

✦ Member-only story

How Probability Calibration Works



Mattia Cinelli · Following

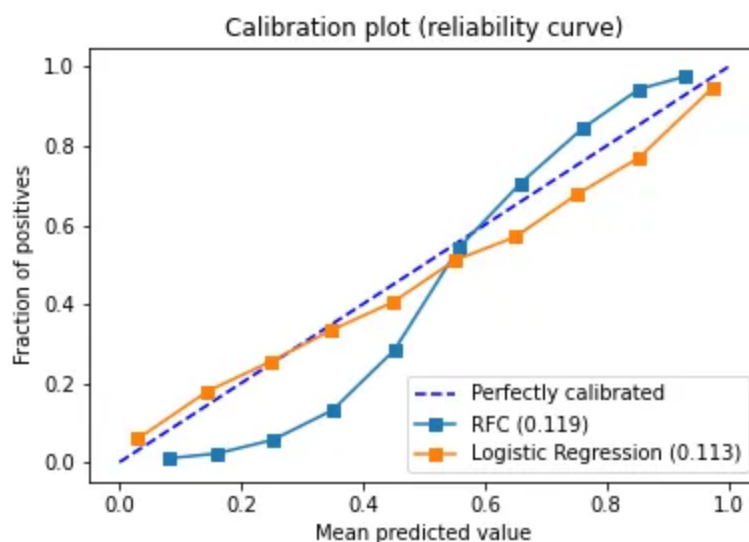
Published in Analytics Vidhya · 6 min read · May 28, 2020



112



3



Probability calibration is the process of calibrating an ML model to return the *true* likelihood of an event. This is necessary when we need the probability of the event in question rather than its classification.

Imagine that you have two models to predict rainy days, Model A and Model B. Both models have an accuracy of 0.8. And indeed, for every 10 rainy days,

both mislabelled two days.

But if we look at the probability connected to each prediction, we can see that Model A reports a probability of 80%, while Model B of 100%.

This means that model B is 100% sure that it will rain, even when it will not, while model A is only 80% sure. It appears that model B is overconfident with its prediction, while model A is more cautious.

And it's this level of confidence in predictions that makes Model A a more reliable model with respect to Model B; Model A is better despite the two models having the same accuracy.

Model B offers a more yes-or-no prediction, while Model A tells us the *true likelihood of the event*. And in real life, when we look at the weather forecast, we get the prediction and its probability, leaving us to decide if, for example, a 30% risk of rain is acceptable or not.

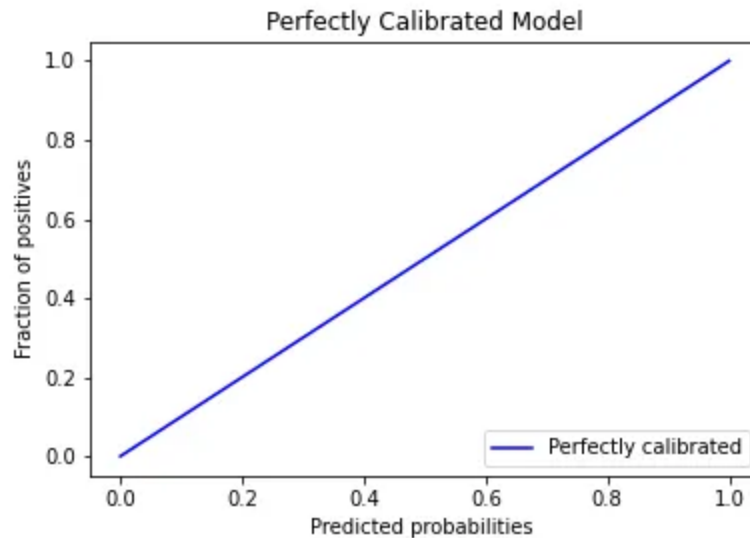
In data science, you will not always care about the probability associated with each prediction, but there are cases in which we want to know the probability.

