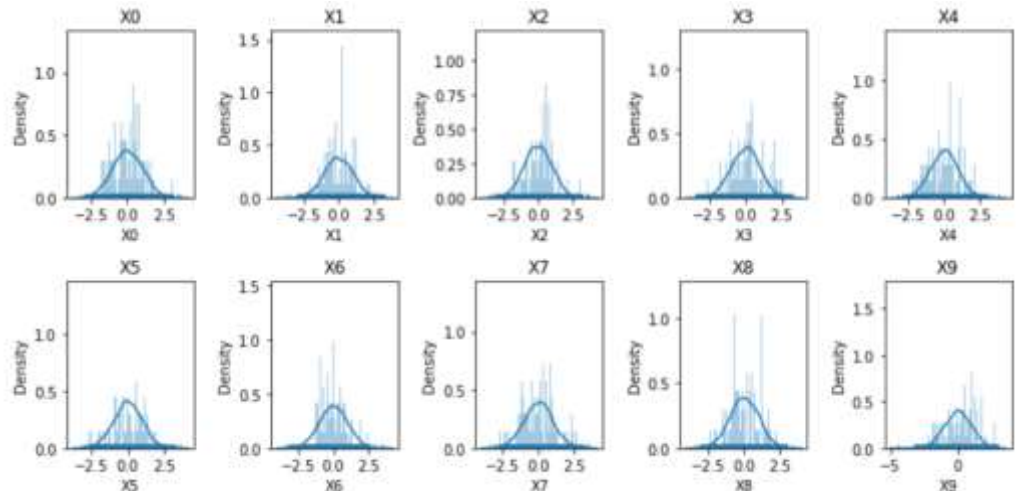


1. Explore the data using visualizations and statistics

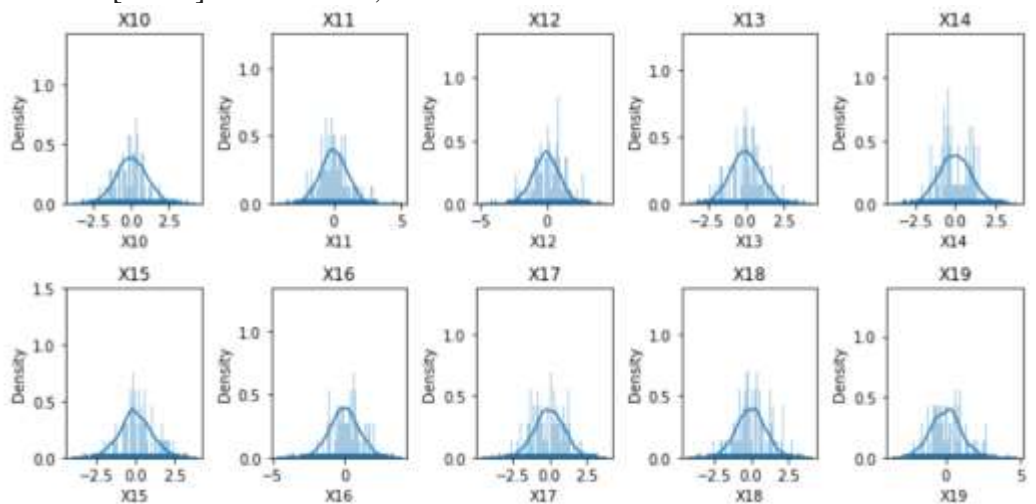
- What are the feature distributions (compute histograms for them)?

The reason why do we need to check the feature distribution is that many ML models are based on the normal distribution assumption. We need to know what type of distributions we are dealing with. After we draw all features distributions, we can see all features are normal distributed which is pretty good.

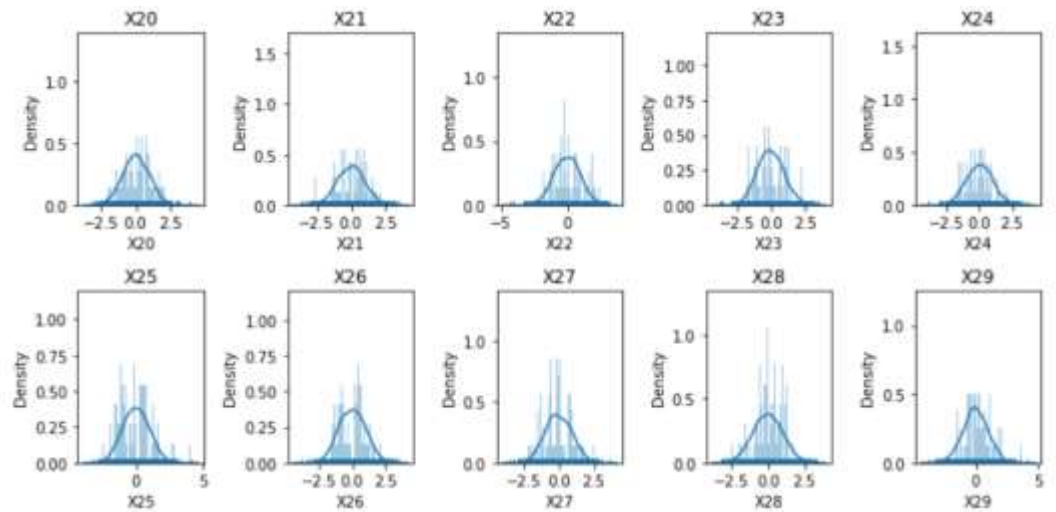
- Columns [0:9] distributions, I set 3000 bins



