

You have **2** free stories left this month. Sign up and get an extra one for free.



Source: Pixabay

DS ALGORITHMS

The 5 Feature Selection Algorithms every Data Scientist should know

Bonus: What makes a good footballer great?



Rahul Agarwal

Jul 27, 2019 · 7 min read ★

Data Science is the study of algorithms.

I grapple through with many algorithms on a day to day basis, so I thought of listing some of the most common and most used algorithms one will end up using in this new DS Algorithm series.

How many times it has happened when you create a lot of features and then you need to come up with ways to reduce the number of features.

We sometimes end up using correlation or tree-based methods to find out the important features.

Can we add some structure to it?

This post is about some of the most common feature selection techniques one can use while working with data.

. . .

Why Feature Selection?

Before we proceed, we need to answer this question. Why don't we give all the features to the ML algorithm and let it decide which feature is important?

So there are three reasons why we don't:

1. Curse of dimensionality — Overfitting



As the dimensionality of the features space increases, the number Configurations can grow exponentially, and thus the number of configurations covered by an observation decreases.

Chris Albon

[Source](#)

If we have more columns in the data than the number of rows, we will be able to fit our training data perfectly, but that won't generalize to the new samples. And thus we learn absolutely nothing.

2. Occam's Razor:

We want our *models to be simple* and explainable. We lose explainability when we have a lot of features.

3. Garbage In Garbage out:

Most of the times, we will have many non-informative features. For Example, Name or ID variables. *Poor-quality input will produce Poor-Quality output.*

Also, a large number of features make a model bulky, time-taking, and harder to implement in production.

• • •

So What do we do?





We select only useful features.

Fortunately, Scikit-learn has made it pretty much easy for us to make the feature selection. There are a lot of ways in which we can think of feature selection, but most feature selection methods can be divided into three major buckets

- **Filter based:** We specify some metric and based on that filter features. An example of such a metric could be correlation/chi-square.
- **Wrapper-based:** Wrapper methods consider the selection of a set of features as a search problem. Example: Recursive Feature Elimination
- **Embedded:** Embedded methods use algorithms that have built-in feature selection methods. For instance, Lasso and RF have their own feature selection methods.

• • •

So enough of theory let us start with our five feature selection methods.

We will try to do this using a dataset to understand it better.

I am going to be using a football player dataset to find out *what makes a good player great?*

Don't worry if you don't understand football terminologies. I will try to keep it at a minimum.

Here is the Kaggle Kernel with the code to try out yourself.

• • •

Some simple Data Preprocessing

We have done some basic preprocessing such as removing Nulls and one hot encoding.
And converting the problem to a classification problem using:

```
y = traindf['Overall'] >= 87
```

Here we use High Overall as a proxy for a great player.

Our dataset(X) looks like below and has 223 columns.

	Crossing	Finishing	ShortPassing	Dribbling	LongPassing	BallControl	Acceleration	SprintSpeed	Agility	Stamina	...	Nationality_Uganda	Nationality_Ukraine	Nationality_United Arab Emirates	Nationality_United States
0	84.0	95.0	90.0	97.0	87.0	96.0	91.0	86.0	91.0	72.0	...	0	0	0	0
1	84.0	94.0	81.0	88.0	77.0	94.0	89.0	91.0	87.0	88.0	...	0	0	0	0
2	79.0	87.0	84.0	96.0	78.0	95.0	94.0	90.0	96.0	81.0	...	0	0	0	0
3	17.0	13.0	50.0	18.0	51.0	42.0	57.0	58.0	60.0	43.0	...	0	0	0	0
4	93.0	82.0	92.0	86.0	91.0	91.0	78.0	76.0	79.0	90.0	...	0	0	0	0

train Data X

• • •

1. Pearson Correlation

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}}$$

This is a filter-based method.

We check the *absolute value of the Pearson's correlation* between the target and numerical features in our dataset. We keep the top n features based on this criterion.

```

1 def cor_selector(X, y,num_feats):
2     cor_list = []
3     feature_name = X.columns.tolist()
4     # calculate the correlation with y for each feature
5     for i in X.columns.tolist():
6         cor = np.corrcoef(X[i], y)[0, 1]
7         cor_list.append(cor)

```

```

8   # replace NaN with 0
9   cor_list = [0 if np.isnan(i) else i for i in cor_list]
10  # feature name
11  cor_feature = X.iloc[:,np.argsort(np.abs(cor_list))[-num_feats:]].columns.tolist()
12  # feature selection? 0 for not select, 1 for select
13  cor_support = [True if i in cor_feature else False for i in feature_name]
14  return cor_support, cor_feature
15 cor_support, cor_feature = cor_selector(X, y,num_feats)
16 print(str(len(cor_feature)), 'selected features')

```

correlation.py hosted with ❤ by GitHub

[view raw](#)

• • •

2. Chi-Squared

This is another filter-based method.