Think, for example, in the medical sector, where a doctor has to decide if a patient labelled as needing an operation would really benefit from it. Or how safe an investment is, and so on.

Therefore, in these situations it is important to have a model with good accuracy but also that is well **calibrated**!

A model would be perfectly calibrated when the predicted probability reflects the true likelihood of the prediction [1].

Thus, like in our Model A example, the probability of 0.8 is equal to the accuracy for all those events with the same prediction probability. Therefore, there will be an $x=y$ relationship between the predicted probabilities and the fraction of positive entries in the data. This can be visualized with the following plot:



The perfectly calibrated line of an ideal model. A model is perfectly calibrated if, for any p , a prediction of a class with confidence p is correct $100 \cdot p\%$ of the time.[2]

Naturally, our models will rarely be so perfect. Therefore, we need to understand how to see the prediction probability, how to numerically quantify it, and how to calibrate our model in order to improve it.

How to find the prediction probability?

Let's now imagine a model for a binary classification:

```
1
2 # Create dataset of classification task with many redundant and few
3 # informative features
4
5 import numpy as np
6 from sklearn import datasets
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.metrics import brier_score_loss, precision_score, recall_score, f1_score
10 np.random.seed(0)
11
12 X, y = datasets.make_classification(n_samples=100000, n_features=20, n_informative=2, n_redundant=
13
14 train_samples = 100
15 X_train = X[:train_samples]
16 X_test = X[train_samples:]
17 y_train = y[:train_samples]
18 y_test = y[train_samples:]
19
20 model_rfc = RandomForestClassifier()
21 model_rfc.fit(X_train, y_train)
22 y_pred_rfc = model_rfc.predict(X_test)
23
24 print("Precision: %0.2f" % precision_score(y_test, y_pred_rfc))
25 print("Recall: %0.2f" % recall_score(y_test, y_pred_rfc))
26 print("F1: %0.2f\n" % f1_score(y_test, y_pred_rfc))
27
28 model_lr = LogisticRegression()
29 model_lr.fit(X_train, y_train)
30 y_pred_lr = model_lr.predict(X_test)
31
32 print("Precision: %0.2f" % precision_score(y_test, y_pred_lr))
33 print("Recall: %0.2f" % recall_score(y_test, y_pred_lr))
34 print("F1: %0.2f" % f1_score(y_test, y_pred_lr))
35
36 # Precision: 0.85
37 # Recall: 0.92
38 # F1: 0.89
39
40 # Precision: 0.84
41 # Recall: 0.86
42 # F1: 0.85
```

The random forest classifier (RFC) got an F1 score of 0.89, which is not bad. The logistic regression performed just a bit worse than RF with a score of 0.85. But how well calibrated are they?

To do so, we can use the out-of-the-box function `predict_proba()`.

Note: not all ML algorithms have this function. For others we need the `decision_function()` (see [Comparison of Calibration of Classifiers](#) for details).

In the case of `RandomForestClassifier`, we can use `predict_proba()` and obtain the predicted class probability for each input sample.

The result would be something like this:

```
model.predict_proba(X_test)[:, 1]  
# array([0.18, 0.84, 0.91, ..., 0.38, 0.54, 0.84])
```

This is not enough if we want to see the relationship between probability and positive prediction. For that, we need to use the function `calibration_curve`.

This function takes in as input the true label used in the test and the result of `predict_proba()`, and the number of bins we want to use to summarise the result (default 10).

`Calibration_curve` returns two results: the fraction of positives, that is, the proportion of positive class samples, per bin, and the mean predicted value, thus the average value of all predicted probability in each bin.

If we plot these two results, we will obtain a **calibration plot**.