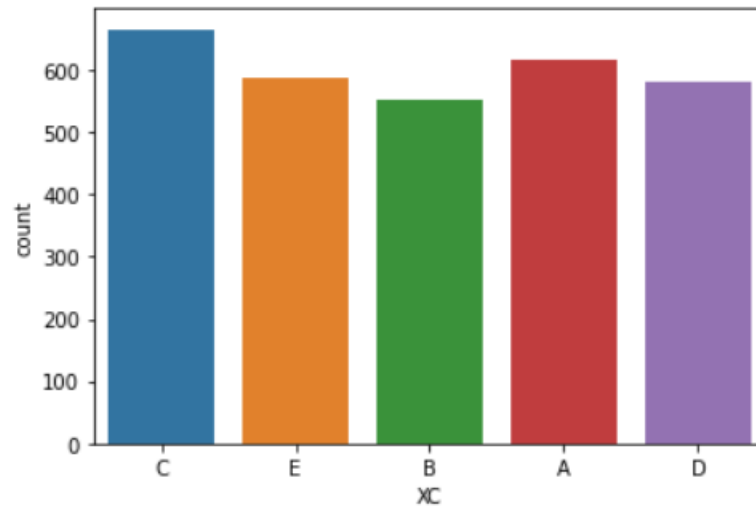
- Columns [10:19] distributions, I set 3000 bins



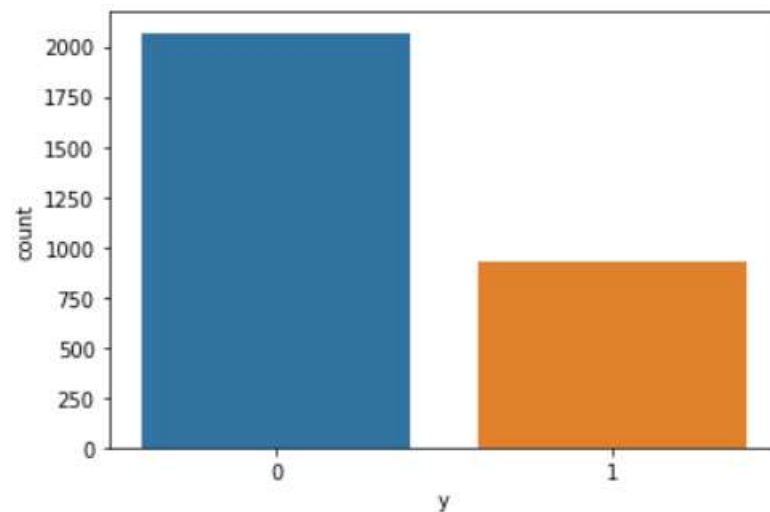
- Columns [20:29] distributions, I set 3000 bins



