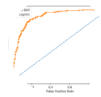


Never miss a tutorial:



Machine Learning Mastery
Making Developers Awesome at Machine Learning

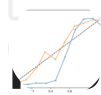
Picked for you:



How to Use ROC Curves and Precision-Recall Curves for Classification in Python

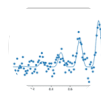
[Click to Take the FREE Probability Crash-Course](#)

Search...



How and When to Use a Calibrated Classification Model with scikit-learn

Probabilistic Model Selection with AIC, BIC, and MDL

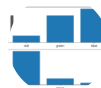


How to Implement Bayesian Optimization from Scratch in Python

Tweet

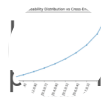
Share

Share



How to Calculate the KL Divergence for Machine Learning

Model selection is the problem of choosing one from among a set of candidate models.



A Gentle Introduction to Cross-Entropy for Machine Learning

Common to choose a model that performs the best on a hold-out test dataset or to estimate model performance using a resampling technique, such as [k-fold cross-validation](#).

An alternative approach to model selection involves using probabilistic statistical measures that attempt to quantify both the model performance on the training dataset and the complexity of the model.

Loving the tutorials?

Examples include the Akaike and Bayesian Information Criterion and the Minimum Description Length.

The [Probability for Machine Learning](#) EBook is

The benefit of these information criterion statistics is that they do not require a hold-out test set, although a limitation is that they do not take the uncertainty of the models into account and may end-up selecting

>> SEE WHAT'S INSIDE

In this post, you will discover probabilistic statistics for machine learning model selection.

After reading this post, you will know:

- Model selection is the challenge of choosing one among a set of candidate models.
- Akaike and Bayesian Information Criterion are two ways of scoring a model based on its log-likelihood and complexity.
- Minimum Description Length provides another scoring method from information theory that can be shown to be equivalent to BIC.

Kick-start your project with my new book [Probability for Machine Learning](#), including *step-by-step tutorials* and the *Python source code* files for all examples.

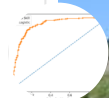
Let's get started.

[Start Machine Learning](#)

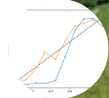
Never miss a tutorial:



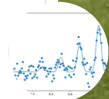
Picked for you:



How to Use ROC Curves and Precision-Recall Curves for Classification in Python



How and When to Use a Calibrated Classification Model with scikit-learn



How to Implement Bayesian Optimization from Scratch in Python



How to Calculate the KL Divergence for Machine Learning

Probabilistic Model Selection

Photo by Guilhem Ver

A Gentle Introduction to Cross-Entropy for Machine Learning

Overview

This tutorial is divided into five parts; they are:

Loving the Tutorials?

1. The Challenge of Model Selection
2. The Probability for Machine Learning EBook is where you'll find the **Really Good** stuff.
3. Akaike Information Criterion
4. Baye >> SEE WHAT'S INSIDE
5. Minimum Description Length

The Challenge of Model Selection

Model selection is the process of fitting multiple models on a given dataset and choosing one over all others.



Model selection: estimating the performance of different models in order to choose the best one.

— Page 222, *The Elements of Statistical Learning*, 2016.

This may apply in unsupervised learning, e.g. choosing a clustering model, or supervised learning, e.g. choosing a predictive model for a regression or classification task. It may also be a sub-task of modeling, such as feature selection for a given model.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

There are many common approaches that may be used for model selection. For example, in the case of supervised learning, the three most common approaches are:

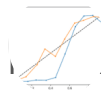


Training, Validation, and Test Datasets.

- Resampling Methods.

Picked for you: Statistics.

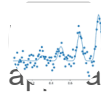
The simplest and most reliable method of model selection involves fitting candidate models on a training set, then on the validation dataset, and selecting a model that performs the best on the test dataset according to a chosen metric, such as accuracy or error. A problem with this approach is that it requires a lot of data.



How and When to Use a Calibrated

Classification Model with scikit-learn

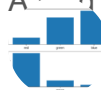
Sampling techniques attempt to achieve the same results as the full dataset approach, although using a small dataset. An example is k-fold cross-validation, where the data is split into k parts, and a model is fit and evaluated on each part. This approach is selected with the best average score across all folds. A problem with this approach is that only model performance is assessed on the training set.



How to Implement Bayesian Optimization

from Scratch in Python

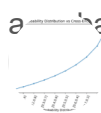
A third approach to model selection attempts to convert the performance of the model into a score, then select the model with the highest score.



How to Calculate the KL Divergence for

Machine Learning

We can refer to this approach as statistical or probabilistic framework.



A Gentle Introduction to Cross-Entropy for

Machine Learning

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Want to Learn Probability for Machine Learning Loving the Tutorials?

Take my free 7-day email crash course now (with sample code). The [Probability for Machine Learning](#) EBook is where you'll find the **Really Good** stuff. Click to sign-up and also get a free PDF Ebook version of the course.

>> SEE WHAT'S INSIDE

Download Your FREE Mini-Course

Probabilistic Model Selection

Probabilistic model selection (or “information criteria”) provides an analytical technique for scoring and choosing among candidate models.

Models are scored both on their performance on the training dataset and based on the complexity of the model.

- **Model Performance.** How well a candidate model has performed on the training dataset.
- **Model Complexity.** How complicated the trained candidate model is after training.

