

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

A Gentle Introduction to Cross-Entropy for Machine Learning



How to Calculate the KL Divergence for Brownlee on October 21, 2019 in **Probability**
Machine Learning

Tweet

Share

Share



How to Implement Bayesian Optimization from Scratch in Python
dated on December 22, 2020

Cross-entropy is commonly used in machine learning as a loss function.



A Gentle Introduction to Cross-Entropy for Machine Learning

Entropy is a measure from the field of information theory, building upon [entropy](#) and generally quantifying the difference between two probability distributions. It is closely related to but is different from [KL divergence](#) that calculates the relative entropy between two probability distributions, whereas cross-entropy can be thought to calculate the total entropy between the distributions.

Cross-entropy is often confused with [logistic loss](#), called [log loss](#). Although the two measures are derived from a different source, when used as loss functions for classification models, both measures calculate the same quantity and can be used interchangeably.

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

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

In this tutorial, you will discover cross-entropy for machine learning.

>> SEE WHAT'S INSIDE

After completing this tutorial, you will know:

- How to calculate cross-entropy from scratch and using standard machine learning libraries.
- Cross-entropy can be used as a loss function when optimizing classification models like logistic regression and artificial neural networks.
- Cross-entropy is different from KL divergence but can be calculated using KL divergence, and is different from log loss but calculates the same quantity when used as a loss function.

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](#)

• **Update Oct/2019:** Gave an example of cross-entropy for identical distributions and updated description for this case (thanks Ron U). Added an example of calculating the entropy of the known class label.



• **Update Nov/2019:** Improved structure and added more explanation of entropy. Added intuition for predicted class probabilities.

• **Picked for you:** **Update Dec/2020:** Tweaked the introduction to information and entropy to be clearer.

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 Calculate the KL Divergence for Machine Learning

How to Implement Bayesian Optimization from Scratch in Python

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**.

a Gentle Introduction to Cross-Entropy for Machine Learning
Photo by [Jerome Bon](#), some rights reserved.

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

Tutorial Overview >> SEE WHAT'S INSIDE

This tutorial is divided into five parts; they are:

1. What Is Cross-Entropy?
2. Cross-Entropy Versus KL Divergence
3. How to Calculate Cross-Entropy
 1. Two Discrete Probability Distributions
 2. Calculate Cross-Entropy Between Distributions
 3. Calculate Cross-Entropy Between a Distribution and Itself
 4. Calculate Cross-Entropy Using KL Divergence
4. Cross-Entropy as a Loss Function
 1. Calculate Entropy for Class Labels

Start Machine Learning

Never miss a tutorial:

2. Calculate Cross-Entropy Between Class Labels and Probabilities
3. Calculate Cross-Entropy Using Keras
4. Intuition for Cross-Entropy on Predicted Probabilities
5. Cross-Entropy Versus Log Loss

1. Log Loss is the Negative Log Likelihood

Picked for you:

2. Log Loss and Cross Entropy Calculate the Same Thing

What Is Cross-Entropy?

Cross-entropy is a measure of the difference between two probability distributions for a given random variable or set of events.

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

[Classification Model with scikit-learn](#)

Y might recall that **information** quantifies the number of bits required to describe the event.

Lower probability events have more information, higher probability events have less information.

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

For example, if we like to describe the “surprise” of an event, it is, meaning it contains more information.

[How to Implement Bayesian Optimization](#)

Low Probability Event (surprising): More information
High Probability Event (unsurprising): Less information

Information $h(x)$ can be calculated for an event x , given its probability $P(x)$.

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

Entropy is the number of bits required to transmit a randomly selected event from a probability distribution. A skewed distribution has a low entropy, whereas a distribution where events have equal probability has a larger entropy.

Loving the Tutorials?

A skewed probability distribution has less “surprise” and in turn a low entropy because likely events dominate. Balanced distribution are more surprising and turn have higher entropy because events are equally likely.

- **Skeved Probability Distribution (unsurprising):** Low entropy.
- **Balanced Probability Distribution (surprising):** High entropy.

Entropy $H(x)$ can be calculated for a random variable with a set of x in X discrete states and their probability $P(x)$ as follows:

- $H(X) = - \sum_{x \in X} P(x) * \log(P(x))$

If you would like to know more about calculating information for events and entropy for distributions see this tutorial:

- [A Gentle Introduction to Information Entropy](#)

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

Cross-entropy builds upon the idea of entropy from information theory and calculates the number of bits required to represent or transmit an average event from one distribution compared to another distribution.



... the cross entropy is the average number of bits needed to encode data coming from a source with distribution p when we use model q ...

Picked for you:

 [How Machine Learning: A Probabilistic Perspective, 2012.](#)


[Recall Curves for Classification in Python](#)

The intuition for this definition comes if we consider a target or underlying probability distribution P and an approximation of the target distribution Q , then the cross-entropy of Q from P is the number of additional

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

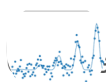
[Classification Model with scikit-learn](#)

The cross-entropy between two probability distribution

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

[Machine Learning](#)

Where $H()$ is the cross-entropy function, P may be the target distribution.