- Distribution of the last column that is the categorical data

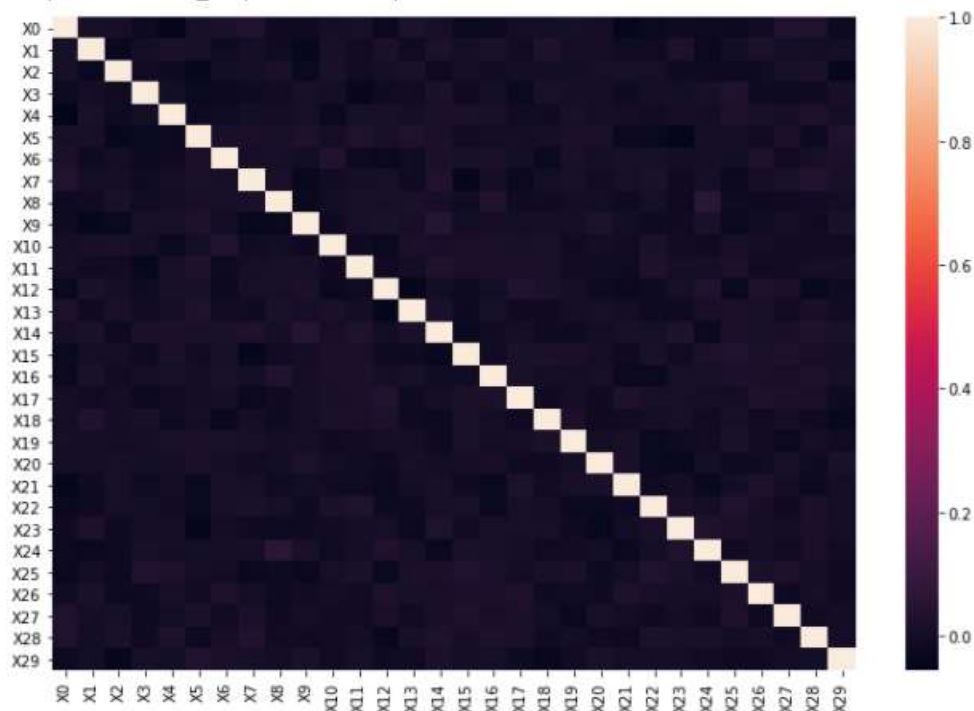


- I think we also need to check the sample distribution. It is quite important too because it affects our predict result. If the data is extremely imbalance, the result will be very bad. Sometimes, we need to think about how to balance the data.



- Is the data correlated?
 - Check Pearson correlation between columns.

I drew a heat map of every feature to check if they are correlated. We can notice that all features are not correlated at all.



I also print out their Pearson coefficients

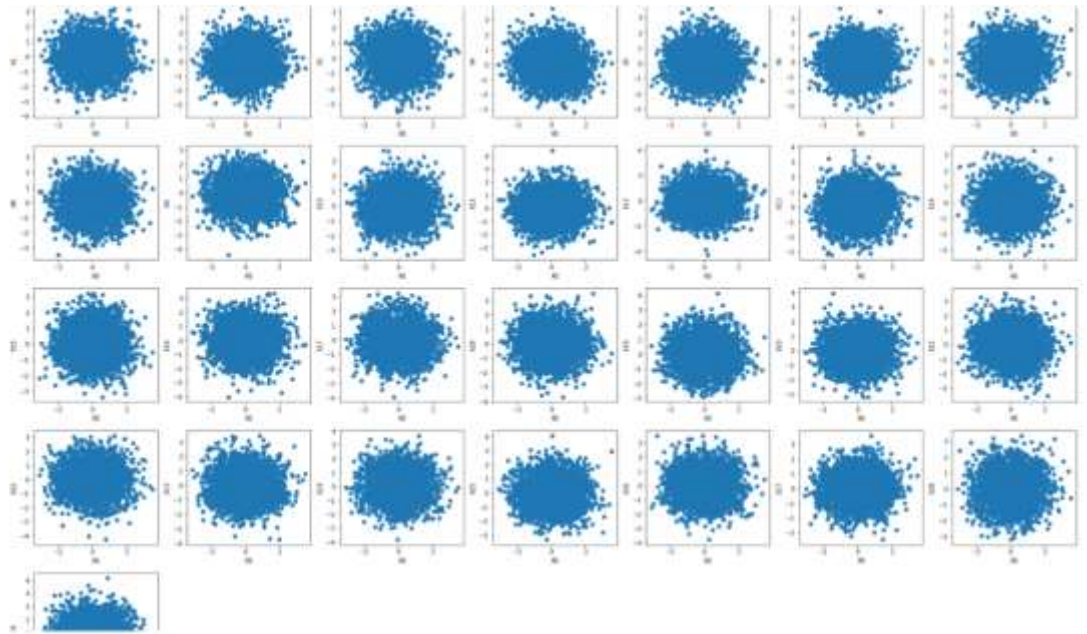
```

# 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29
0 1.00000 0.01196 0.01900 -0.04660 -0.00319 0.01147 0.01634 0.00861 0.01004 -0.01662 0.00106 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
1 0.01196 1.00000 0.03672 0.05908 0.01477 0.01080 -0.01600 0.00361 0.01005 0.00077 0.01174 0.00620 0.01000 0.01022 0.01201 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
2 0.03672 0.05908 1.00000 -0.01029 -0.01020 -0.03000 -0.00457 0.00357 0.01400 -0.00740 0.01443 0.00320 0.00470 0.00005 -0.00073 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
3 0.01477 0.01080 -0.01029 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
4 -0.04660 0.01477 -0.01020 0.00000 1.00000 -0.01530 -0.01530 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
5 0.01147 0.01005 -0.00740 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
6 0.01634 0.00861 0.00361 0.00357 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
7 0.01005 0.00077 0.01400 -0.00740 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
8 0.01174 0.00620 0.01443 0.00320 0.00470 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
9 0.01000 0.01022 0.01201 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
10 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
11 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
12 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
13 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
14 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
15 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
16 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
17 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
18 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
19 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
20 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
21 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
22 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000 0.00000
23 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000 0.00000
24 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000 0.00000
25 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000 0.00000
26 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000 0.00000
27 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 1.00000
28 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
29 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000