Calibration plot

Let's see an example of a calibration plot:

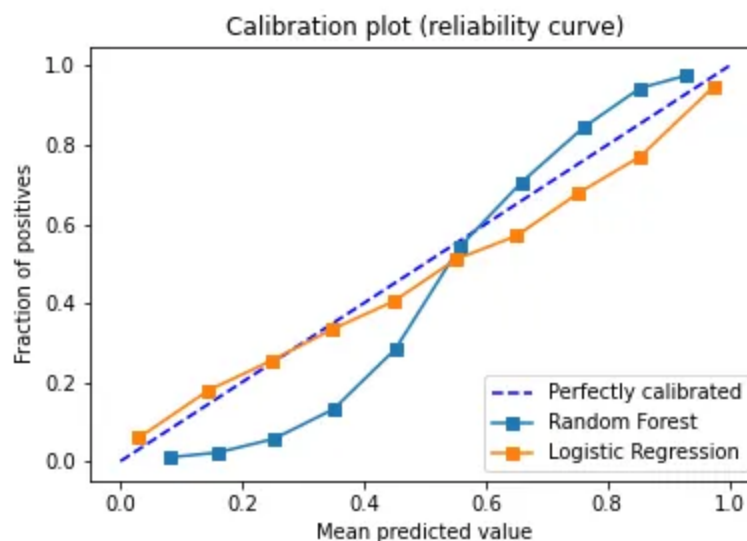
```

1  import matplotlib.pyplot as plt
2  from sklearn.calibration import CalibratedClassifierCV, calibration_curve
3
4  fig = plt.figure()
5  ax = fig.add_subplot()
6  ax.plot([0, 1], [0, 1], "b--", label="Perfectly calibrated")
7  ax.set_ylabel("Fraction of positives")
8  ax.set_xlabel("Mean predicted value")
9  ax.set_title('Calibration plot (reliability curve)')
10
11 prob_pos_rfc = model_rfc.predict_proba(X_test)[: , 1]
12 fraction_of_positives_rfc, mean_predicted_value_rfc = calibration_curve(y_test, prob_pos_rfc, n_bin
13 ax.plot(mean_predicted_value_rfc, fraction_of_positives_rfc, "s-", label="%s" % ('Random Forest'
14
15 prob_pos_lr = model_lr.predict_proba(X_test)[: , 1]
16 fraction_of_positives_lr, mean_predicted_value_lr = calibration_curve(y_test, prob_pos_lr, n_bin
17 ax.plot(mean_predicted_value_lr, fraction_of_positives_lr, "s-", label="%s" % ('Logistic Regress
18
19 ax.legend(loc="lower right")
20 plt.show()
21 fig.savefig('02_rf_lr.png', tight=True, quality=100)

```

probability_calibration_2.py hosted with ❤ by GitHub

[view raw](#)



The plot above is commonly referred to as the **calibration plot or reliability diagram (or curve)**. In our example, it contains calibration curves for the random forest and logistic regression classifiers, as well as the diagonal dotted line that represents the perfectly calibrated model.

On the x-axis, is the mean predicted value for each bin, and on the y-axis is the fraction of positives.

The random forest has its typical sigmoid shape; the RFC is overconfident in its prediction, pushing the probability toward zero and one. Logistic regression on the other hand usually returns a well calibrated prediction. Therefore, it is usually plotted together with other methods even if it is not directly in use.

Note: When a model goes **below the diagonal**, the model is *over-forecasting*. **Above the diagonal**, the model is *under-forecasting*.

Now, we have a visual representation of the calibration curves, but how can we decide which model is better calibrated? For this we need the Brier Score!

Brier score

The Brier score loss ([wiki](#)) is a single specific number that we can refer to, to assess how good our model is. It is a loss score so the smaller it is, the more calibrated the model.