Model performance may be evaluated using a probabilistic framework of maximum likelihood estimation. Model

Start Machine Learning

degrees of freedom or parameters in the model.

Never miss a tutorial:



‘information criteria’ have been proposed that attempt to correct for the bias of maximum likelihood by the addition of a penalty term to compensate for the over-fitting of more complex models.

Picked for you:



Page 36, [Pattern Recognition and Machine Learning](#), 2006.

[Recall Curves for Classification in Python](#)

A benefit of probabilistic model selection methods is that a test dataset is not required, meaning that all of the data can be used to fit the model, and the final model that will be used for prediction in the



[How and When to Use a Calibrated Classification Model with scikit-learn](#)

A limitation of probabilistic model selection methods is that they are calculated across a range of different types of model.

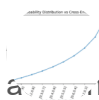


[How to Implement Bayesian Optimization from Scratch in Python](#)



It should be noted that the AIC statistic is calculated across a range of different types of model. (as opposed to comparisons of models)

— Page 493, [Applied Predictive Modeling](#), 2013.



A Gentle Introduction to Cross-Entropy for Machine Learning

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE



Such criteria do not take account of the uncertainty in the model parameters, however, and in practice they tend to favour overly simple models.

Loving the Tutorials?

The [Probability for Machine Learning](#) EBook is

— Page 33, [Pattern Recognition and Machine Learning](#), 2006.

There are [SEE WHAT'S INSIDE](#) es to estimating how well a given model fits a dataset and how complex the model is. And each can be shown to be equivalent or proportional to each other, although each was derived from a different framing or field of study.

They are:

- Akaike Information Criterion (AIC). Derived from frequentist probability.
- Bayesian Information Criterion (BIC). Derived from Bayesian probability.
- Minimum Description Length (MDL). Derived from information theory.

Each statistic can be calculated using the log-likelihood for a model and the data. Log-likelihood comes from Maximum Likelihood Estimation, a technique for finding or optimizing the parameters of a model in response to a training dataset.

In [Maximum Likelihood Estimation](#), we wish to maximize the conditional probability of observing the data (X) given a specific probability distribution and its parameters (θ), stated formally as:

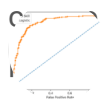
- $P(X; \theta)$

Start Machine Learning

Where X is, in fact, the joint probability distribution of all observations from the problem domain from 1 to n .



The joint probability distribution can be restated as the multiplication of the conditional probability for observing each example given the distribution parameters. Multiplying many small probabilities together



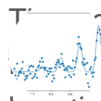
unstable, as ROC Curves can be non-monotonic. Recall Curves for Classification in Python

- $\sum_{i=1}^n \log(P(x_i; \theta))$



How and When to Use a Calibrated Frequentist Model

the frequentist use of log in the likelihood function, it is commonly referred to as a log-likelihood function.

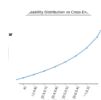


How to Calculate the KL Divergence for Machine Learning

log-likelihood function for Bayesian Optimization from Scratch in Python

How to Calculate the KL Divergence for Machine Learning

Akaike Information Criterion



A Gentle Introduction to Cross-Entropy for Akaike Information Criterion, or AIC for short, Machine Learning

It is named for the developer of the method, [Hirotugu Akaike](#), and may be shown to have a basis in information theory and frequentist-based inference.

Loving the Tutorials?



This is derived from a frequentist framework, and cannot be interpreted as an approximation to the marginal likelihood.

— Page 1 >> SEE WHAT'S INSIDE [Probabilistic Perspective](#), 2012.

The AIC statistic is defined for logistic regression as follows (taken from “[The Elements of Statistical Learning](#)”):

- $AIC = -2/N * LL + 2 * k/N$

Where N is the number of examples in the training dataset, LL is the log-likelihood of the model on the training dataset, and k is the number of parameters in the model.

The score, as defined above, is minimized, e.g. the model with the lowest AIC is selected.



To use AIC for model selection, we simply choose the model giving smallest AIC over the set of models considered.

— Page 231, [The Elements of Statistical Learning](#), 2016.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Compared to the BIC method (below), the AIC statistic penalizes complex models less, meaning that it may put more emphasis on model performance on the training dataset, and, in turn, select more complex model.



Picked for you: We see that the penalty for AIC is less than for BIC. This causes AIC to pick more complex models.



How to Use ROC Curves and Precision-

Recall Curves for Classification in Python
 102, Machine Learning: A Probabilistic Perspective, 2012.

Bayesian Information Criterion



How and When to Use a Calibrated

Classification Model with scikit-learn

Bayesian Information Criterion, or BIC for short

It is named for the field of study from which it was



How to Implement Bayesian Optimization

Appropriate for models fit under the maximum

from Scratch in Python

The BIC statistic is calculated for logistic regression



How to Calculate the KL Divergence for

Machine Learning

$BIC = -2 * LL + \log(N) * k$

Where $\log()$ has the base-e called the natural log



A Gentle Introduction to Cross-Entropy for

of examples in the training dataset, and k

Machine Learning

The score as defined above is minimized, e.g. the model with the lowest BIC is selected.

The quantity calculated is different from AIC, although can be shown to be proportional to the AIC.

Unlike the AIC, the BIC penalizes the model more for its complexity, meaning that more complex

The Probability for Machine Learning EBook is
 models will have a worse (larger) score and will, in turn, be less likely to be selected.
 where you'll find the **Really Good** stuff.