```

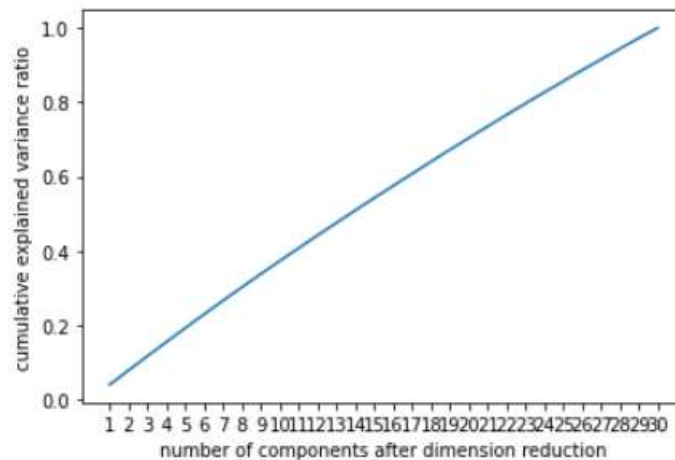
- Try looking at the scatterplot of the dimensions

For this part, it has a lot of plots, so here, I only take part of the plots. The graph shown here checks the correlation of first column and the other columns. We can combine the heat map and this graph, we can draw a conclusion that features are not correlated.



In notebook, we can also run all columns to check any two features are correlated or not.

- Can you remove unnecessary features?
 - Alternatively, can you reduce the dimensionality of the data with PCA? We can check the graph of cumulative explained variance of ratio. We can see that there is no elbow point. So, it is unnecessary to use PCA to reduce the dimensionality but we can check the other method such as lasso or ridge.



- When I tried L1 norm that is feature selection, it can actually reduce the dimension. It reduces to only 20 dimensions.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
0	0.985100	-2.332203	1.044273	1.141718	-1.400005	1.043318	0.723653	0.001453	0.947043	0.048547	1.385370	0.217052	0.028831	0.511403	-0.358514	0.887934	0.210942	0.0	-1.0	0.0	0.0
1	0.938818	1.838048	-0.062190	0.235558	1.520009	-1.411871	0.040112	1.440449	-1.330019	-2.218894	-0.129052	2.303508	0.754007	0.828244	0.200188	1.028799	1.150271	0.0	0.0	0.0	-1.0
2	-1.424465	0.006188	0.115991	1.100002	-1.233711	0.599079	0.242434	0.549367	-0.884077	0.007487	0.845298	1.470555	-1.388526	-2.139636	0.150602	-2.317249	0.153460	0.0	0.0	0.0	1.0
3	-1.038157	0.330467	-0.028101	0.012123	0.573185	0.183184	-0.116562	0.670690	-1.151814	0.278623	2.185787	0.801025	0.305298	-0.968107	0.225715	-0.865773	0.020332	0.0	0.0	0.0	1.0
4	0.177091	1.004133	-0.724015	-0.508380	-0.524431	1.841500	-1.895282	-0.318822	-0.750372	-0.413788	0.868173	0.510647	1.080991	0.883104	0.008553	0.862388	0.803111	0.0	0.0	0.0	1.0
...																					
2995	0.127214	-0.386251	-1.188522	0.846578	-0.152534	-0.182938	-0.080189	-2.278804	1.258887	-0.488413	-1.287280	0.030096	-1.382081	-0.281345	0.255234	-1.145785	0.873201	1.0	0.0	0.0	0.0
2996	0.427700	-1.260329	2.206259	-0.801528	0.406828	-0.408890	-0.157087	-1.79582	-0.781312	-1.815481	2.045879	0.350805	0.558441	-1.433444	-0.865400	0.501944	1.172911	0.0	1.0	0.0	0.0
2997	0.893540	0.403887	-0.118989	0.005295	-0.607212	-0.747290	0.006322	0.064390	1.131446	-2.186588	-0.615433	1.998071	0.024523	-0.304820	0.316073	0.243841	-1.841217	0.0	0.0	0.0	0.0
2998	0.151291	-1.800683	-1.352597	2.221015	0.002388	0.383042	0.325477	0.173198	-0.191733	-0.060871	0.798170	1.570336	0.886517	0.849427	-0.048052	0.142821	1.957203	0.0	0.0	0.0	0.0
2999	-1.351918	-0.818818	-0.417885	-0.888738	-1.255391	0.199008	-0.707243	-0.118135	0.940941	0.588864	0.218222	1.920291	-0.081711	1.105963	1.075838	-1.082485	0.803792	1.0	0.0	0.0	0.0