It is computed with the following method, `brier_score_loss(y_test, prob_pos)`. Let's see an example:

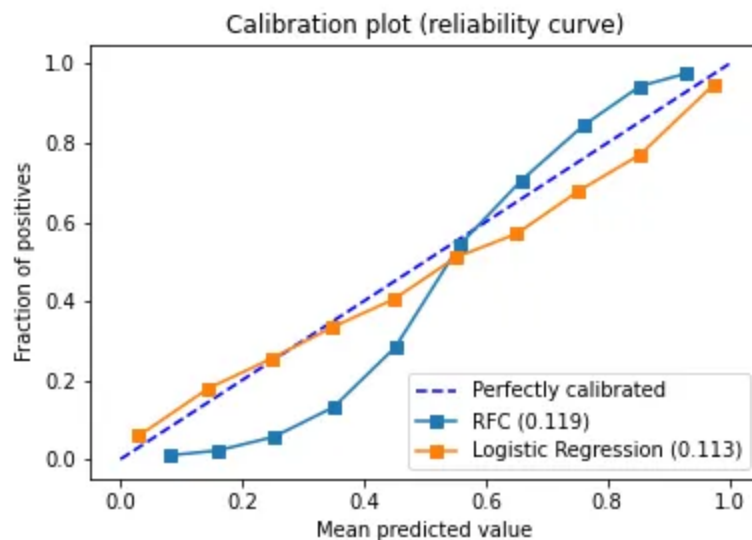
```

1  fig = plt.figure()
2  ax = fig.add_subplot()
3  ax.plot([0, 1], [0, 1], "b--", label="Perfectly calibrated")
4  ax.set_ylabel("Fraction of positives")
5  ax.set_xlabel("Mean predicted value")
6  ax.set_title('Calibration plot (reliability curve)')
7
8  model_rfc_score = brier_score_loss(y_test, y_pred_rfc, pos_label=y.max())
9  fraction_of_positives, mean_predicted_value = calibration_curve(y_test, prob_pos_rfc, n_bins=10)
10
11 ax.plot(mean_predicted_value, fraction_of_positives, "s-", label="%s (%0.3f)" % ('RFC', model_rfc_score))
12 print("Brier score RFC: %0.2f" % brier_score_loss(y_test, y_pred_rfc))
13
14
15 model_lr_score = brier_score_loss(y_test, prob_pos_lr, pos_label=y.max())
16 fraction_of_positives, mean_predicted_value = calibration_curve(y_test, prob_pos_lr, n_bins=10)
17
18 ax.plot(mean_predicted_value, fraction_of_positives, "s-", label="%s (%0.3f)" % ('Logistic Regression', model_lr_score))
19 print("Brier score LR: %0.2f" % brier_score_loss(y_test, y_pred_lr))
20
21 ax.legend(loc="lower right")
22 plt.show()
23 fig.savefig('03_rf_lr.png', tight=True, quality=100)
24 # Brier score RFC: 0.12
25 # Brier score LR: 0.11

```

probability_calibration_3.py hosted with ❤ by GitHub

[view raw](#)



In our example, it returns a value of 0.12 for random forest and 0.11 for logistic regression.

Now we have a system to evaluate the calibration curve. Let's see how we can improve it. For this we need to **calibrate** the model.

Probability calibration

The probability calibration of a model is a re-scaling of the model, it can be done using the scikit function CalibratedClassifierCV

There are two arguments of the function we have to consider: the **methods** and the **validations**:

1. Two methods of calibration:

- *Sigmoid scaling or Platt's method*. This is suitable for models with a sigmoid curve (like RFC) and it works well with small datasets.
- *Isotonic regression* — a non-parametric approach. It is a more powerful calibration method but it tends to overfit and is not advised for small datasets.