In this method, we calculate the chi-square metric between the target and the numerical variable and only select the variable with the maximum chi-squared values.

CHI-SQUARED FOR FEATURE SELECTION

To use χ^2 for feature selection, we calculate χ^2 between each feature and the target, and select the desired number of features with the best χ^2 scores.

The intuition is that if a feature is independent to the target it is uninformative for classifying observations.

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

of observations
 in class i ↗
 ↓
 # of expected
 observations in class i
 if there was no relationship
 between the feature
 and target.

[Source](#)

Let us create a small example of how we calculate the chi-squared statistic for a sample.

So let's say we have 75 Right-Forwards in our dataset and 25 Non-Right-Forwards. We observe that 40 of the Right-Forwards are good, and 35 are not good. Does this signify that the player being right forward affects the overall performance?

	Good	NotGood		
Not RightForward	20	5	25	
Rightforward	40	35	75	
	60	40	100	

	Good	NotGood		
Not RightForward	15	10	25	
Rightforward	45	25	75	
	60	40	100	

Observed and Expected Counts

We calculate the chi-squared value:

To do this, we first find out the values we would expect to be falling in each bucket if there was indeed independence between the two categorical variables.

This is simple. We multiply the row sum and the column sum for each cell and divide it by total observations.

so Good and NotRightforward Bucket Expected value = 25(Row Sum)*60(Column Sum)/100(Total Observations)

Why is this expected? Since there are 25% notRightforwards in the data, we would expect 25% of the 60 good players we observed in that cell. Thus 15 players.

Then we could just use the below formula to sum over all the 4 cells:

$$\chi_c^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

I won't show it here, but the chi-squared statistic also works in a hand-wavy way with non-negative numerical and categorical features.

We can get chi-squared features from our dataset as:

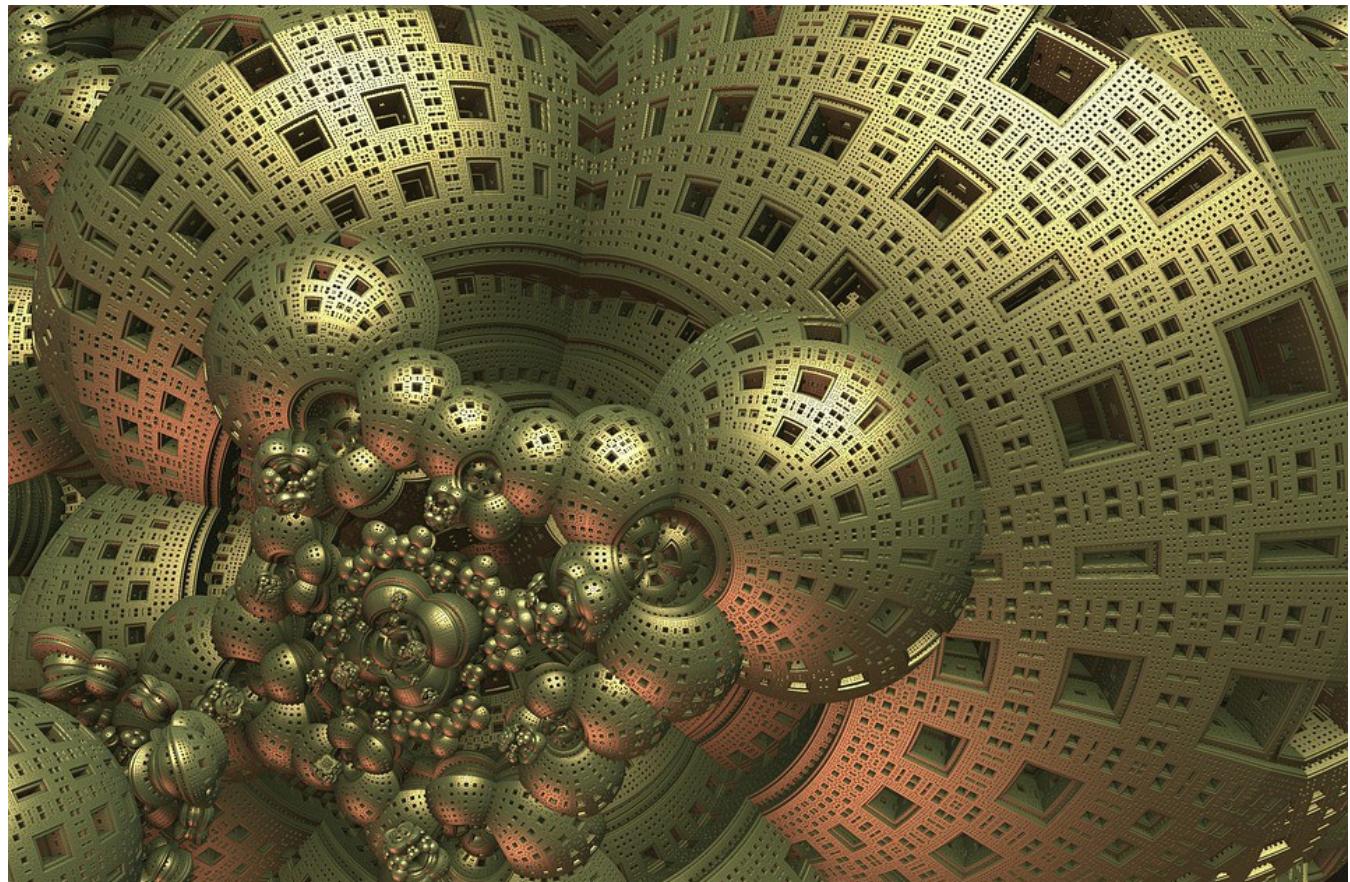
```
1  from sklearn.feature_selection import SelectKBest
2  from sklearn.feature_selection import chi2
3  from sklearn.preprocessing import MinMaxScaler
4  X_norm = MinMaxScaler().fit_transform(X)
5  chi_selector = SelectKBest(chi2, k=num_feats)
6  chi_selector.fit(X_norm, y)
7  chi_support = chi_selector.get_support()
8  chi_feature = X.loc[:,chi_support].columns.tolist()
9  print(str(len(chi_feature)), 'selected features')
```