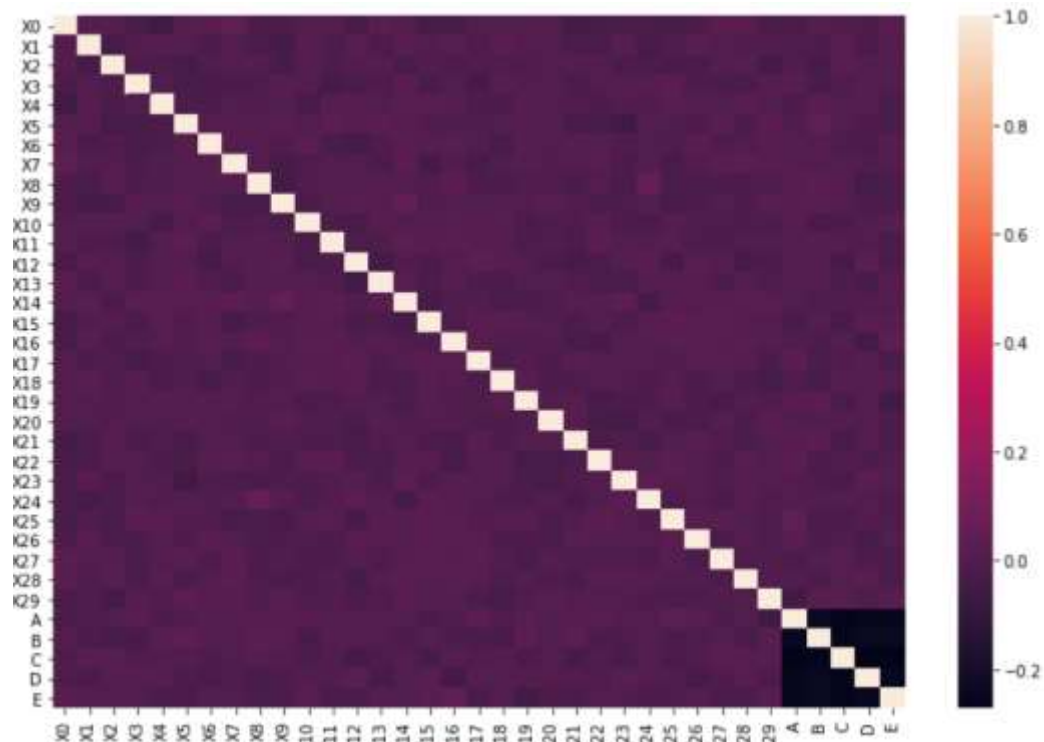
2000 rows x 21 columns

2. Preprocess the data

- One-hot encode categorical variable

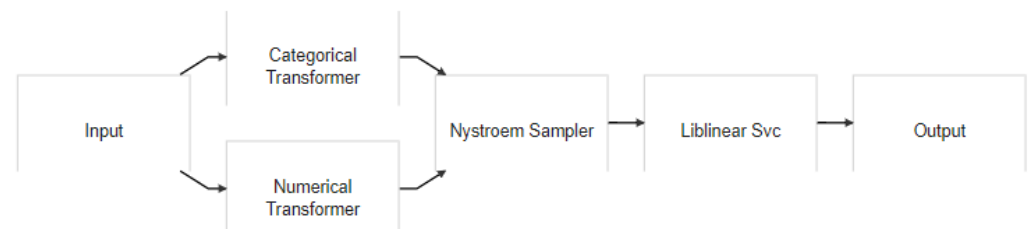
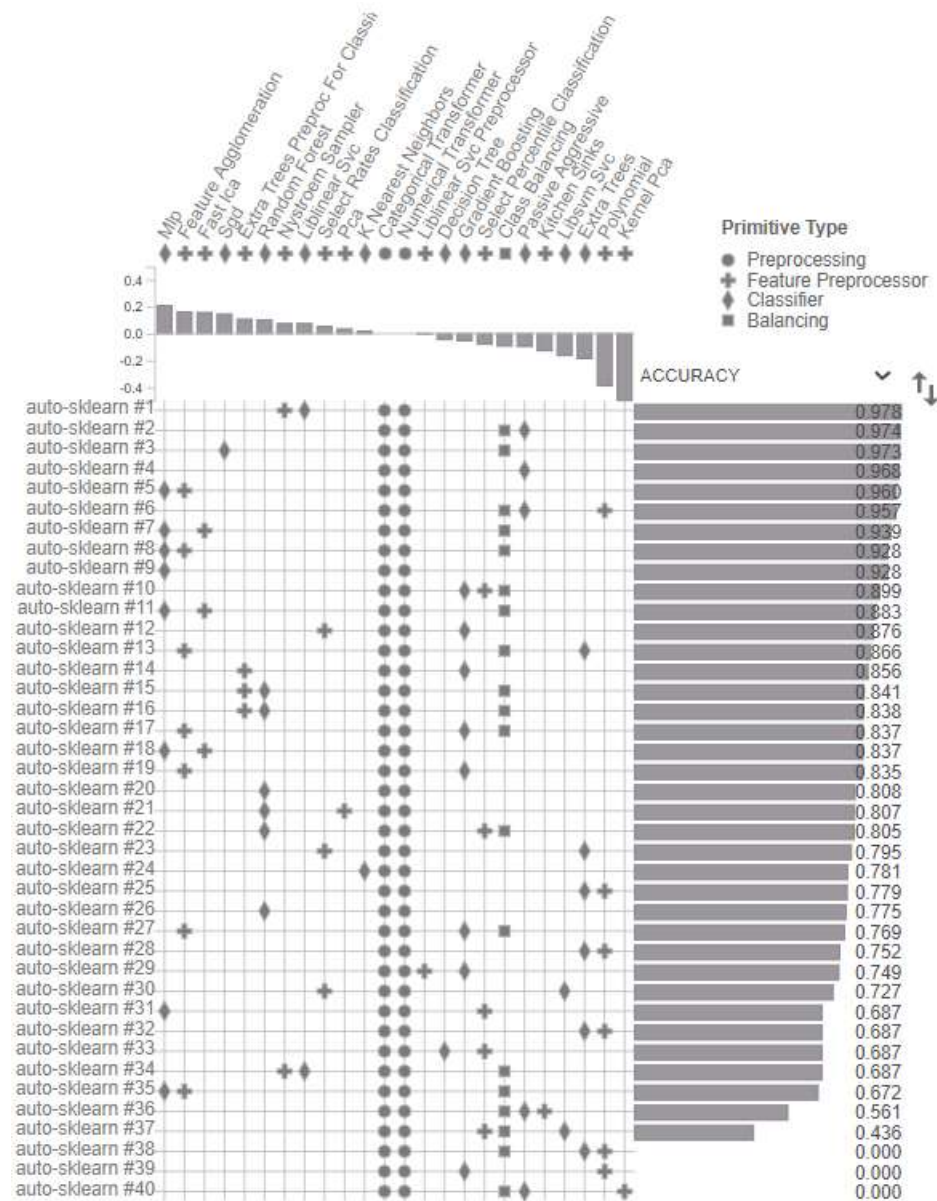
	A	B	C	D	E
0	0.0	0.0	1.0	0.0	0.0
1	0.0	0.0	0.0	0.0	1.0
2	0.0	0.0	0.0	0.0	1.0
3	0.0	0.0	0.0	0.0	1.0
4	0.0	0.0	0.0	0.0	1.0