2. Two possible validations:

- **Before**, using classical cross validation:

```
1 X_train, X_test, y_train, y_test = train_test_split(...)\n2 model = ...\n3 calibrator = CalibratedClassifierCV(model, cv=3)\n4 calibrator.fit(trainX, trainy)\n5 predictions = calibrator.predict(testX)
```

probability_calibration_before.py hosted with ❤ by GitHub

[view raw](#)

- **After** the model is trained using “prefit” model. This is usually more suitable for larger datasets:

```
1 X_train, X_train2, y_train, X_test, X_test2, y_test = ...\n2 model = ...\n3 model.fit(trainX, trainy)\n4 calibrator = CalibratedClassifierCV(model, cv='prefit')\n5 calibrator.fit(valX, valy)\n6 predictions = calibrator.predict(testX)
```

probability_calibration_after.py hosted with ❤ by GitHub

[view raw](#)

Let's see an example:

Cross Validation

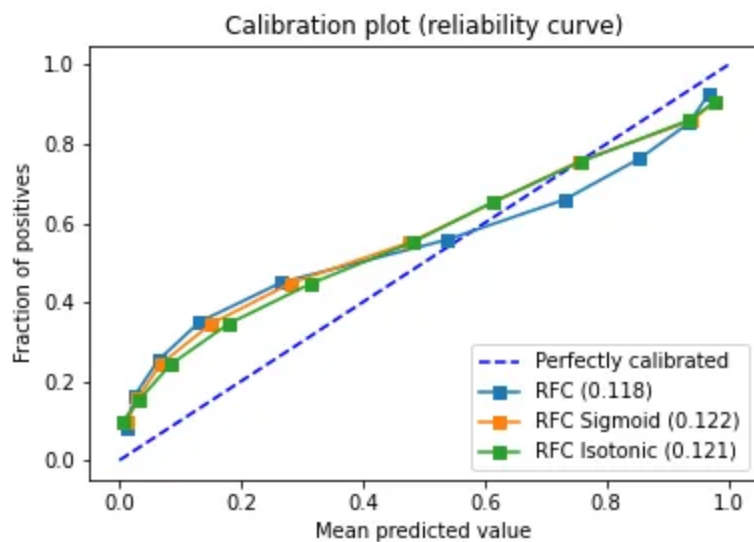
```

1  def get_calibration_curve_values(model, X_train, y_train, X_test, y_test):
2      model.fit(X_train, y_train)
3      y_pred = model.predict(X_test)
4      prob_pos = model.predict_proba(X_test)[: , 1]
5      model_score = brier_score_loss(y_test, y_pred, pos_label=y.max())
6
7      print("F1: %0.2f" % f1_score(y_test, y_pred))
8      fraction_of_positives, mean_predicted_value = calibration_curve(y_test, prob_pos, n_bins=10)
9
10     return(fraction_of_positives, mean_predicted_value, model_score)
11
12  fig = plt.figure()
13  ax = fig.add_subplot()
14  ax.plot([0, 1], [0, 1], "b--", label="Perfectly calibrated")
15  ax.set_ylabel("Fraction of positives")
16  ax.set_xlabel("Mean predicted value")
17  ax.set_title('Calibration plot (reliability curve)')
18
19  model_rfc = RandomForestClassifier()
20  mpv, fp, score = get_calibration_curve_values(model_rfc, X_train, y_train, X_test, y_test)
21  ax.plot(mpv, fp, "s-", label="%s (%1.3f)" % ('RFC', score))
22
23  calibrator_sigmoid = CalibratedClassifierCV(RandomForestClassifier(), cv=5, method='sigmoid')
24  mpv, fp, score = get_calibration_curve_values(calibrator_sigmoid, X_train, y_train, X_test, y_te
25  ax.plot(mpv, fp, "s-", label="%s (%1.3f)" % ('RFC Sigmoid', score))
26
27  calibrator_isotonic = CalibratedClassifierCV(RandomForestClassifier(), cv=5, method='isotonic')
28  mpv, fp, score = get_calibration_curve_values(calibrator_sigmoid, X_train, y_train, X_test, y_te
29  ax.plot(mpv, fp, "s-", label="%s (%1.3f)" % ('RFC Isotonic', score))
30
31  ax.legend(loc="lower right")
32  plt.show()
33  fig.savefig('04_before.png', tight=True, quality=100)
34  # F1: 0.89
35  # F1: 0.88
36  # F1: 0.88

```

probability_calibration_4.py hosted with ❤ by GitHub

[view raw](#)



The F1 score of this method has not changed but the Brier score has indeed improved.