chi_square.py hosted with ❤ by GitHub

[view raw](#)

• • •

3. Recursive Feature Elimination



This is a wrapper based method. As I said before, wrapper methods consider the selection of a set of features as a search problem.

From `sklearn` Documentation:

*The goal of recursive feature elimination (RFE) is to select features by **recursively considering smaller and smaller sets of features**. First, the estimator is trained on the initial set of features and the importance of each feature is obtained either through a `coef_` attribute or through a `feature_importances_` attribute. Then, the least important features are pruned from current set of features. That procedure is recursively repeated on the pruned set until the desired number of features to select is eventually reached.*

As you would have guessed, we could use any estimator with the method. In this case, we use `LogisticRegression`, and the RFE observes the `coef_` attribute of the `LogisticRegression` object

```

1  from sklearn.feature_selection import RFE
2  from sklearn.linear_model import LogisticRegression
3  rfe_selector = RFE(estimator=LogisticRegression(), n_features_to_select=num_feats, step=10, vert
4  rfe_selector.fit(X_norm, y)
5  rfe_support = rfe_selector.get_support()
6  rfe_feature = X.loc[:,rfe_support].columns.tolist()
7  print(str(len(rfe_feature)), 'selected features')

```

rfe.py hosted with ❤ by GitHub

[view raw](#)

4. Lasso: SelectFromModel



- Lasso regression uses L1 norm as regularizer.

$$\alpha \sum_{i=1}^k |w_i|$$

alpha' != ~parameters'

- Unlike ridge regression, lasso's norm regularizer drives parameters to zero.
- Higher the value of alpha, the fewer features have non-zero values.

Chris Albon

[Source](#)

This is an Embedded method. As said before, Embedded methods use algorithms that have built-in feature selection methods.

For example, Lasso and RF have their own feature selection methods. Lasso Regularizer forces a lot of feature weights to be zero.