Loving the Tutorials? *Nk >> SEE WHAT'S INSIDE [...] this penalizes model complexity more heavily.*

— Page 217, [Pattern Recognition and Machine Learning](#), 2006.

Importantly, the derivation of BIC under the Bayesian probability framework means that if a selection of candidate models includes a true model for the dataset, then the probability that BIC will select the true model increases with the size of the training dataset. This cannot be said for the AIC score.

... given a family of models, including the true model, the probability that BIC will select the correct model approaches one as the sample size $N \rightarrow \infty$.

— Page 235, [The Elements of Statistical Learning](#), 2016.

A downside of BIC is that for smaller, less representative training datasets, it is more likely to choose models that are too simple.

Minimum Description Length

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.**

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

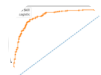
Start Machine Learning

The [Minimum Description Length](#), or MDL for short, is a method for scoring and selecting a model. **Never miss a tutorial:**

It is named for the field of study from which it was derived, namely information theory.



Information theory is concerned with the representation and transmission of information on a noisy channel, and as such, measures quantities like entropy, which is the average number of bits required to represent an event from a random variable or probability distribution. **Picked for you:**



[How to Use ROC Curves and Precision-](#)

[Recall Curves for Classification in Python](#)

When you want to transmit both the predictions (or more precisely, their probability distributions) and the model used to generate them. Both the predicted target variable and the model can be described in terms of the number of bits required to transmit them on a



[How and When to Use a Calibrated](#)

[Classification Model with scikit-learn](#)

The Minimum Description Length is the minimum number of bits required to represent the data and the



[How to Implement Bayesian Optimization](#)

[from Scratch in Python](#)

The Minimum Description Length (MDL) principle minimizes the sum of these two descriptions

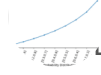


[How to Calculate the KL Divergence for](#)

[Machine Learning](#)

[Page 173, Machine Learning, 1997.](#)

The MDL statistic is calculated as follows (taken from [A Gentle Introduction to Cross-Entropy for](#)



[Machine Learning](#)

$MDL = L(h) + L(D | h)$

Where h is the model, D is the predictions made by the model, $L(h)$ is the number of bits required to represent the model, and $L(D | h)$ is the number of bits required to represent the predictions from the model on the training dataset. **Loving the Tutorials?**

The [Probability for Machine Learning](#) EBook is

The score as defined above is minimized, e.g. the model with the lowest MDL is selected.

The number of bits required to encode D (denoted $L(D | h)$) and the number of bits required to encode h can be calculated as the negative log-likelihood; for example (taken from "[The Elements of Statistical Learning](#)"):

- $MDL = -\log(P(\theta)) - \log(P(y | X, \theta))$

Or the negative log-likelihood of the model parameters (θ) and the negative log-likelihood of the target values (y) given the input values (X) and the model parameters (θ).

This desire to minimize the encoding of the model and its predictions is related to the notion of [Occam's Razor](#) that seeks the simplest (least complex) explanation: in this context, the least complex model that predicts the target variable.



The MDL principle takes the stance that the best theory for a body of data is one that minimizes the size of the theory plus the amount of information necessary to specify the exceptions relative to the theory ...

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

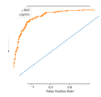
Page 198, [Data Mining: Practical Machine Learning Tools and Techniques](#), 4th edition, 2016.

Never miss a tutorial:

The MDL calculation is very similar to BIC and can be shown to be equivalent in some situations.

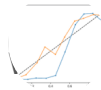


Picked for you: Hence the BIC criterion, derived as approximation to log-posterior probability, can also be viewed as a device for (approximate) model choice by minimum description length.



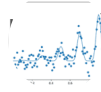
How to Use ROC Curves and Precision-Recall Curves for Classification
 236 The Elements of Statistical Learning, 2016.

Worked Example for Linear Regression



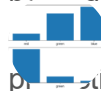
How and When to Use a Calibrated Classification Model with scikit-learn
 I make the calculation of AIC and BIC conc

In this section, we will use a test problem and fit a model. The AIC and BIC metrics.



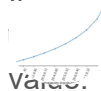
How to Implement Bayesian Optimization from Scratch in Python

Importantly, the specific functional form of AIC and BIC have been derived, making the example relatively straightforward. It is important to either find an appropriate model or look into deriving the calculation.



How to Calculate the KL Divergence for Machine Learning

In this example, we will use a test regression problem. The problem will have two input variables and one output variable.



A Gentle Introduction to Cross-Entropy for Machine Learning

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

```
1 ...
2 # generate dataset
3 X, y = make_regression(n_samples=100, n_features=2, noise=0.1)
4 # define and fit the model on all data
```

We will fit a LinearRegression() model on the entire dataset directly.

```
1 ... >> SEE WHAT'S INSIDE
2 # define and fit the model on all data
3 model = LinearRegression()
4 model.fit(X, y)
```

Once fit, we can report the number of parameters in the model, which, given the definition of the problem, we would expect to be three (two coefficients and one intercept).

```
1 ...
2 # number of parameters
3 num_params = len(model.coef_) + 1
4 print('Number of parameters: %d' % (num_params))
```

The likelihood function for a linear regression model can be shown to be identical to the least squares function; therefore, we can estimate the maximum likelihood of the model via the mean squared error metric.