Prefit

Open in app ↗



Search Medium



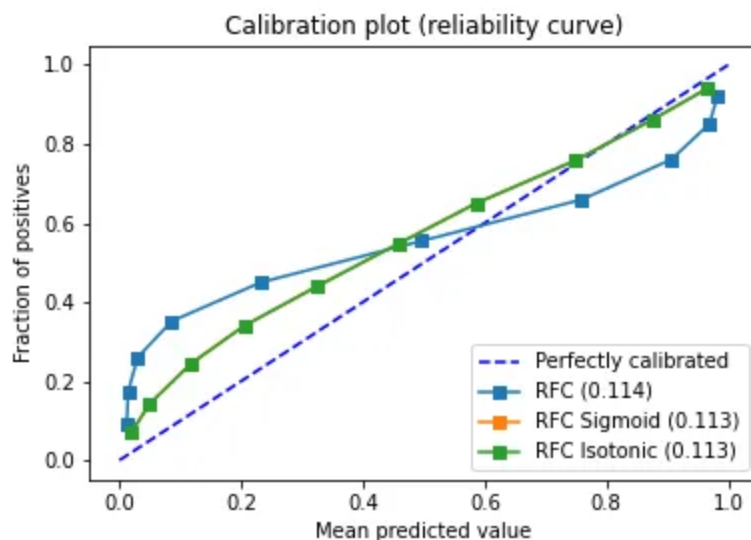
Write



```
2
3  fig = plt.figure()
4  ax = fig.add_subplot()
5  ax.plot([0, 1], [0, 1], "b--", label="Perfectly calibrated")
6  ax.set_ylabel("Fraction of positives")
7  ax.set_xlabel("Mean predicted value")
8  ax.set_title('Calibration plot (reliability curve)')
9
10 X_train, X_train2, y_train, y_train2 = train_test_split(X_train, y_train, test_size=0.5)
11
12 model_rfc = RandomForestClassifier()
13 mpv, fp, score = get_calibration_curve_values(model_rfc, X_train, y_train, X_test, y_test)
14 ax.plot(mpv, fp, "s-", label="%s (%1.3f)" % ('RFC', score))
15
16 calibrator_sigmoid = CalibratedClassifierCV(model_rfc, cv='prefit', method='sigmoid')
17 mpv, fp, score = get_calibration_curve_values(calibrator_sigmoid, X_train2, y_train2, X_test, y_test)
18 ax.plot(mpv, fp, "s-", label="%s (%1.3f)" % ('RFC Sigmoid', score))
19
20 calibrator_isotonic = CalibratedClassifierCV(model_rfc, cv='prefit', method='isotonic')
21 mpv, fp, score = get_calibration_curve_values(calibrator_isotonic, X_train2, y_train2, X_test, y_test)
22 ax.plot(mpv, fp, "s-", label="%s (%1.3f)" % ('RFC Isotonic', score))
23
24 ax.legend(loc="lower right")
25 plt.show()
26 fig.savefig('05_after.png', tight=True, quality=100)
27 # F1: 0.89
28 # F1: 0.89
29 # F1: 0.89
```

probability_calibration_5.py hosted with ❤ by GitHub

[view raw](#)



With this method not only the Brier score is improved but also the F1.

Note: In this particular case, the F1 has improved, but it is not always the case! most of the time accuracy and other metrics would decrease in performance. Be aware of this trade-off.

Besides, sigmoid and isotonic calibration overlap, this is rarely the case as well.

Sources

In order to understand the concept of the Calibrate Probability I had to look at several source, here a list of the ones I rely the most.