After I added categorical data, we can check the heat map again to see if any features are correlated or not. We can see that they are not correlated.

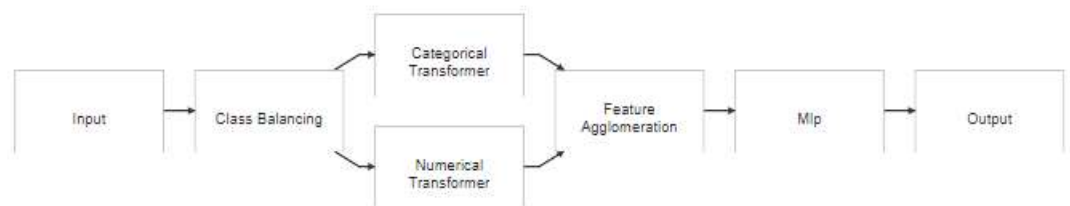
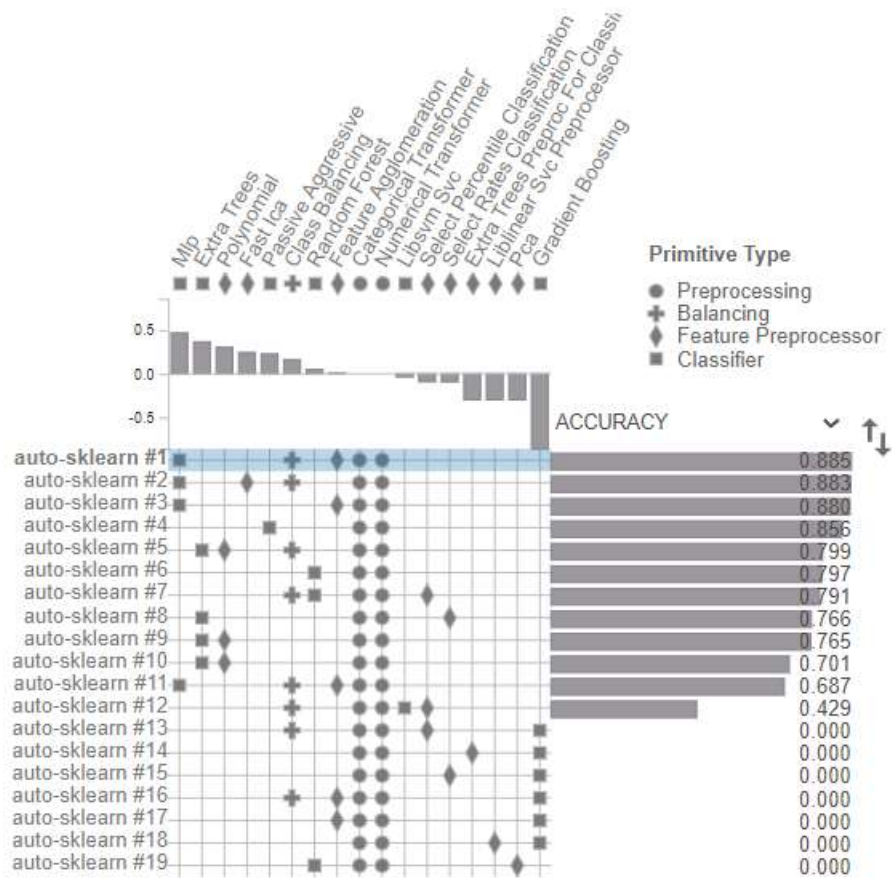


