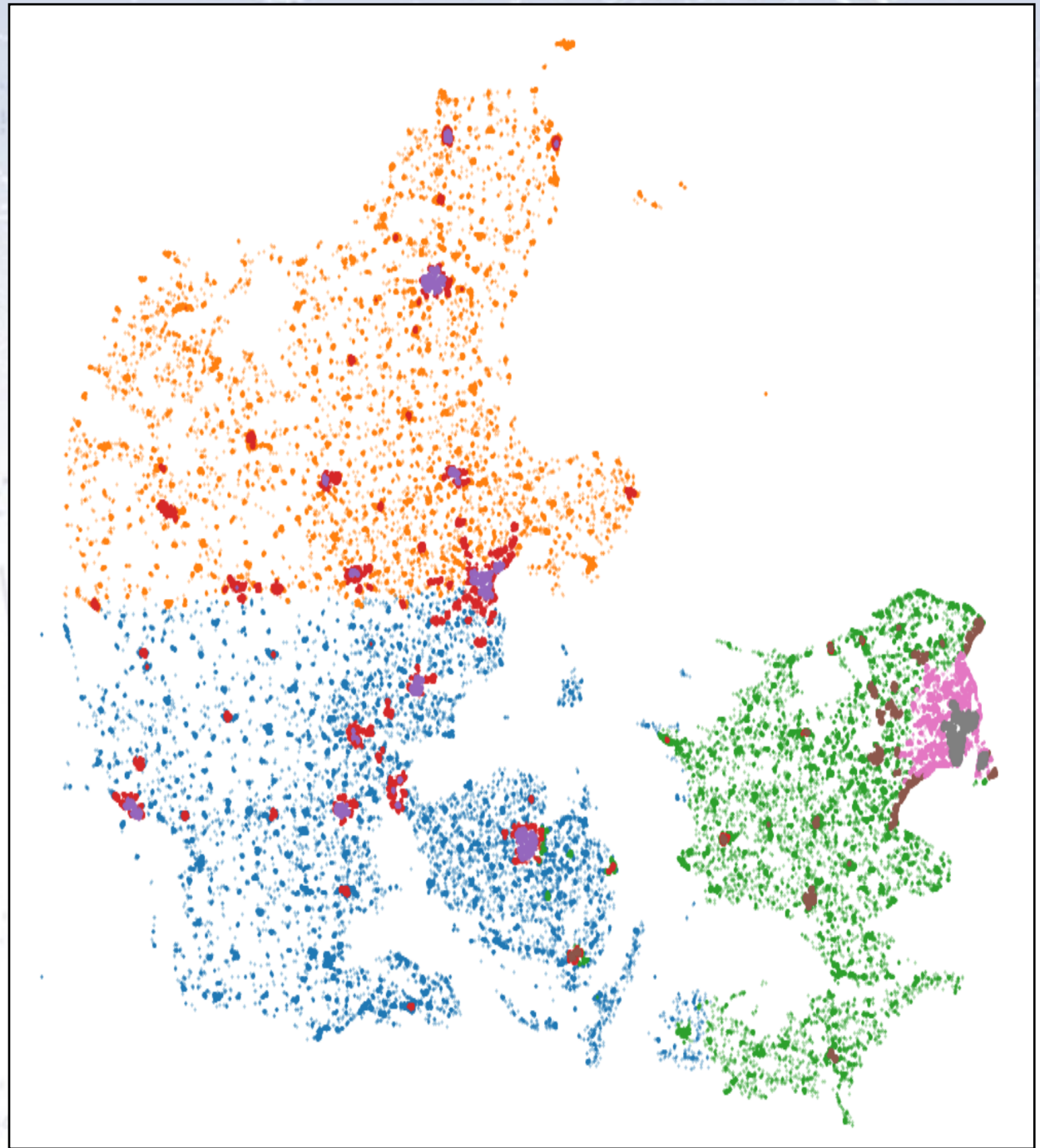


Clustering housing

While postal codes are good, they are not very useful in clustering Denmark.

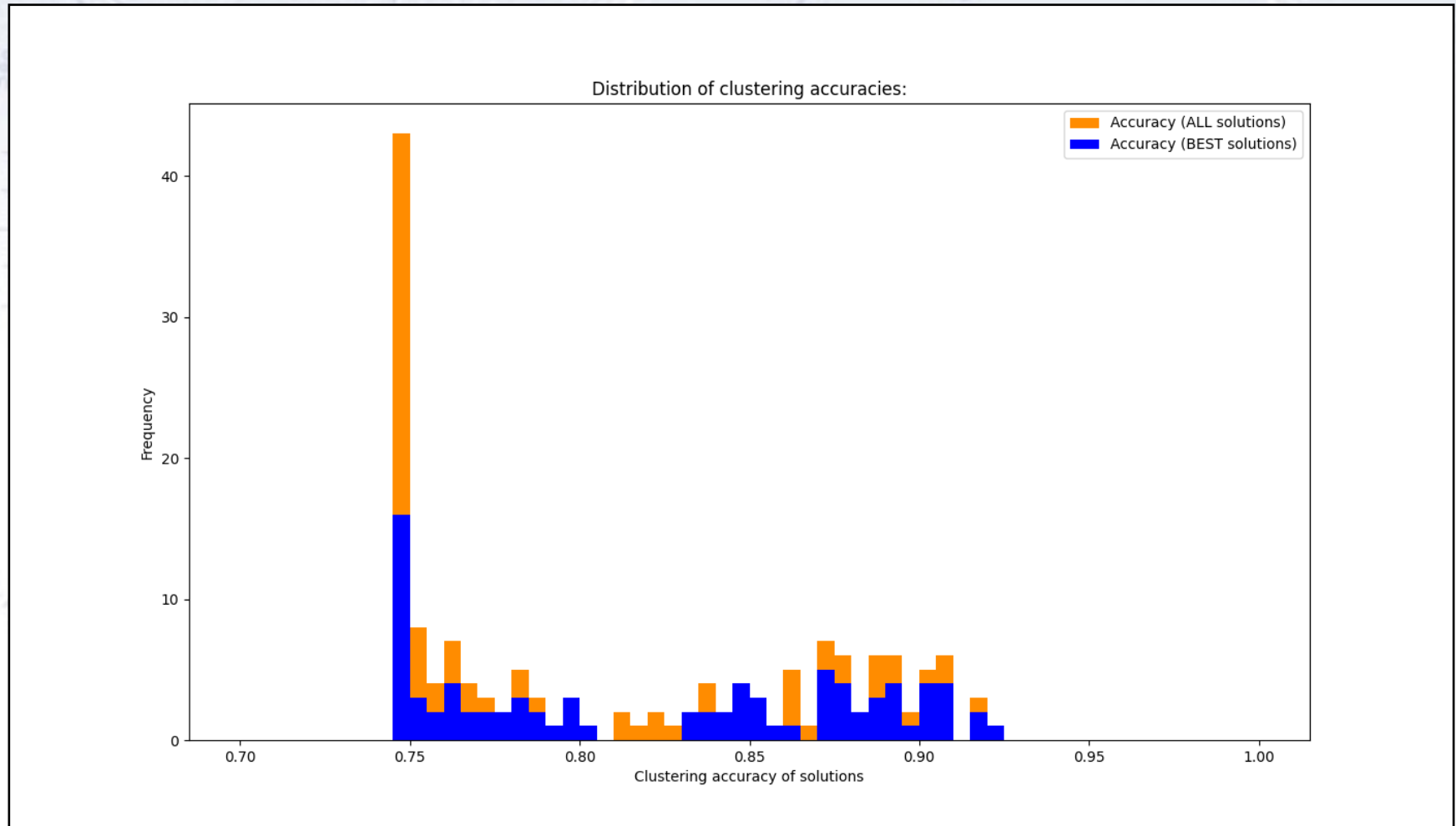
However, using just a few variables (x , y , density, price/m²), one can cluster villas in Denmark very efficiently.

In this way, one can follow trends for a type of house much better.