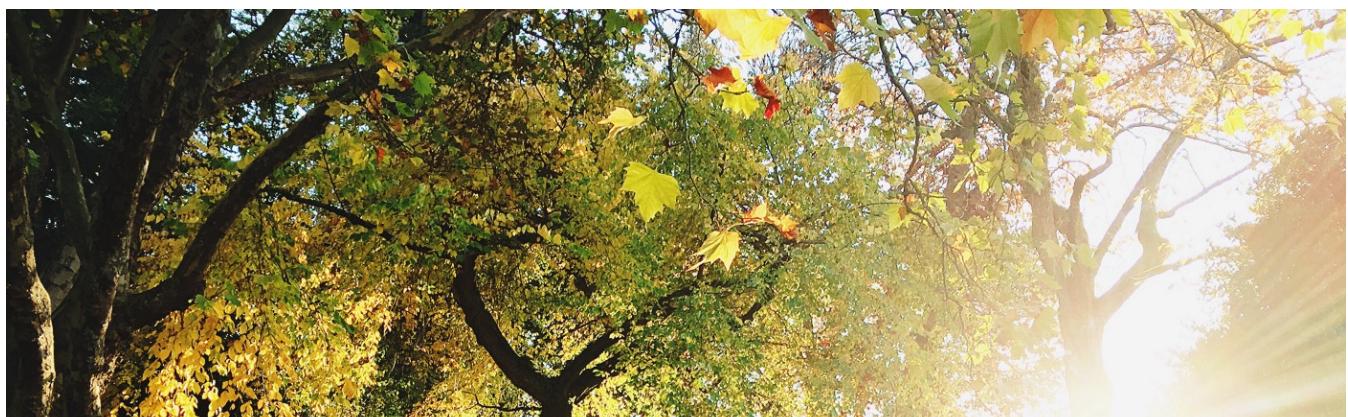
Here we use Lasso to select variables.

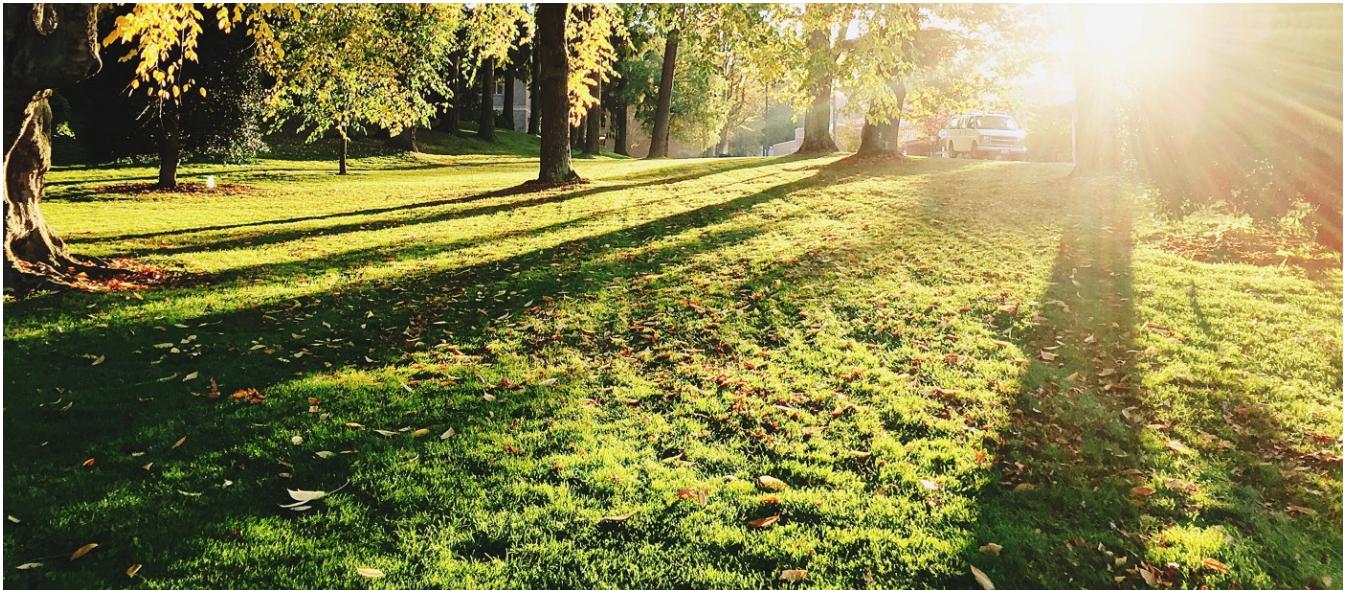
```

1 from sklearn.feature_selection import SelectFromModel
2 from sklearn.linear_model import LogisticRegression
3
4 embeded_lr_selector = SelectFromModel(LogisticRegression(penalty="l1"), max_features=num_feats)
5 embeded_lr_selector.fit(X_norm, y)
6
7 embeded_lr_support = embeded_lr_selector.get_support()
8 embeded_lr_feature = X.loc[:,embeded_lr_support].columns.tolist()
9 print(str(len(embeded_lr_feature)), 'selected features')

```

5. Tree-based: SelectFromModel





This is an Embedded method. As said before, Embedded methods use algorithms that have built-in feature selection methods.

We can also use RandomForest to select features based on feature importance.