First, the model can be used to estimate an outcome for each example in the training dataset, then the `mean_squared_error()` scikit-learn function can be used to calculate the mean squared error for the model.

Start Machine Learning


```

1 ...
2 # predict the training set
3 yhat = model.predict(X)
4 # calculate the error
5 mse = mean_squared_error(y, yhat)
6 print('MSE: %.3f' % mse)

```

Picking this all together, the complete example of defining the dataset, fitting the model, and reporting the number of parameters and maximum likelihood estimate of the model is listed below.

How to Use ROC Curves and Precision-Recall Curves for Classification in Python

```

1 # generate dataset and fit a linear regression model
2 from sklearn.datasets import make_regression
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5 # generate dataset
6 X, y = make_regression(n_samples=100, n_features=2, noise=0.1)
7 # define and fit the model on all data
8 model = LinearRegression()
9 model.fit(X, y)
10 # number of parameters
11 num_params = len(model.coef_) + 1
12 print('Number of parameters: %d' % (num_params))
13 # predict the training set
14 yhat = model.predict(X)
15 # calculate the error
16 mse = mean_squared_error(y, yhat)
17 print('MSE: %.3f' % mse)

```

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Note: Your results may vary given the stochastic nature of the data and numerical precision. Consider running the example multiple times to get a better idea of the results.

Running the example first reports the number of parameters in the model as 3, as we expected, then reports the MSE as about 0.01.

Loving the Tutorials?

```

1 Number of parameters: 3
2 MSE: 0.010

```

where you'll find the **Really Good** stuff.

Next, we can adapt the example to calculate the AIC for the model.

>> SEE WHAT'S INSIDE

Skipping the derivation, the AIC calculation for an ordinary least squares linear regression model can be calculated as follows (taken from “A New Look At The Statistical Identification Model”, 1974.):

- $AIC = n * LL + 2 * k$

Where n is the number of examples in the training dataset, LL is the log-likelihood for the model using the natural logarithm (e.g. the log of the MSE), and k is the number of parameters in the model.

The `calculate_aic()` function below implements this, taking n , the raw mean squared error (`mse`), and k as arguments.

```

1 # calculate aic for regression
2 def calculate_aic(n, mse, num_params):
3     aic = n * log(mse) + 2 * num_params
4     return aic

```

The example can then be updated to make use of this new function and calculate the AIC for the model.

The complete example is listed below.

Start Machine Learning

```

1 # calculate akaike information criterion for a linear regression model
2 from math import log
3 from sklearn.datasets import make_regression
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6
7 # calculate aic for regression
8 def calculate_aic(n, mse, num_params):
9     aic = n * log(mse) + 2 * num_params
10    return aic
11
12 # generate dataset
13 X, y = make_regression(n_samples=100, n_features=2, noise=0.1)
14 # define and fit the model on all data
15 model = LinearRegression()
16 model.fit(X, y)
17 # number of parameters
18 num_params = len(model.coef_) + 1
19 print('Number of parameters: %d' % (num_params))
20 # predict the training set
21 yhat = model.predict(X)
22 # calculate the error
23 mse = mean_squared_error(y, yhat)
24 print('MSE: %.3f' % mse)
25 # calculate the aic
26 aic = calculate_aic(len(y), mse, num_params)
27 print('AIC: %.3f' % aic)

```

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Your results may vary given the stochastic nature of the data. A Gentle Introduction to Cross-Entropy for Machine Learning

In this case, the AIC is reported to be a value of about -451.616. This value can be minimized in order to choose better models.

Loving the Tutorials?

The Probability for Machine Learning EBook is

```

1 Number of parameters: 3
2 MSE: 0.010
3 AIC: -451.616

```

>> SEE WHAT'S INSIDE

We can also update the example with the calculation of BIC instead of AIC.

Skipping the derivation, the BIC calculation for an ordinary least squares linear regression model can be calculated as follows (taken from here):

- $BIC = n * LL + k * \log(n)$

Where n is the number of examples in the training dataset, LL is the log-likelihood for the model using the natural logarithm (e.g. log of the mean squared error), and k is the number of parameters in the model, and $\log()$ is the natural logarithm.

The `calculate_bic()` function below implements this, taking n , the raw mean squared error (`mse`), and k as arguments.

```

1 # calculate bic for regression
2 def calculate_bic(n, mse, num_params):
3     bic = n * log(mse) + num_params * log(n)
4     return bic

```

The example can then be updated to make use of

Start Machine Learning

del.

The complete example is listed below.

Never miss a tutorial:

```
1 # calculate bayesian information criterion for a linear regression model
2 from math import log
3 from sklearn.datasets import make_regression
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6
7 # calculate bic for regression
8 def calculate_bic(n, mse, num_params):
9     bic = n * log(mse) + num_params * log(n)
10    return bic
11
12 # generate dataset
13 X, y = make_regression(n_samples=100, n_features=2, noise=0.1)
14 # define and fit the model on all data
15 model = LinearRegression()
16 model.fit(X, y)
17 # number of parameters
18 num_params = len(model.coef_) + 1
19 print('Number of parameters: %d' % num_params)
20 # predict the training set
21 yhat = model.predict(X)
22 # calculate the error
23 mse = mean_squared_error(y, yhat)
24 print('MSE: %.3f' % mse)
25 # calculate the bic
26 bic = calculate_bic(len(y), mse, num_params)
27 print('BIC: %.3f' % bic)
```

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Running the example reports the number of parameters, the MSE, and the BIC. A Gentle Introduction to Cross-Entropy for Machine Learning

Now Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

Loving the Tutorials?