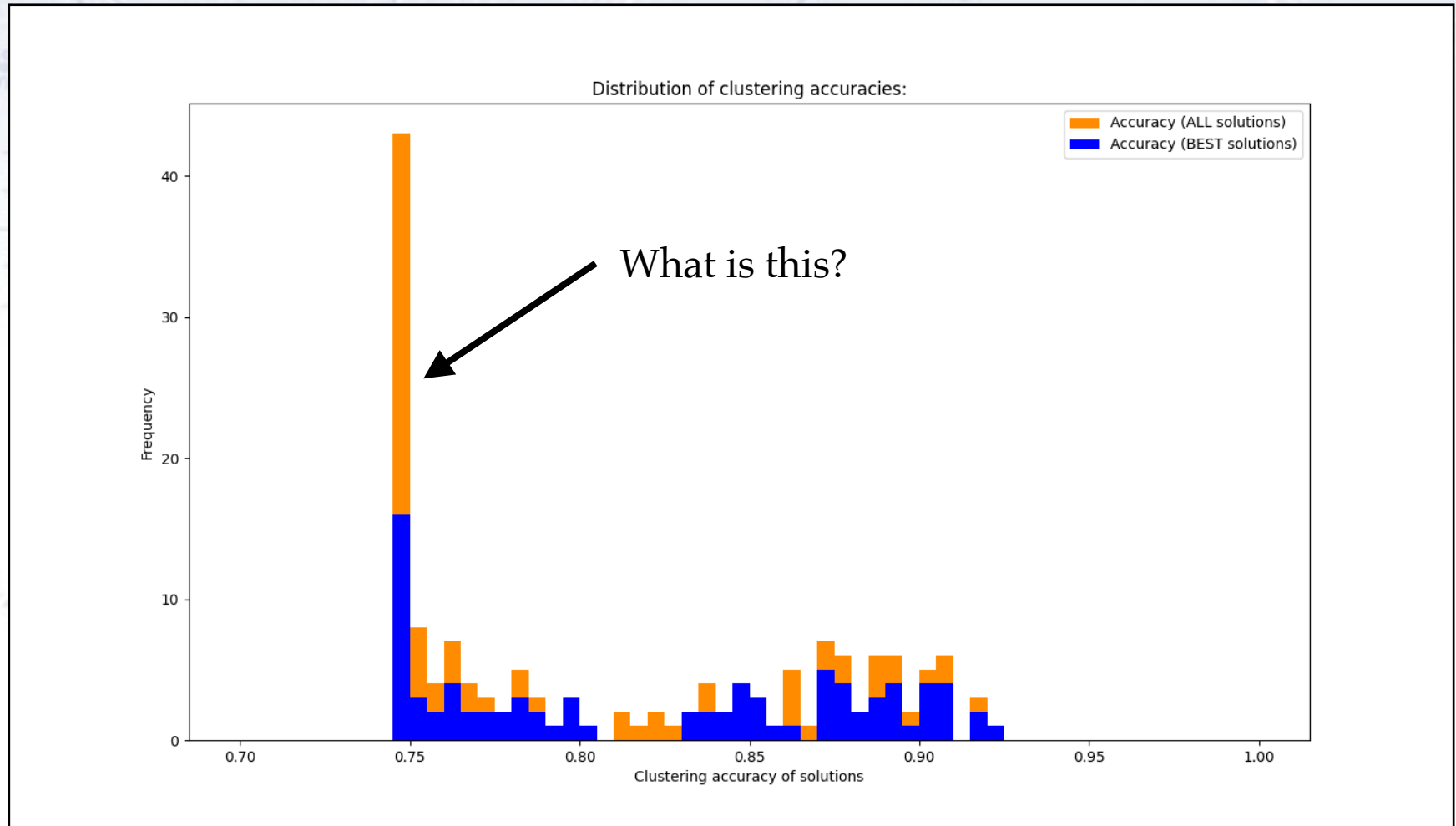
Clustering accuracy distribution

The accuracy of the clustering (when assigned either electron or not) was:



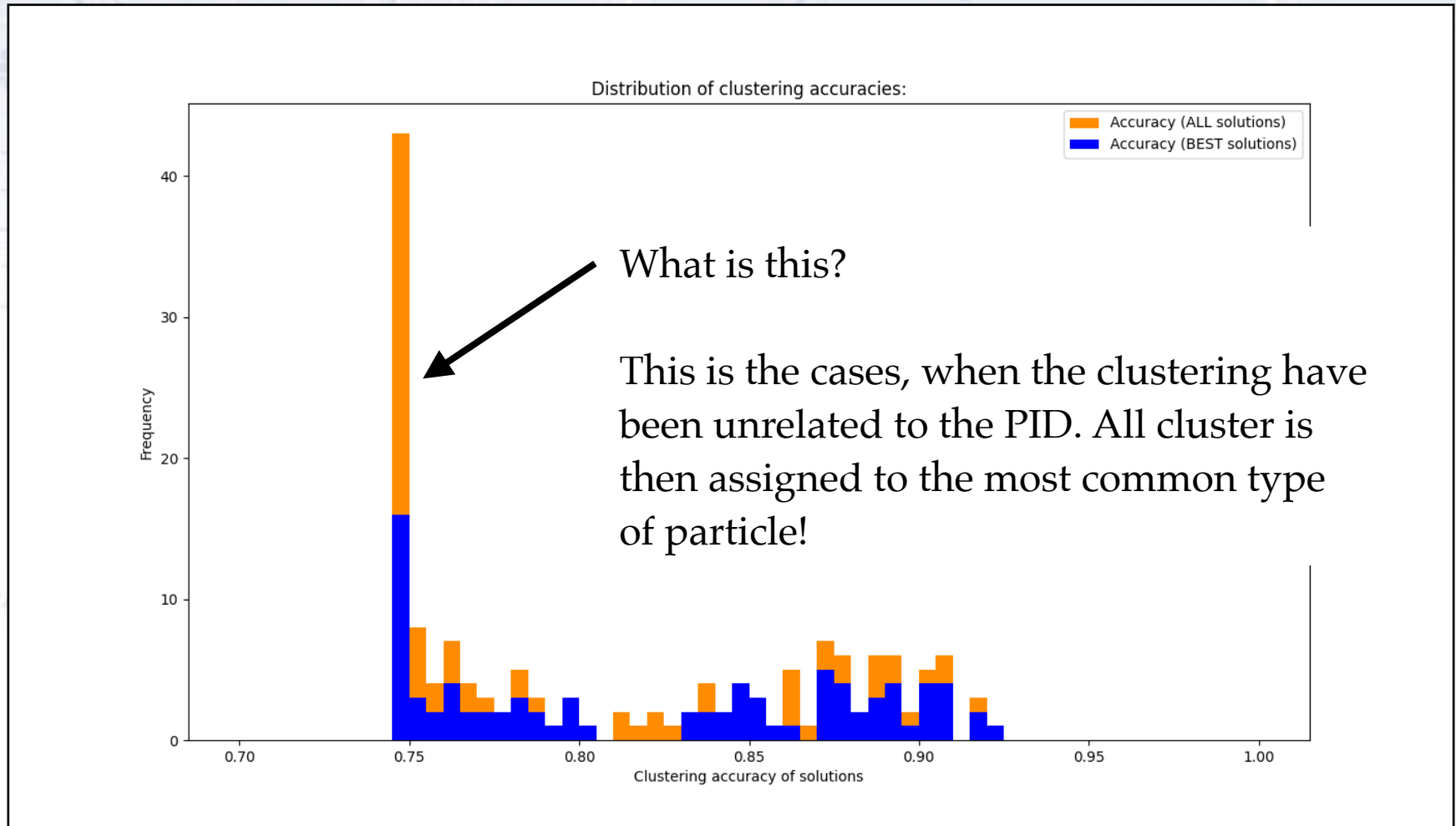
Clustering accuracy distribution

The accuracy of the clustering (when assigned either electron or not) was:



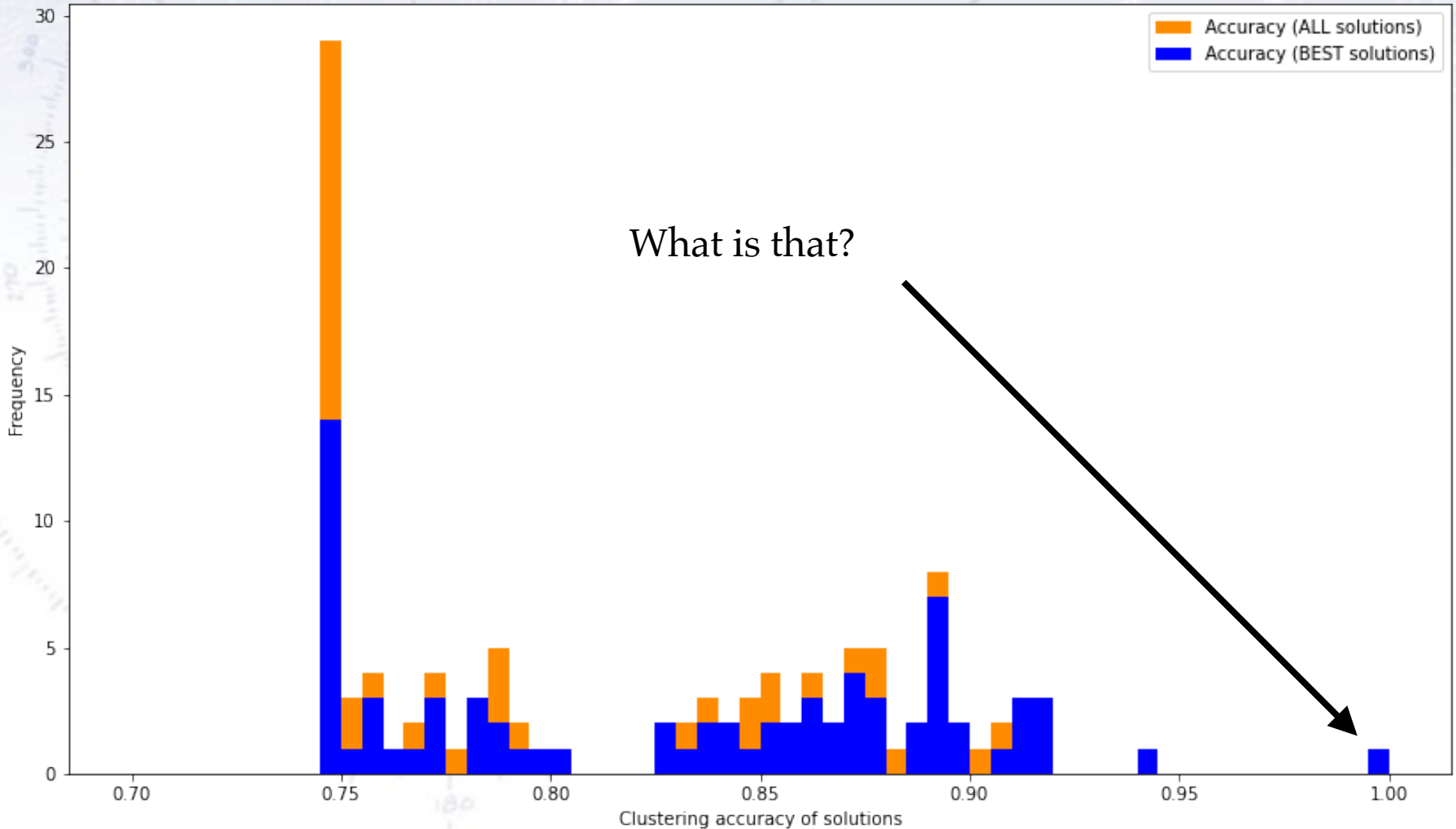
Clustering accuracy distribution

The accuracy of the clustering (when assigned either electron or not) was:



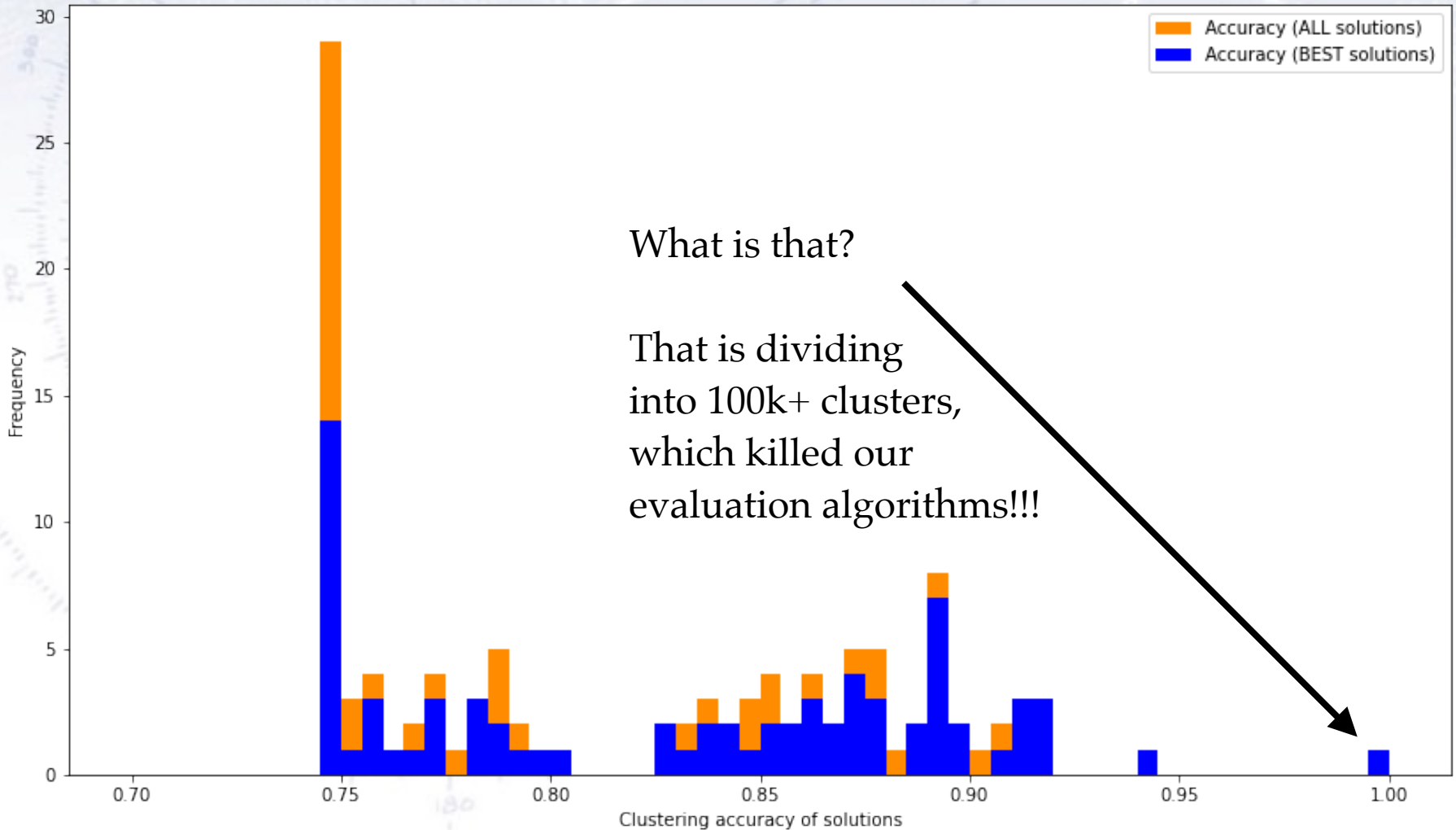
...last year's solutions!

Distribution of clustering accuracies:



...last year's solutions!

Distribution of clustering accuracies:

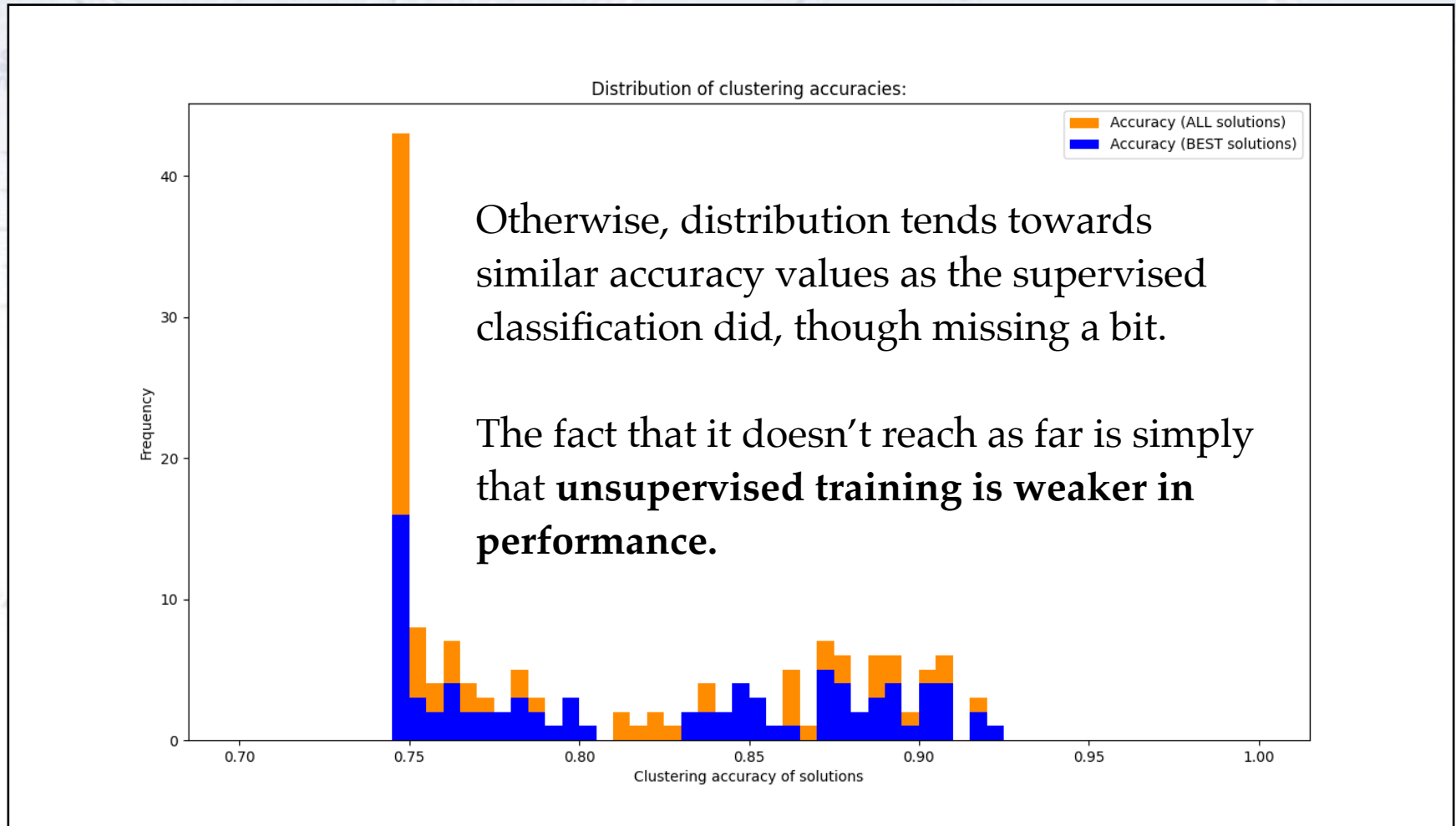


What is that?

That is dividing
into 100k+ clusters,
which killed our
evaluation algorithms!!!

Clustering accuracy distribution

The accuracy of the clustering (when assigned either electron or not) was:





Scoring your solutions

How do we grade your projects?

Final Score:

You submitted a full solution, from which you get: **67 points**

Your choice of methods based on your description was scored as follows [0, 6]:

Your solution entailed N different algorithms, which gives you a score of [0, 6]:

Your best performance for classification gave: $\max(0, (-\log(\text{CrossEntropy} - 0.12)) \times 1.4)$:

Your variable choice for classification was scored $4 \times (\text{VarFreq}(\text{you}) / \text{VarFreq}(\text{top}))$:

Your classification had 0 penalties, totalling to:

Your best performance for regression gave: $\max(0, -\log(\text{MAD}((E-T)/T)/7500-1) \times 1.8)$:

Your variable choice for regression was scored $5 \times (\text{VarFreq}(\text{you}) / \text{VarFreq}(\text{top}))$:

Your regression had 0 penalties, totalling to:

Your best performance for clustering gave: $\max(0, (\text{Accuracy} - 0.75) \times 20)$:

Your variable choice for clustering was scored $(\text{VarFreq}(\text{you}) / \text{VarFreq}(\text{top}))$:

Your clustering had 0 penalties, totalling to:

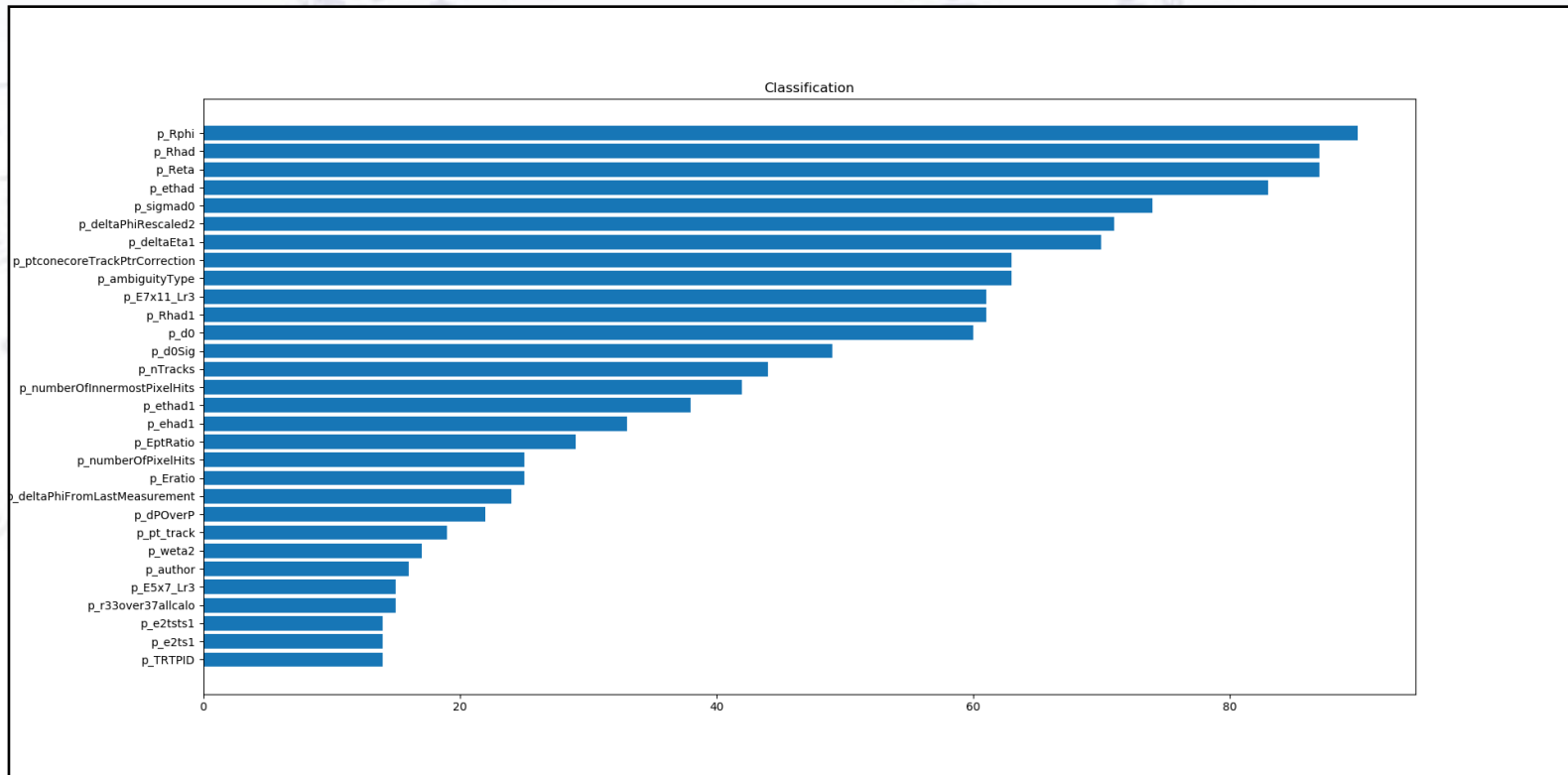
Thus your total number of points was:

Your variable choice

Assuming, that the variable frequency reflected the actual ranking very well, your variable choice was scored as follows (factors were 4, 5, and 1):

$$8 \times \left(\sum Freq(\text{Your variables}) / \sum Freq(\text{Top variables}) \right)$$