1. [How to Calibrate Probabilities for Imbalanced Classification](#)
2. [Model Calibration — is your model ready for the real world? — Inbar Naor — PyCon Israel 2018](#)
3. [Safe Handling Instructions for Probabilistic Classification | SciPy 2019 | Gordon Chen](#)
4. [Applied ML 2020-10 — Calibration, Imbalanced data](#)

5. [plot_calibration_curve](#)

6. [Classifier calibration](#)

7. [Calibration of Models — Nupur Agarwal](#)

I wrote this article as a way to organise and test my knowledge on this topic. you find errors, or you find it useful, or have suggestions please leave a comment.

For the code in jupyter click [here](#).

Originally published at <https://mattiacinelli.com> on May 27, 2020.

[Python](#)[Probability](#)[Calibration](#)[Machine Learning](#)[Data Science](#)

More from the list: "ML"

Curated by Himanshu Birla



Kyosuke... in Towards Dat...

Probability Calibration for Imbalanced Dataset



Jason Yo... in Towards Dat...

Why Calibrators? Part 1 of the Series on Probabilit...



Jaideep Ray

**Probability c
why it matter**



★ · 8 min read · Oct 20, 2019

7 min read · Oct 4, 2020

3 min read · Ju

[View list](#)



Written by Mattia Cinelli

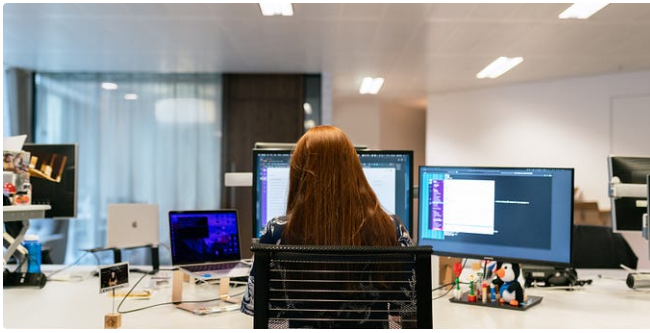
321 Followers · Writer for Analytics Vidhya

Data Scientist

Following



More from Mattia Cinelli and Analytics Vidhya



Mattia Cinelli in Towards Data Science

SOLID Coding in Python

An overview of the 5 SOLID principles of coding that would improve your Python...

🌟 · 8 min read · Jun 29, 2021



1.3K



12



...



259



1



...



Kia Eisinga in Analytics Vidhya

How to create a Python library

Ever wanted to create a Python library, albeit for your team at work or for some open...

7 min read · Jan 27, 2020



2K



24



...



270



2



...

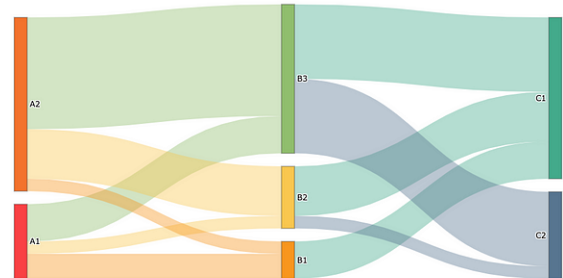


Sajal Rastogi in Analytics Vidhya

Wifi -Hacking using PyWifi 🗝

.

4 min read · Feb 6, 2021



Mattia Cinelli in Analytics Vidhya

How to do a Sankey Plot in Python

If you have been more than five seconds on r/datsbeautiful/, you will have probably...

🌟 · 4 min read · May 3, 2020



270



2

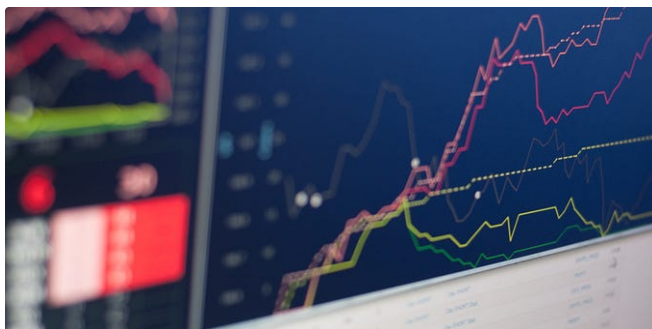


...