In this case, the BIC is reported to be a value of about -450.020, which is very close to the AIC value of -451.616. Again, this value can be minimized in order to choose better models. The Probability for Machine Learning EBook is where you'll find the **Really Good** stuff.

```
1 Number of parameters: 3
2 MSE: 0.010
3 BIC: -450.020
```

Further Reading

This section provides more resources on the topic if you are looking to go deeper.

Books

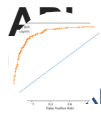
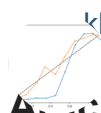
- Chapter 7 Model Assessment and Selection, [The Elements of Statistical Learning](#), 2016.
- Section 1.3 Model Selection, [Pattern Recognition and Machine Learning](#), 2006.
- Section 4.4.1 Model comparison and BIC, [Pattern Recognition and Machine Learning](#), 2006.
- Section 6.6 Minimum Description Length Principle, [Machine Learning](#), 1997.
- Section 5.3.2.4 BIC approximation to log marginal likelihood, [Machine Learning: A Probabilistic Perspective](#), 2012.
- [Applied Predictive Modeling](#), 2013.
- Section 28.3 Minimum description length (MDL), [Information Theory, Inference and Learning Algorithms](#), 2003.

Start Machine Learning

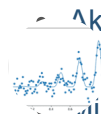
• Section 5.10 The MDL Principle, [Data Mining: Practical Machine Learning Tools and Techniques](#), 4th edition, 2016.




• [A New Look At The Statistical Identification Model](#), 1974.

 How to Use ROC Curves and Precision-Recall Curves for Classification in Python
[learn.datasets.make_regression](#) API.
 • [sklearn.linear_model.LinearRegression](#) API.
[sklearn.metrics.mean_squared_error](#) API.
 How and When to Use a Calibrated Classification Model with scikit-learn

Articles

 [Akaike information criterion, Wikipedia](#)
[How to Implement Bayesian Optimization from Scratch in Python](#)
[Minimum description length, Wikipedia](#)

 [How to Calculate the KL Divergence for Machine Learning](#)

In this post, you discovered probabilistic statistics

 [A Gentle Introduction to Cross-Entropy for Machine Learning](#)

- Model selection is the challenge of choosing one among a set of candidate models.
 - Akaike and Bayesian Information Criterion are two ways of scoring a model based on its log-likelihood and complexity.
 - Minimum Description Length provides another scoring method from information theory that can be
- The [Probability for Machine Learning](#) Ebook is where you'll find the **Really Good** stuff.

Do you have any questions?

Ask your >> SEE WHAT'S INSIDE s below and I will do my best to answer.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Get a Handle on Probability for Machine Learning!

Develop Your Understanding of Probability

...with just a few lines of python code

Discover how in my new Ebook:
[Probability for Machine Learning](#)

It provides **self-study tutorials** and **end-to-end projects** on:

Bayes Theorem, Bayesian Optimization, Distributions, Maximum Likelihood, Cross-Entropy, Calibrating Models and much more...

Finally Harness Uncertainty in Your Projects

Skip the Academic

Start Machine Learning

Never miss a tutorial:

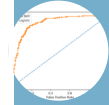
Probability for
Machine Learning



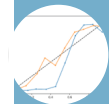
Discover how to harness
Uncertainty With Python

Picked for you:

Jason Brownlee



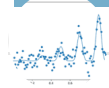
How to Use ROC Curves and Precision-Recall Curves for Classification in Python



How and When to Use a Calibrated Classification Model with scikit-learn



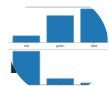
How to Implement Bayesian Optimization from Scratch in Python



Tweet

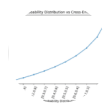
Share

Share



How to Calculate the KL Divergence for

On This Topic



A Gentle Introduction to Cross-Entropy for Machine Learning

Loving the Tutorials?

The [Probability for Machine Learning](#) EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

SEE WHAT'S INSIDE

Start Machine Learning

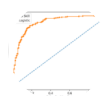


You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

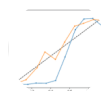
How to Develop a Probabilistic Model of Breast...

Start Machine Learning

Never miss a tutorial:**Picked for you:**

How to Use ROC Curves and Precision-Recall Curves for Classification in Python

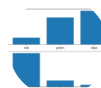
Model Selection Tips From Competitive Machine Learning



How and When to Use a Calibrated Classification Model with scikit-learn



How to Implement Bayesian Optimization from Scratch in Python



A Gentle Introduction to How to Calculate the KL Divergence for Machine Learning



A Gentle Introduction to Cross-Entropy for Machine Learning

Loving the Tutorials?

Feature Selection to Improve Accuracy and Decrease...

The [Probability for Machine Learning](#) EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Feature Importance and Feature Selection With...

About Jason Brownlee

Jason Brownlee, PhD is a machine learning specialist who teaches developers how to get results with modern machine learning methods via hands-on tutorials.