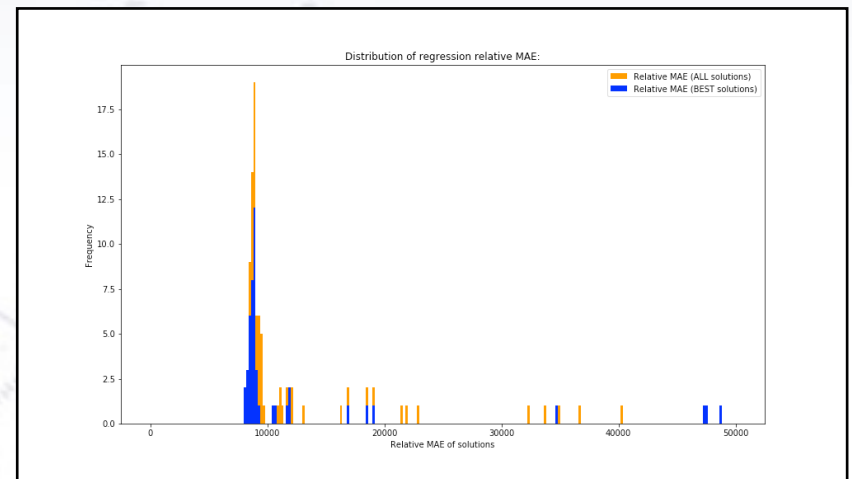
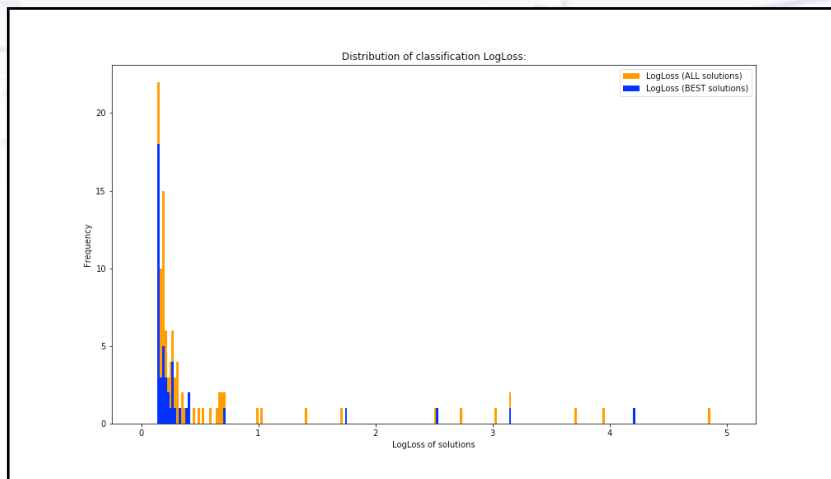
...so if you picked the top variables, you would get full points.



Performance scoring

As mentioned, performance isn't everything, and we certainly didn't want it to be for the small project. Getting close to the information limit is just great.

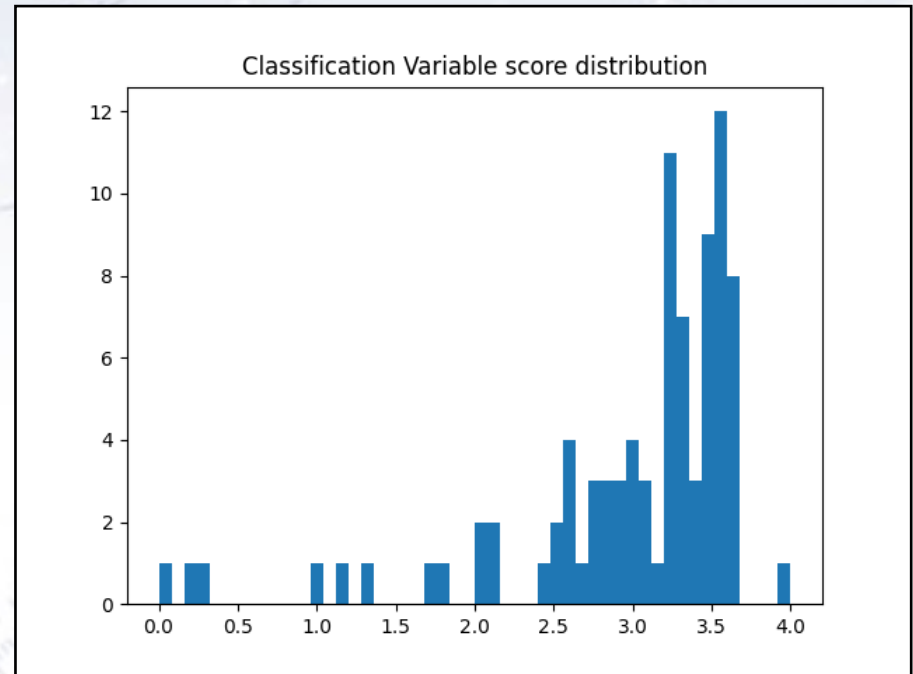
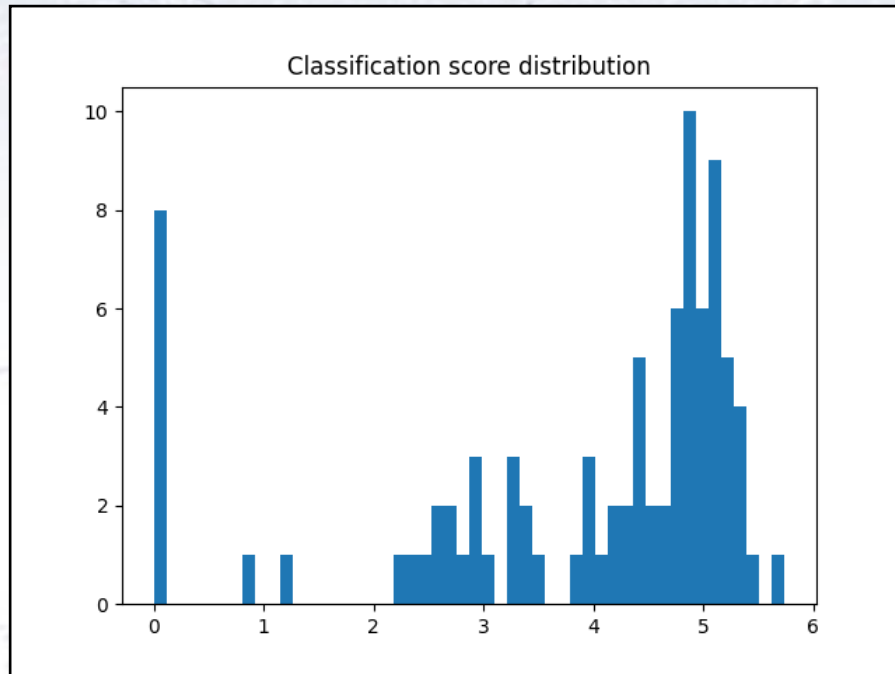
This was reflected by using a logarithmic scoring, which turned your best key performance parameter into a score in the (open) range $[0,5+]$:



In all of this, you could of course not get negative points for an accepted solution!

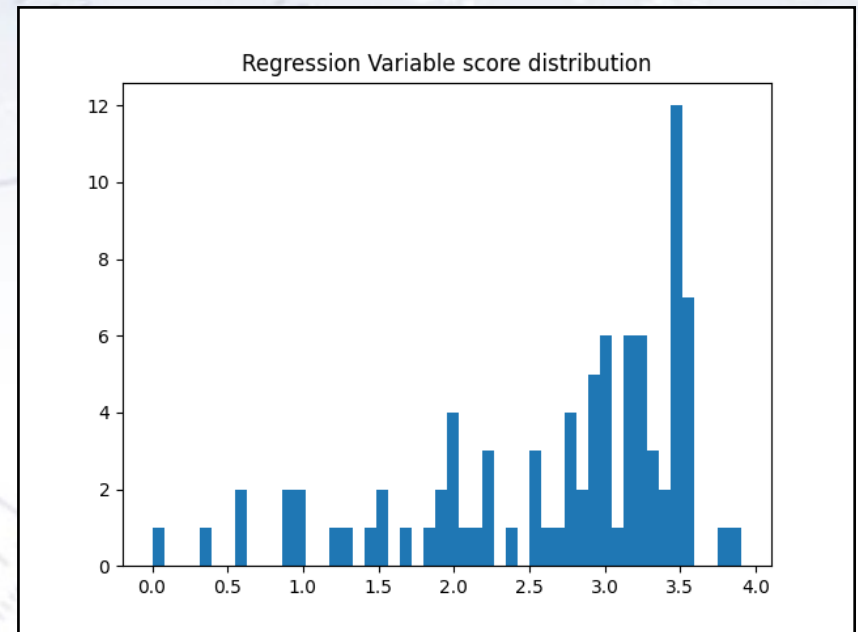
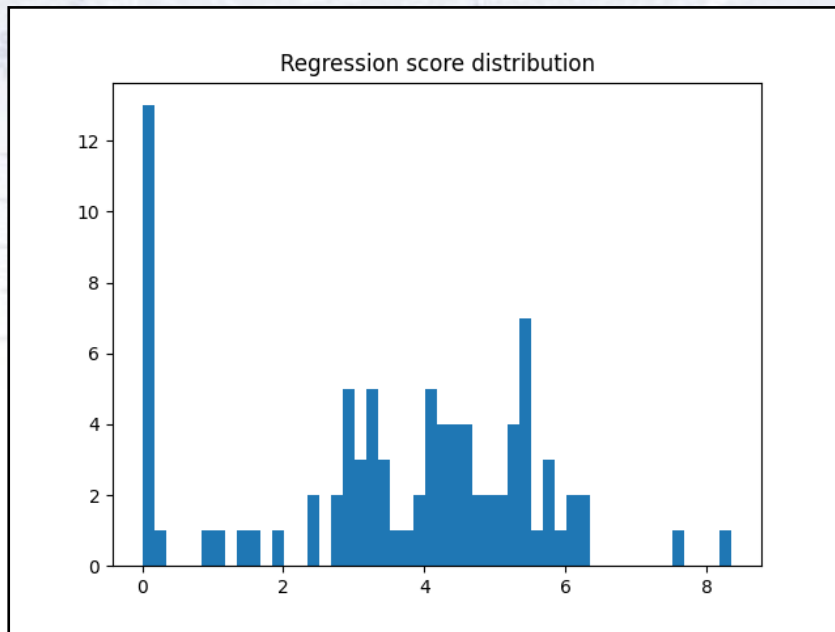
The resulting score distributions