See all from Mattia Cinelli

See all from Analytics Vidhya

Recommended from Medium



Chandra Prakash Bathula

Machine Learning Concept 68: Platt's Scaling

Platt's Scaling:

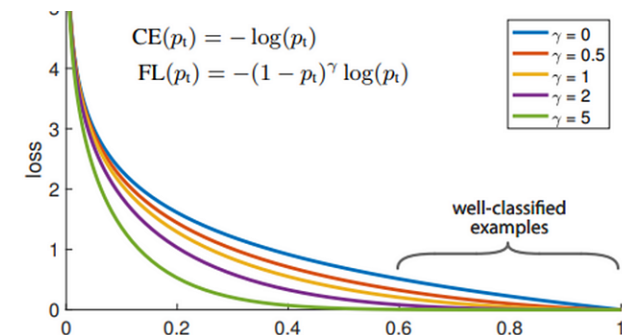
3 min read · Apr 13



5



...



Uman Niyaz in Data Science @ Ecom Express

Focal loss for handling the issue of class imbalance

Text classification is widely used in various industries to tackle business challenges by...

9 min read · Jun 12

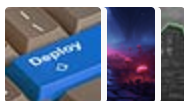


53



...

Lists



Predictive Modeling w/ Python

20 stories · 473 saves



Practical Guides to Machine Learning

10 stories · 544 saves



Coding & Development

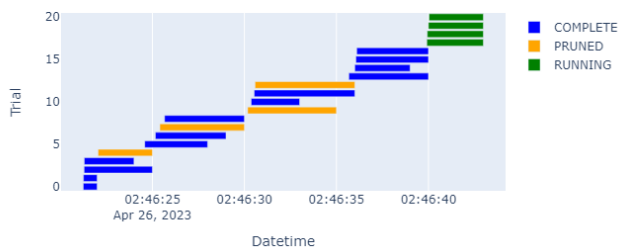
11 stories · 207 saves



New_Reading_List

174 stories · 143 saves

Timeline Plot



Hiroki Takizawa in Optuna

Introducing Timeline Plot: a new feature of Optuna v3.2

The newest version of Optuna, Optuna v3.2, an open source software tool for black box...

4 min read · Jun 23



14

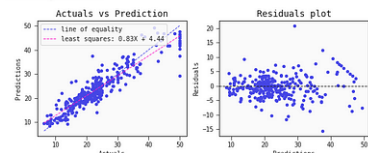


...

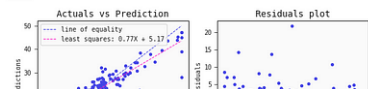


...

Training Metrics



Test



Casper Skern Wilstrup

Symbolic Regression: a Simple and Friendly Introduction

Symbolic Regression is like a treasure hunt for the perfect mathematical equation to...

3 min read · May 5



18



1



...



94



...



Juan Broglio

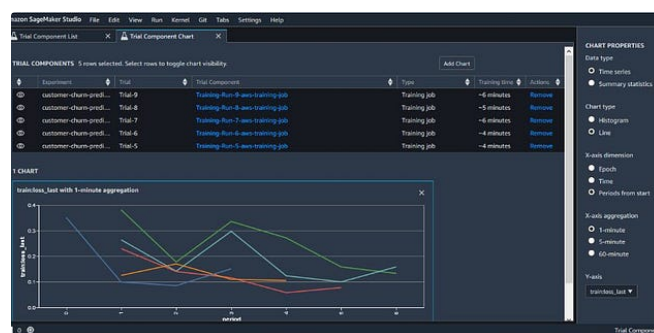
Gamma Regression vs Linear Regression (in Python)

General Linear Models and Gamma Regression

4 min read · Aug 9



...



Yugank Aman

Top MLOps Tools to Manage Machine Learning Lifecycle

Businesses continue transforming their operations to increase productivity and...

10 min read · May 30



94



...

See more recommendations