[View all posts by Jason Brownlee](#) →

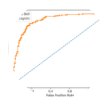
< [A Gentle Introduction to Logistic Regression With Maximur](#)

Start Machine Learning

Never miss a tutorial:



41 Responses to Probabilistic Model Selection with AIC, BIC, and MDL



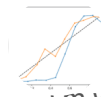
How to Use ROC Curves and Precision-Recall Curves for Classification in Python

Elie Kawerk

November 1, 2019 at 7:00 am #

REPLY ↩

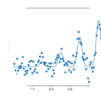
Hi Jason,



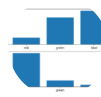
How and When to Use a Calibrated Classification Model with scikit-learn

Thank you for this nice post!

I'm wondering how to deal with non-parametric models strongly on the data (tree-based methods).

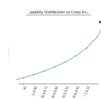


How to Implement Bayesian Optimization from Scratch in Python



How to Calculate the KL Divergence for Machine Learning

Good question.



A Gentle Introduction to Cross-Entropy for Machine Learning

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Pikachu November 15, 2019 at 12:54 am #

REPLY ↩

Loving the Tutorials?

Hi Jason,

The Probability for Machine Learning EBook is

where you'll find the **Really Good** stuff

Instead of using mse, I have used log loss value in the equation of AIC and found good result on my dataset. ' ' ' the value with the AIC value found from statsmodel of python (logit).

Here I ' ' ' value performed better than mse. Can you please explain why has that happened?

Jason Brownlee November 15, 2019 at 7:55 am #

REPLY ↩

Log loss is for classification, e.g. logistic regression. The example in the tutorial is linear regression, e.g. predicting a numerical value and log loss is inappropriate.

Pikachu November 18, 2019 at 2:37 pm #

REPLY ↩

For classification algorithms (listed below) other than Logistic Regression, should we always use Log Loss for calculating the AIC?

List of other classification algorithms:

– k Nearest Neighbor

Start Machine Learning

Never miss a tutorial:

- SVC

- Naive Bayes

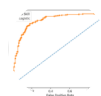


- Linear Discriminant Analysis

- Decision Tree

- Random Forest

Picked for you:



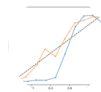
How to Use ROC Curves and Precision-Recall Curves for Classification in Python

Jason Brownlee

November 19, 2019 at 7:37 am #

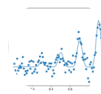
REPLY ↩

I don't think so.



How and When to Use a Calibrated Classification Model with scikit-learn

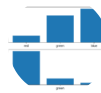
Each algorithm will require its own AIC calculation to be derived, at least that is my understanding.



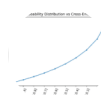
How to Implement Bayesian Optimization from Scratch in Python

James Bowery

January 7, 2020 at 4:22 pm #



Has anyone tried implementing MDL for k-How to Calculate the KL Divergence for Machine Learning



A Gentle Introduction to Cross-Entropy for Machine Learning

No, sorry.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Loving the Tutorials?

Jihoon Jang

February 4, 2020 at 3:56 am #

REPLY ↩

The Probability for Machine Learning EBook is

where you'll find the **Really Good** stuff.

If I use

>> SEE WHAT'S INSIDE

the number of parameters to calculate AIC ?

Thank you

Jason Brownlee

February 4, 2020 at 7:58 am #

REPLY ↩

model.summary() can access the number of parameters.

Jihoon Jang

February 4, 2020 at 8:54 pm #

REPLY ↩

Hi Jason 😊

I have a one more question.

As you told me, I just run "model.summary()",

then it said "'MLPRegressor' object has no

Start Machine Learning

How can I fix this problem?
Never miss a tutorial:

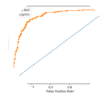
Thank you !



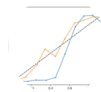
Picked for you:

Jason Brownlee February 5, 2020 at 8:07 am #

REPLY ↩



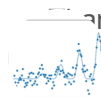
How to Use ROC Curves and Precision-Recall Curves for Classification in Python
 Perhaps your model has no layers?



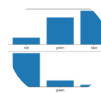
How and When to Use a Calibrated Classification Model with scikit-learn
 February 4, 2020 at 3:59 am #

REPLY ↩

Or do I just insert number of hyper-param



How to Implement Bayesian Optimization from Scratch in Python



How to Use the Akaike Information Criterion for Machine Learning
 February 4, 2020 at 7:50 am #
 I believe AIC requires a specialized ca



A Gentle Introduction to Cross-Entropy for Machine Learning

Jihoon Jang February 4, 2020 at 8:27 pm #

REPLY ↩

Hi Jason,

Loving the Tutorials?

Thank you 😊

The Probability for Machine Learning EBook is
 Is a specialized calculation the number of parameters using model.summary() ?
 where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Jason Brownlee February 5, 2020 at 8:07 am #

REPLY ↩

No, I mean you will need to check the literature for how to calculate the metric for an MLP in an appropriate manner.

CMHennings February 28, 2020 at 12:08 am #

REPLY ↩

Jason, I'm finding your information and code examples most helpful as I work on my MS degree. Thank you for the time and effort it takes to compose these posts!!

To adapt the linear regression example for logistic regression, the calculation for AIC and BIC (line 9) requires adjustment, correct?

Earlier in this post you define the AIC and BIC calculations for Logistic Regression as:

$$\text{AIC} = -2/N * \text{LL} + 2 * k/N$$

$$\text{BIC} = -2 * \text{LL} + \log(N) * k$$

Start Machine Learning

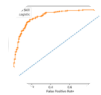
My understanding is line 9 needs replacement with these equations and LL should be replaced with the logistic regression log-likelihood calculation described in your "Gentle Introduction to Logistic



Regression" post:

$\text{log-likelihood} = \log(\text{yhat}) * y + \log(1 - \text{yhat}) * (1 - y)$

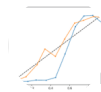
Picked for you: Am I on the right track?