Score distributions for **classification** performance and variable choice:



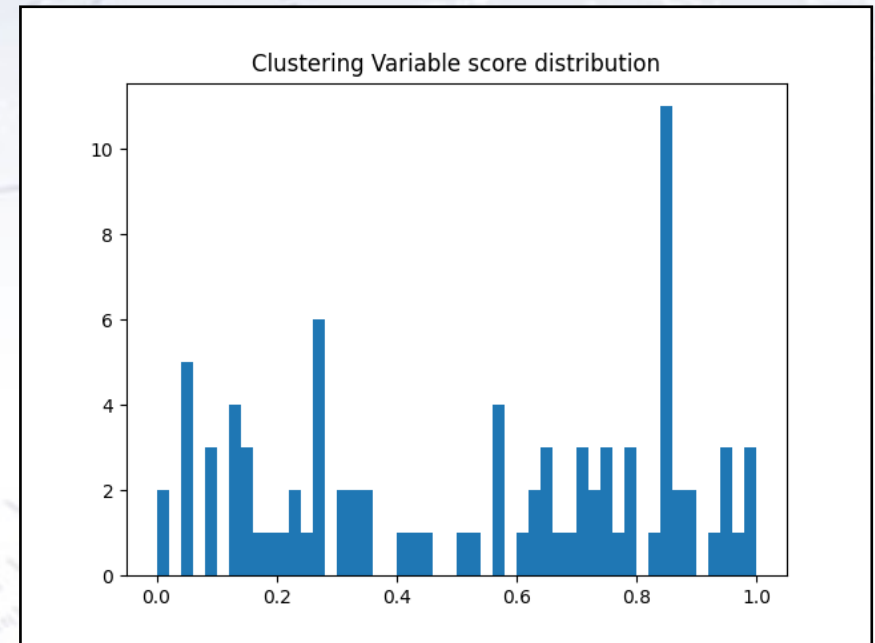
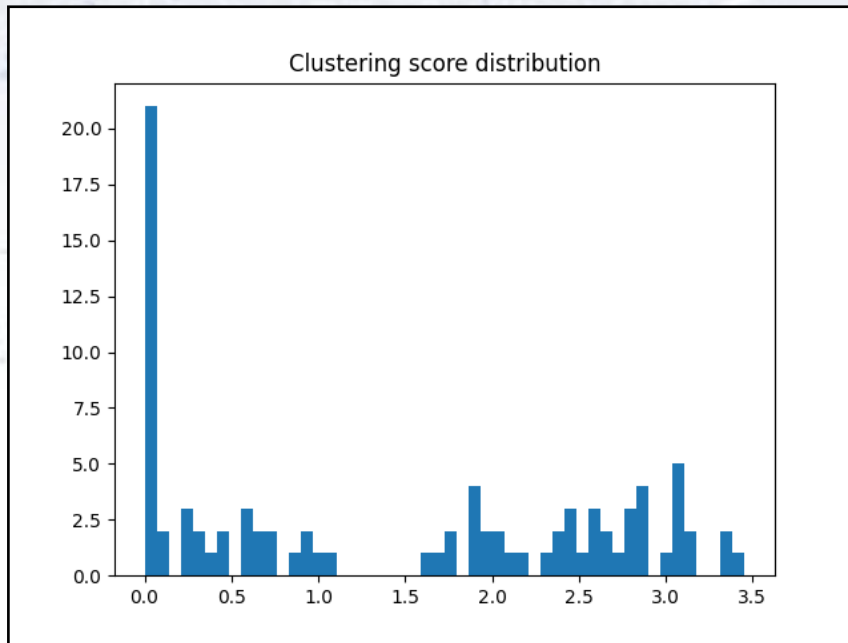
The resulting score distributions

Score distributions for **regression** performance and variable choice:



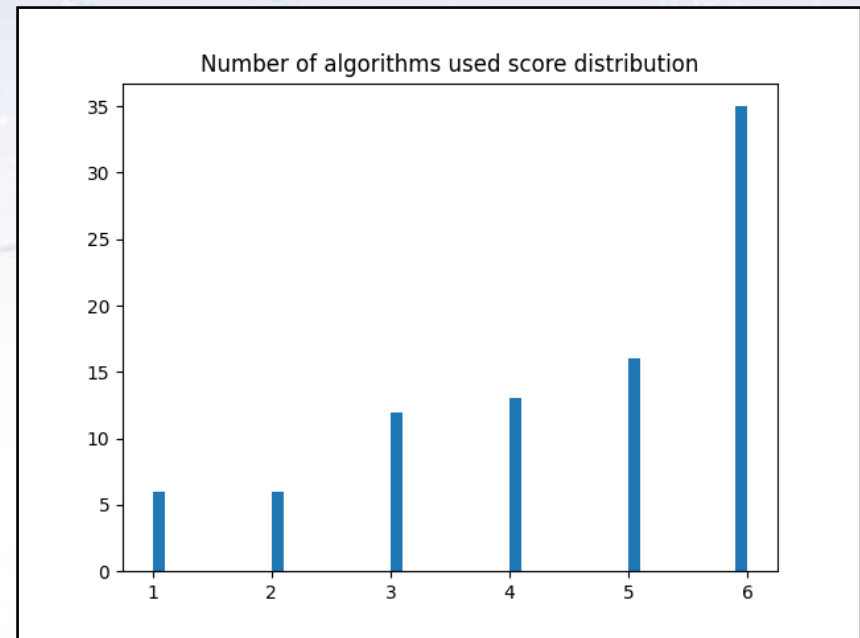
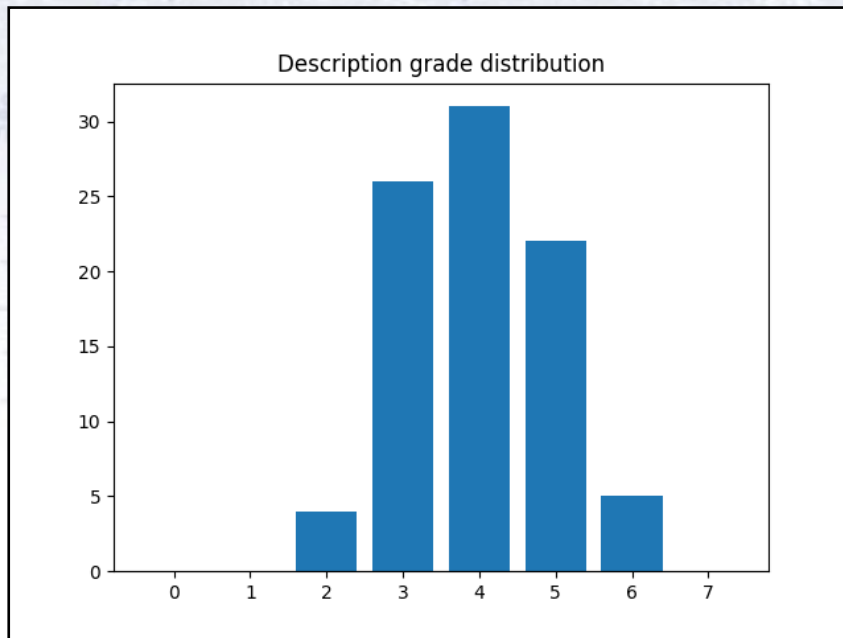
The resulting score distributions

Score distributions for **clustering** performance and variable choice:



The resulting score distributions

The scores for descriptions and number of different algorithms (that work!) are:



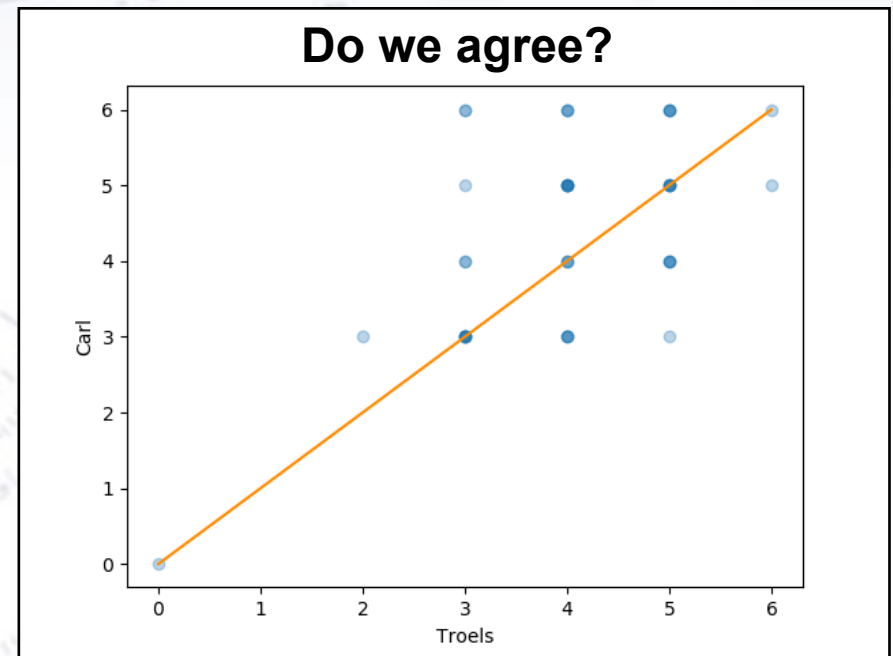
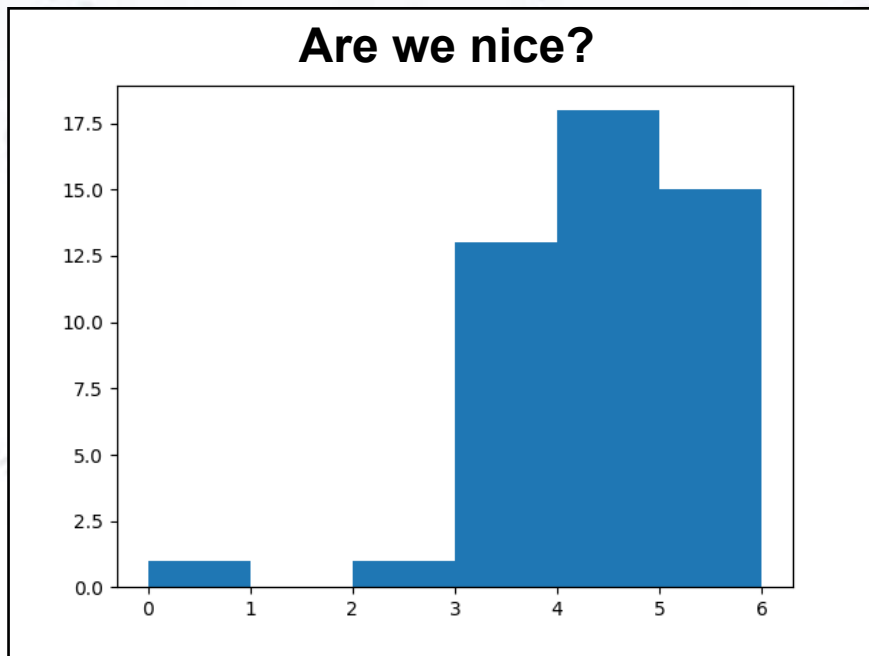
I read several of the “lower scoring” descriptions, but must say that I found them “reasonably acceptable”, so in general the level was high (but don’t do transformation of variables, when using a BDT!).

On algorithms, it was great to see that you both stuck with what you knew, but also explored new algorithms and got them working.

Your description reports

We read through your descriptions, and did a manual scoring (the only) based on choice of algorithms, hyperparameter optimisation, and data division (e.g. cross validation). Each yielded a score of 0-2, giving a total score of 0-6 points.

Numbers from 2021 (where Carl and I did it):



As you can see, we were generally satisfied. The descriptions were short and to the point, and give some insight into your line of thinking and working.



Reporting back to you

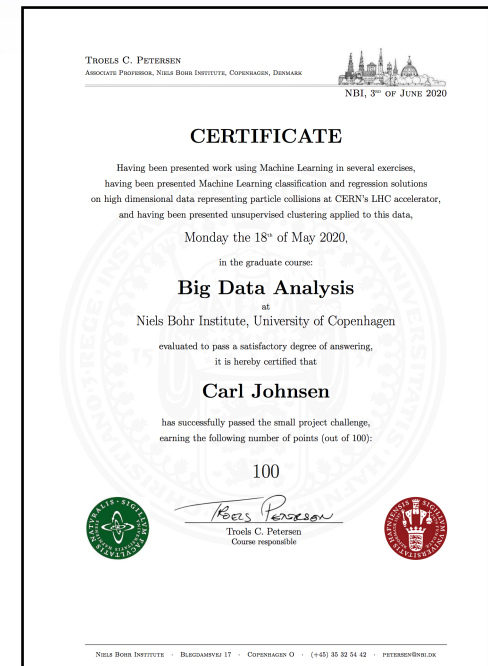
Feedback to you

We have created a small report back to you, which consists of:

- A certificate - for you to be proud of handing in...
- A summary - for you to know how you did...
- A solution scoring with key numbers and illustrations - for you to understand how your model performed.

These are (hopefully) being mailed to you by all of us right now. Please sit down after class and look through them.

Also, don't hesitate to discuss them with your peers. Perhaps you have already done this (great), **but this feedback and reflection is the process through which you learn the most...** please use it.



Classification report

By now you should know what all the different plots and number are...

The solution gave the following metrics:

Metric	Equation	Value
Accuracy	<code>sklearn.metrics.accuracy_score</code>	0.940735
AUC	<code>sklearn.metrics.auc</code>	0.976952
Cross entropy	<code>sklearn.metrics.log_loss</code>	0.153488

The solution produced the following plots:

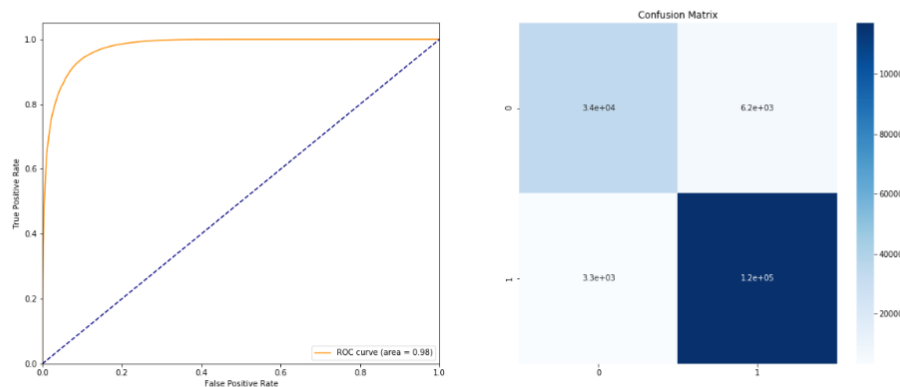


Figure 1: **Left:** ROC curve for the tensorflow2 implementation. The orange curve should be as close to the upper left corner as possible. **Right:** Confusion matrix for the tensorflow2 implementation. The diagonal squares ((0,0) and (1,1)) should have the higher values.

Regression report

The solution gave the following metrics:

Metric	Equation	Value
MAE - Absolute	<code>sklearn.metrics.mean_absolute_error</code>	6953.2194
MAE - Relative	$\sum \frac{ y_p - y_t }{y_t}$	9060.6884
RMS	$\sqrt{\text{mean}((y_p - y_t)^2)}$	14261.8800
RMS 98th percentile	$\sqrt{\text{mean}((y_p - y_t)^2)}$	9238.8301
RMS 90th percentile	$\sqrt{\text{mean}((y_p - y_t)^2)}$	6074.5612
RMS 70th percentile	$\sqrt{\text{mean}((y_p - y_t)^2)}$	4586.3129