- Check if you need to normalize the data
We can normalize the data to the same range
- ## 3. Solve the classification problem using Auto-Sklearn. Try different dataset.
- With / without categorical feature

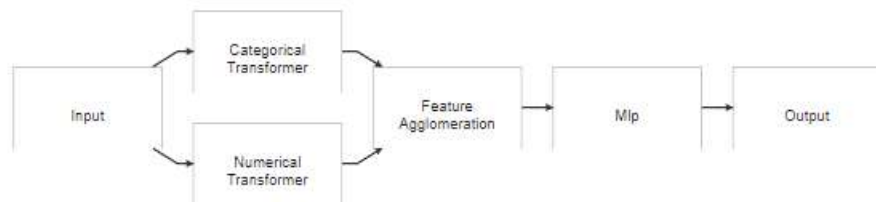
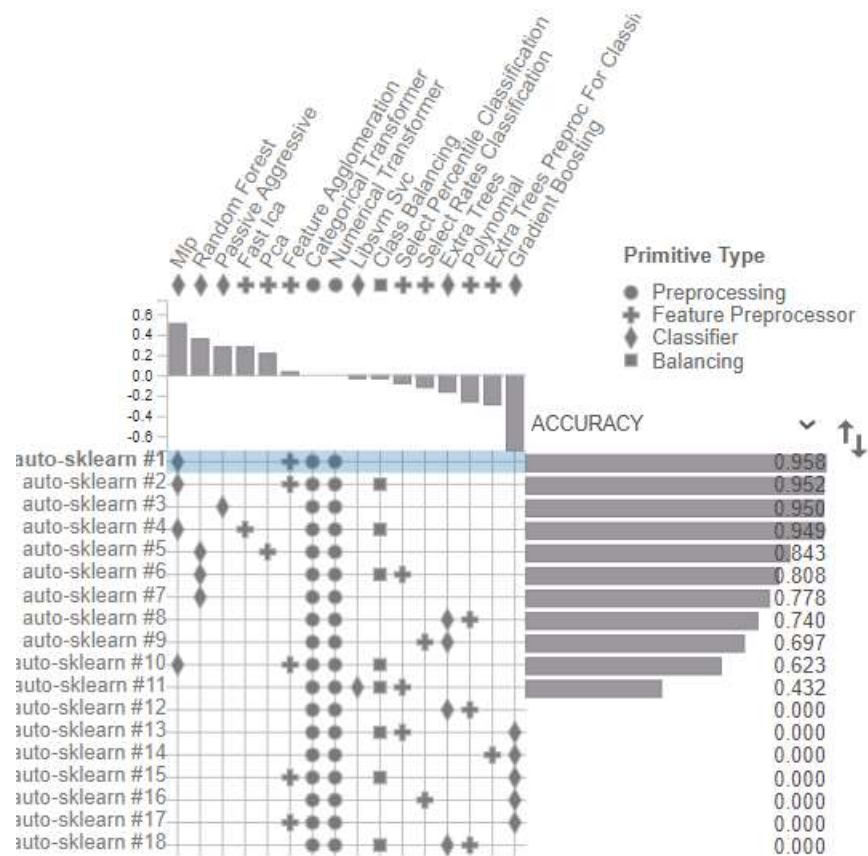
- With categorical data: accuracy is 0.98111112