How to Use ROC Curves and Precision-Recall Curves for Classification in Python

Jason Brownlee February 28, 2020 at 6:12 am #

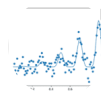
REPLY ↩



How and when to use a formula for the AIC and BIC metrics.

Classification Model with scikit-learn

I tried to provide standard calculations and link approaches that I have seen described.



How to Implement Bayesian Optimization from Scratch in Python

How to Calculate the KL Divergence for Machine Learning

El BIC es el mejor criterio para seleccionar



A Gentle Introduction to Cross-Entropy for Machine Learning

Jason Brownlee March 4, 2020 at 1:32 pm #

REPLY ↩

It really depends on your project goals. There is never a best anything for all cases.
Loving the Tutorials?

The **Probability for Machine Learning** EBook is where you'll find the **Really Good** stuff.

Jacob March 28, 2020 at 2:22 am #

REPLY ↩

>> SEE WHAT'S INSIDE

Thank you for the useful article. What I miss is how can MSE stand in place of L despite the fact that a model is better if it has smaller MSE and not larger, like when we deal with L?

Grzegorz Kępiś April 22, 2020 at 12:29 am #

REPLY ↩

Hello again Jason, thank you for good lecture!

Question: Probabilistic model selection include complexity penalty along to error prediction minimization. I may ask, why don't we just focus on test error score? I guess that the reasons for this are:

- 1) We prefer simple models (easier to interpret)
- 2) Simpler models normally require less memory, less train/test execution time.

Is it correct and maybe you can add something to this list?

Regards!

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

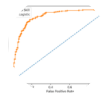
Never miss a tutorial **Jason Brownlee** April 22, 2020 at 5:58 am #

REPLY ↩



but more important: simpler explanations are more likely correct and generalize than complex explanations.

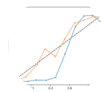
Picked for you:



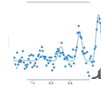
Dr. James T. Walker April 22, 2020 at 7:20 pm #

REPLY ↩

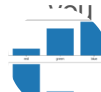
Hello: I am an Auxologist trying to develop a model for describing human height growth data from birth to maturity (0 to 21 years). My model assumes that human growth is due to the combination of nine logistic growth components. When I presented the paper at an international conference, I was told that I could use the AIC and BIC method for selecting the best model. I have a question: what is needed for these data. Do you agree? Each data set has a different number of components. The AIC values vs n shows a u-shaped curve, showing that the best model is the one with the lowest AIC value. I fit the components to a data set containing measurements), the AIC values and plots change, Can I take this approach? Why does the AIC change when I change the number of components?



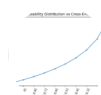
How and When to Use a Calibrated Classification Model with scikit-learn



How to Implement Bayesian Optimization from Scratch in Python



How to Calculate the KL Divergence for Machine Learning



A Gentle Introduction to Cross-Entropy for Machine Learning

Perhaps. Although you are preparing a descriptive rather than predictive model, e.g. statistics rather than machine learning.

The metrics change when the number of elements change because the number of elements impact the complexity of the model.

The **Probability for Machine Learning** EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE pm #

REPLY ↩

Hi Jason,

Nicely articulated indeed!

I got some ambiguity here, btw. Computing statsmodels's `aic(3026)` on sklearn boston dataset is showing different result than this manual aic computation(1565).

Any help would be much appreciated!

Jason Brownlee May 31, 2020 at 6:21 am #

REPLY ↩

Perhaps there is a difference in the implementation.

Nkue June 23, 2020 at 8:56 pm #

REPLY ↩

Start Machine Learning

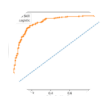
Hi Jason

Never miss a tutorial:

Thank you for this blog. Could you please provide R code for the calculation of the MDL for linear



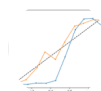
regards

Picked for you:

How to Use ROC Curves and Precision-Recall Curves for Classification in Python

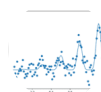
Jason Brownlee June 24, 2020 at 6:31 am #

REPLY ↩



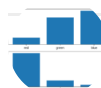
How and when to use a generalized Classification Model with scikit-learn

Thanks for the suggestion, perhaps in the future.



How to Implement Bayesian Optimization from Scratch in Python

Or Python code



How to Calculate the KL Divergence for Machine Learning

Ghizlan September 22, 2020 at 12:22 am #



A Gentle Introduction to Cross-Entropy for Machine Learning

ve one question about using AIC to calculate the goodness of fit for neural network and random forest models ? If we can use it, which package to use in R

Best regards

Loving the Tutorials?The [Probability for Machine Learning](#) EBook iswhere you'll find the **Really Good** stuff.

Jason Brownlee September 22, 2020 at 6:49 am #

REPLY ↩

>> SEE WHAT'S INSIDE

...and which one package off hand, perhaps try a google search.

gizlane September 22, 2020 at 7:34 am #

REPLY ↩

Thank you Jason

Jason Brownlee September 22, 2020 at 7:45 am #

REPLY ↩

You're welcome.