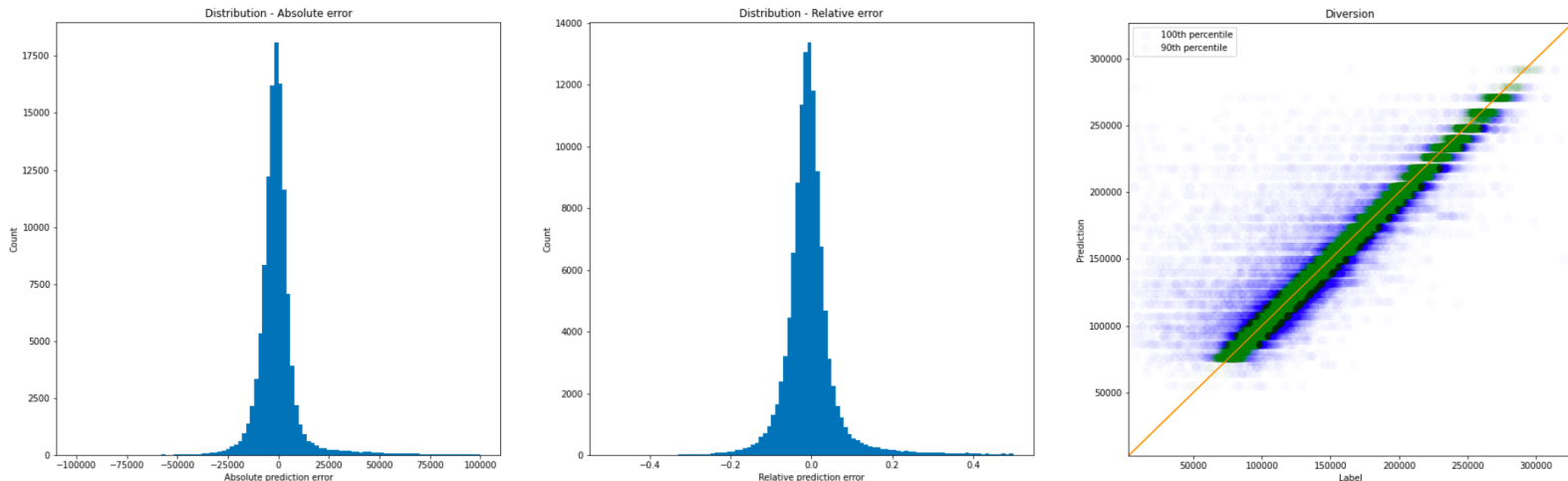


Figure 2: **Upper:** Distribution plots for the xgboost1 implementation. The plots are for absolute error (*Left*) and relative error (*Right*). Both plots should have a tall narrow curve, centered around 0. **Lower:** Diversion plot for the xgboost1 implementation. The dots should be scattered close to the line - especially for the 90th percentile.

Clustering report

The clustering report is necessarily not very detailed, as unsupervised learning carries a great deal of uncertainty on what you're doing.

However, remember the remark by Alexander Nielsen about t-SNE & UMAP, but applied more generally:

"I always start by throwing a clustering algorithm at data, just to see what structures turn up, if any.

Even the latter result tells me something valuable for the further analysis."

clustering - KMeans

The solution produced the following metrics:

Metric	Equation	Value
Accuracy	<code>sklearn.metrics.accuracy_score</code>	0.7492

To compute the accuracy, the following mapping was used, based on the clusters resemblance to electron classification:

Cluster	0	1	2	3
is electron	1	1	1	1

The solution provided the following plot:

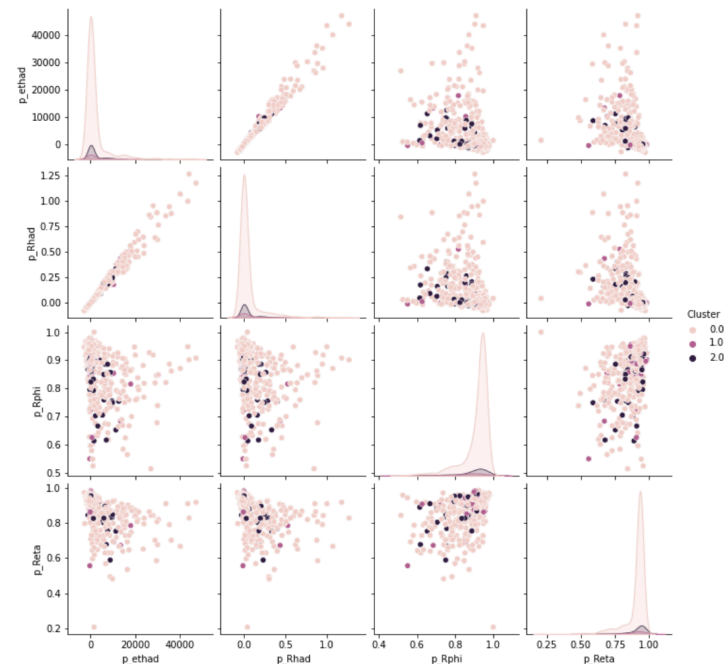
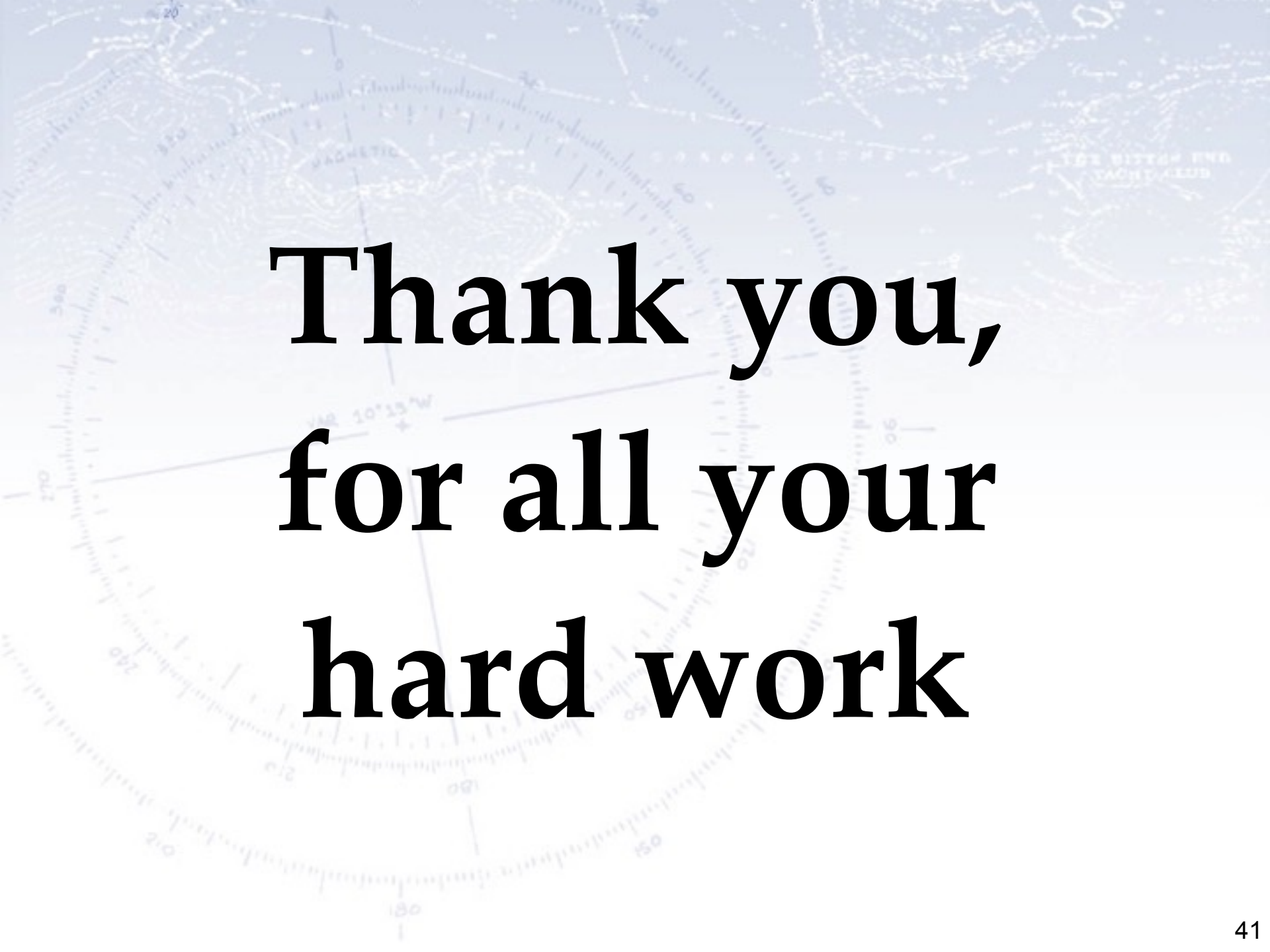


Figure 6: Pairplot for the KMeans implementation. The variables chosen are the top 4 most used variables for clustering. There should be a clear distinction of the clusters.

The background features a light blue map with a prominent compass rose. The word 'MAGNETIC' is visible on the map. In the upper right, there is text that reads '152 BITTEN END TACHT/ALUB'. The map shows various lines of latitude and longitude, along with some faint text and symbols.

**Thank you,
for all your
hard work**