We calculate feature importance using node impurities in each decision tree. In Random forest, the final feature importance is the average of all decision tree feature importance.

```

1  from sklearn.feature_selection import SelectFromModel
2  from sklearn.ensemble import RandomForestClassifier
3
4  embeded_rf_selector = SelectFromModel(RandomForestClassifier(n_estimators=100), max_features=num_features)
5  embeded_rf_selector.fit(X, y)
6
7  embeded_rf_support = embeded_rf_selector.get_support()
8  embeded_rf_feature = X.loc[:,embeded_rf_support].columns.tolist()
9  print(str(len(embeded_rf_feature)), 'selected features')

```

[rf.py](#) hosted with ❤ by GitHub

[view raw](#)

We could also have used a LightGBM. Or an XGBoost object as long it has a `feature_importances_` attribute.

```

1  from sklearn.feature_selection import SelectFromModel
2  from lightgbm import LGBMClassifier
3
4  lgbc=LGBMClassifier(n_estimators=500, learning_rate=0.05, num_leaves=32, colsample_bytree=0.2,
5                      reg_alpha=3, reg_lambda=1, min_split_gain=0.01, min_child_weight=40)

```

```
6  
7     embeded_lgb_selector = SelectFromModel(lgbc, max_features=num_feats)  
8     embeded_lgb_selector.fit(X, y)  
9  
10    embeded_lgb_support = embeded_lgb_selector.get_support()  
11    embeded_lgb_feature = X.loc[:,embeded_lgb_support].columns.tolist()  
12    print(str(len(embeded_lgb_feature)), 'selected features')
```

Bonus



Why use one, when we can have all?

The answer is sometimes it won't be possible with a lot of data and time crunch.

But whenever possible, why not do this?

```
1  # put all selection together  
2  feature_selection_df = pd.DataFrame({'Feature':feature_name, 'Pearson':cor_support, 'Chi-2':chi_<pre></pre>  
3  'Random Forest':embeded_rf_support, 'LightGBM':embeded_lgb_support})
```

```

4 # count the selected times for each feature
5 feature_selection_df['Total'] = np.sum(feature_selection_df, axis=1)
6 # display the top 100
7 feature_selection_df = feature_selection_df.sort_values(['Total','Feature'] , ascending=False)
8 feature_selection_df.index = range(1, len(feature_selection_df)+1)
9 feature_selection_df.head(num_feats)

```

full_selector.py hosted with ❤ by GitHub

[view raw](#)

	Feature	Pearson	Chi-2	RFE	Logistics	Random Forest	LightGBM	Total
1	Reactions	True	True	True	True	True	True	6
2	LongPassing	True	True	True	True	True	True	6
3	Volleys	True	True	True	False	True	True	5
4	ShortPassing	True	True	True	False	True	True	5
5	Finishing	True	True	True	False	True	True	5
6	BallControl	True	True	True	False	True	True	5
7	Strength	True	False	True	False	True	True	4
8	Position_CM	False	True	True	True	False	True	4
9	Nationality_Uruguay	True	True	True	True	False	False	4
10	Nationality_Slovenia	True	True	True	True	False	False	4
11	Nationality_Belgium	True	True	True	True	False	False	4
12	LongShots	True	True	False	False	True	True	4
13	FKAccuracy	True	True	False	False	True	True	4
14	Balance	False	False	True	True	True	True	4
15	Aggression	False	False	True	True	True	True	4
16	Weak Foot	True	False	False	False	True	True	3
17	Stamina	True	False	False	False	True	True	3
18	SprintSpeed	True	False	False	False	True	True	3
19	ShotPower	True	False	False	False	True	True	3
20	Preferred Foot_Right	False	False	True	True	False	True	3

We check if we get a feature based on all the methods. In this case, as we can see `Reactions` and `LongPassing` are excellent attributes to have in a high rated player. And as expected `Ballcontrol` and `Finishing` occupy the top spot too.

• • •

Conclusion

Feature engineering and feature selection are critical parts of any machine learning pipeline.

We strive for accuracy in our models, and one cannot get to a good accuracy without revisiting these pieces again and again.

In this article, I tried to explain some of the most used feature selection techniques as well as my workflow when it comes to feature selection.

I also tried to provide some intuition into these methods, but you should probably try to see more into it and try to incorporate these methods into your work.

Do read my post on feature engineering too if you are interested.

• • •

If you want to learn more about Data Science, I would like to call out this *excellent course* by Andrew Ng. This was the one that got me started. Do check it out.

• • •

Thanks for the read. I am going to be writing more beginner-friendly posts in the future too. Follow me up at **Medium** or Subscribe to my **blog** to be informed about them. As always, I welcome feedback and constructive criticism and can be reached on Twitter @mlwhiz.