Mansi September 25, 2020 at 5:30 am #

REPLY ↩

Hi Jason.

Start Machine Learning

I've built two models.
Never miss a tutorial:
 Model 1- AIC= 8906



Model 2 AIC= -9501

Is it right to compare negative AIC with positive AIC? Also which AIC above, proves a better model?

Picked for you:



How to Use ROC Curves and Precision-Recall Curves for Classification in Python
Jason Brownlee September 25, 2020 at 6:41 am #

REPLY ↩

Sorry, I don't interpret results.



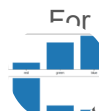
How and When to Use a Calibrated Classification Model with scikit-learn

Matt November 30, 2020 at 1:34 am #



How to Implement Bayesian Optimization from Scratch in Python
Jason

Thanks for this great post.



For purposes of calculating BIC for a linear regression
 How to Calculate the KL Divergence for Machine Learning
 ber of terms (in the example above, 1 for the in
 include the variance, also being estimated in the fit

I'm wondering because this <https://en.wikipedia.org>
 A Gentle Introduction to Cross-Entropy for Machine Learning
 ms to get to the variance.



Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.**
 Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Jason Brownlee November 30, 2020 at 6:37 am #
Loving the Tutorials?

REPLY ↩

Perhaps, I based my description on the textbooks listed at the end of the tutorial.
 The Probability for Machine Learning Ebook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Jon March 14, 2021 at 12:02 pm #

REPLY ↩

How do you calculate AIC and BIC for Logistic Regression Models in Python?

Jason Brownlee March 15, 2021 at 5:51 am #

REPLY ↩

The above calculations will help directly.

Neetika May 26, 2021 at 1:23 am #

REPLY ↩

It is really nicely articulated article. I had initially struggled to understand these concepts, but your article made it crystal clear. I wanted to implement new criteria for model selection via GLM based approach – stepwise forward regression using R or Python. Could you please suggest what parameters I

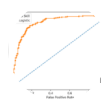
Start Machine Learning

can consider for defining criteria. Also in case you have sample code for GLM or stepwise forward regression, it would be great help.



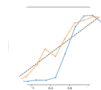
Picked for you: **Jason Brownlee** May 26, 2021 at 5:55 am #

REPLY ↩



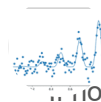
You're welcome
How to Use ROC Curves and Precision-Recall Curves for Classification in Python

Perhaps you can implement the algorithm from a paper or textbook or start with an existing implementation.



How and When to Use a Calibrated Classification Model with scikit-learn

Martin Zwanzig June 7, 2021 at 6:12 pm #



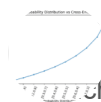
How to Implement Bayesian Optimization from Scratch in Python

Do the methods listed above (make_regression, ...), the following principle is ignored:



"The theory of AIC requires that the log-likelihood of models not fitted by maximum likelihood, their AIC is calculated using the log-likelihood of the model fitted by maximum likelihood (package 'stats'))

In other words: When linear models are not fitted by maximum likelihood, their AIC is calculated using the log-likelihood of the model fitted by maximum likelihood (package 'stats'))



A Gentle Introduction to Cross-Entropy for Machine Learning

Other measures such as AIC and BIC. Such measures also cannot be used to compare models fitted to a different response (or when the response has been transformed in one case but not the other).

@ Dr. James T. Walker: Also not to the same response considering a different sample size!

Loving the Tutorials?

The **Probability for Machine Learning** EBook is

where you'll find the **Really Good** stuff.

Jason Brownlee June 8, 2021 at 7:14 am #

REPLY ↩

>> SEE WHAT'S INSIDE

... introduces a new regression dataset, learn more here:

https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make_regression.html

It's just the context for the tutorial.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Leave a Reply

Start Machine Learning

Never miss a tutorial:

Name (required)



Email (will not be published) (required)

Picked for you:

Website



How to Use ROC Curves and Precision-Recall Curves for Classification in Python

SUBMIT COMMENT



How and When to Use a Calibrated Classification Model with scikit-learn



Welcome!

I'm Jason Brownlee PhD

and I help developers get results with

Read more

How to Implement Bayesian Optimization from Scratch in Python



How to Calculate the KL Divergence for Machine Learning



A Gentle Introduction to Cross-Entropy for Machine Learning

Start Machine Learning



You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Loving the Tutorials?

The [Probability for Machine Learning](#) EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Start Machine Learning

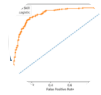
Never miss a tutorial:



Picked for you:

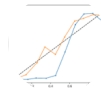
© 2021 Machine Learning Mastery. All Rights Reserved.

[LinkedIn](#) | [Twitter](#) | [Facebook](#) | [Newsletter](#) | [RSS](#)

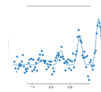


[How to Use ROC Curves and Precision-](#)

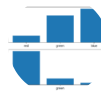
[Recall Curves for Classification in Python](#)



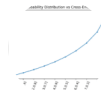
[How and When to Use a Calibrated Classification Model with scikit-learn](#)



[How to Implement Bayesian Optimization from Scratch in Python](#)



[How to Calculate the KL Divergence for Machine Learning](#)



[A Gentle Introduction to Cross-Entropy for Machine Learning](#)

Start Machine Learning



You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Loving the Tutorials?

The [Probability for Machine Learning](#) EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Start Machine Learning