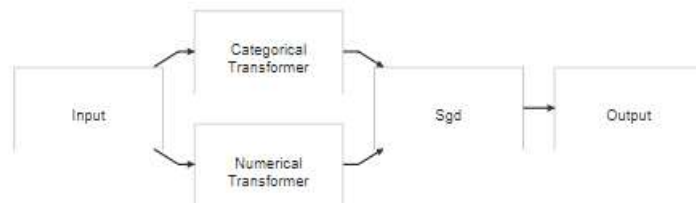
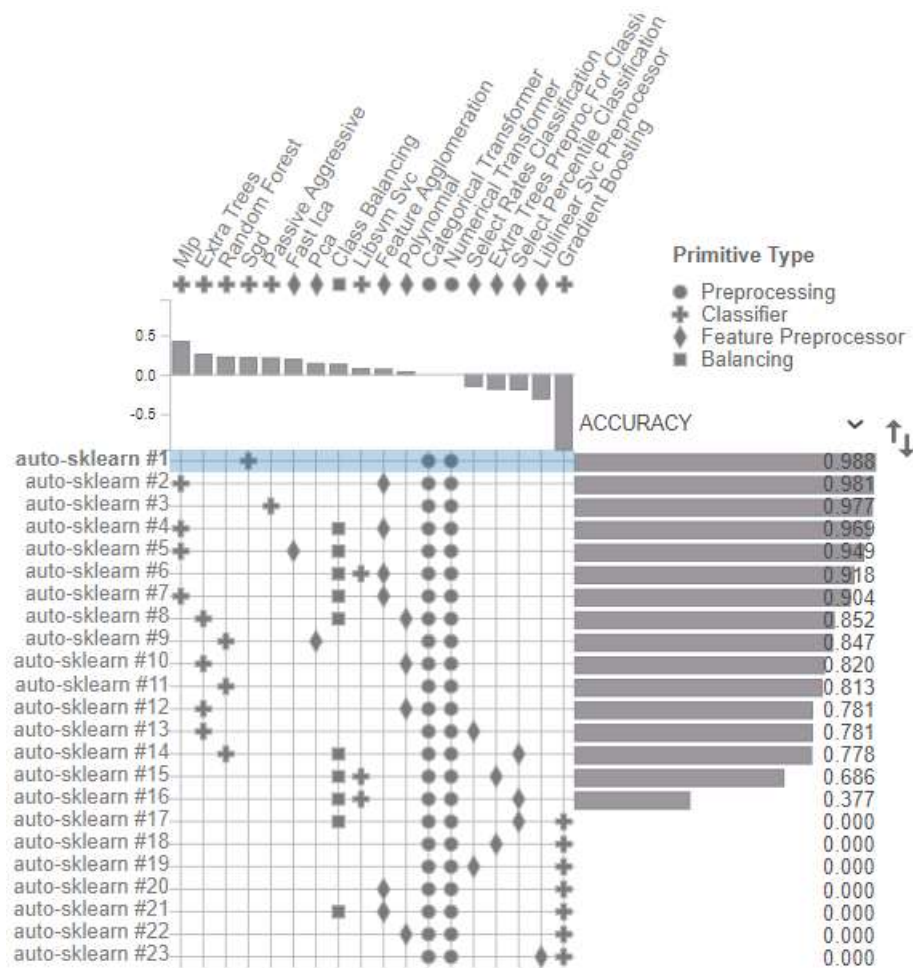
- Without categorical data: accuracy is 0.8844444444444445. It is a very bad result without categorical data.



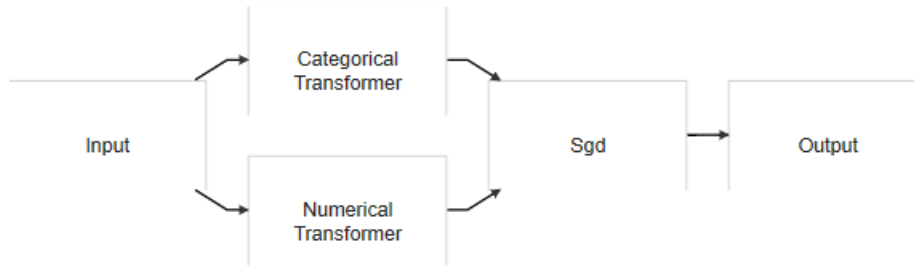
- With / without normalization
 - With normalization: 0.96444444444



- Without normalization: this dataset is same as the data with categorical data because that data is not normalized either.
- With / without feature selection / PCA
 - With feature selection: accuracy is 0.988888888888889



- Without feature selection: this dataset is same as the data with categorical data because this data did not do the feature selection.
4. Explore the models using PipelineProfiler. What primitives perform well for this task?
 Primitive name: `auto_sklearn.primitives.classifier.sgd`.
- Finally, we got the highest accuracy that is 0.9889 with feature selection data. It uses categorical transformer and Numerical transformer to deal with categorical data and numerical data respectively. Then uses the SGD classifier to classify the data. The graph below shows the process of dealing with data.



5. Select the best model out of your experiments (split the data in a 70/30 Training / Validation set)

In autML training, we split the data 70/30 training and validation set.

Finally, we will choose the SGD model to classify the data with parameters:

Primitive Name: `auto_sklern.primitives.classifier.sgd`

```

{
  alpha:0.00023,
  Average:True,
  Fit_intercept:True,
  Learning_rate:optimal,
  Loss:log,
  Penalty:l1,
  Tol: 0.00001
}
  
```

To deal with the categorical data:

Primitive name: `auto_sklern.primitives.data_preprocessing.categorical_transformer`
with paramters:

```

{
  Categorical_encoding: one_hot_encoding,
  Categorical_coalescence: no_coalescence
}
  
```

To deal with the numerical data:

Primitive Name: `auto_sklern.primitives.data_preprocessing.numerical_transformer`
with parameters:

```

{
  Imputation_strategy: mean,
  Rescaling: standarize
}
  
```