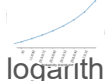
 [How to Implement Bayesian Optimization](#)

[from Scratch in Python](#)

entropy can be calculated using the probabilistic

- $H(P, Q) = - \sum_{x \in X} P(x) * \log(Q(x))$

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

 [Machine Learning](#)
 $P(x)$ is the probability of the event x in P , $Q(x)$ is the probability of event x in Q and \log is the base-2 logarithm, meaning that the results are in bits. If the base-e or natural logarithm is used instead, the result will have the units called nats.

This calculation is for discrete probability distributions, although a similar calculation can be used for continuous probability distributions using the integral across the events instead of the sum.

The result will be a positive number measured in bits and will be equal to the entropy of the distribution if the two probability distributions are identical.

Note: this is the joint probability, or more specifically, the **joint entropy** between P and Q . This is misleading as we are scoring the difference between probability distributions with cross-entropy. Whereas, joint entropy is a different concept that uses the same notation and instead calculates the uncertainty across two (or more) random variables.

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

Take my free 7-day email crash course now (with sample code).

Click to sign-up and also get a free PDF Ebook version of the course.

Start Machine Learning

Never miss a tutorial:

Download Your FREE Mini-Course



Cross-Entropy Versus KL Divergence



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

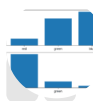
Cross-entropy is related to divergence measures, such as the [Kullback-Leibler](#), or [KL](#), [Divergence](#) that quantifies how much one distribution differs from another.



How and When to Use a Calibrated

Classification Model with scikit-learn

Typically, the KL divergence measures a very similar number of extra bits required to represent a message



How to Calculate the KL Divergence for

In other words, the KL divergence is the average

data, due to the fact that we used distribution

p.



How to Implement Bayesian Optimization

from Scratch in Python

Chapter 58, Machine Learning: A Probabilistic Perspective

As such, the KL divergence is often referred to as the

A Gentle Introduction to Cross-Entropy for

Machine Learning



Cross-Entropy: Average number of total bits to represent an event from Q instead of P.

- **Relative Entropy** (*KL Divergence*): Average number of extra bits to represent an event from Q instead of P.

KL divergence is the negative sum of probability of each event in P multiplied by the log of the probability of the event in Q over the probability of the event in P. Typically, log base-2 so that the result is measured in bits.

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

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

$$KL(P \parallel Q) = - \sum_{x \in X} P(x) \cdot \log(Q(x) / P(x))$$

The value [>> SEE WHAT'S INSIDE](#) divergence for a given event.

As such, we can calculate the cross-entropy by adding the entropy of the distribution plus the additional entropy calculated by the KL divergence. This is intuitive, given the definition of both calculations; for example:

- $H(P, Q) = H(P) + KL(P \parallel Q)$

Where $H(P, Q)$ is the cross-entropy of Q from P, $H(P)$ is the entropy of P and $KL(P \parallel Q)$ is the divergence of Q from P.

Entropy can be calculated for a probability distribution as the negative sum of the probability for each event multiplied by the log of the probability for the event, with

Start Machine Learning

• $H(P) = - \sum_{x \in X} p(x) * \log(p(x))$

Never miss a tutorial:

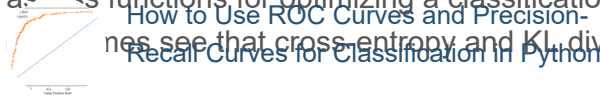
Like KL divergence, cross-entropy is not symmetrical, meaning that:



• $H(P, Q) \neq H(Q, P)$

Picked for you:

As we will see later, both cross-entropy and KL divergence calculate the same quantity when they are used as loss functions for optimizing a classification predictive model. It is under this context that you might



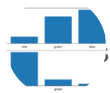
How to Use ROC Curves and Precision-Recall Curves for Classification in Python

For a lot more detail on the KL Divergence, see the tutorial:



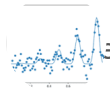
How and When to Use a Calibrated Classification Model with scikit-learn

How to Calculate Cross-Entropy



How to Calculate the KL Divergence for Machine Learning

Two Discrete Probability Distributions



How to Implement Bayesian Optimization for a Random Variable with three discrete events

We may have two different probability distributions for



A Gentle Introduction to Cross-Entropy for

```
1 # ... Machine Learning
2 # define distributions
3 events = ['red', 'green', 'blue']
4 p = [0.10, 0.40, 0.50]
5 q = [0.80, 0.15, 0.05]
```

We can plot a bar chart of these probabilities to compare them directly as probability histograms.

Loving the Tutorials?

The complete example is listed below.

The Probability for Machine Learning EBook is

```
1 # where you'll find the Really Good stuff.
2 from matplotlib import pyplot
3 # define distributions
4 events = ['red', 'green', 'blue']
5 p = [0.10, 0.40, 0.50]
6 q = [0.80, 0.15, 0.05]
7 print('P=%.3f Q=%.3f' % (sum(p), sum(q)))
8 # plot first distribution
9 pyplot.subplot(2,1,1)
10 pyplot.bar(events, p)
11 # plot second distribution
12 pyplot.subplot(2,1,2)
13 pyplot.bar(events, q)
14 # show the plot
15 pyplot.show()
```

Running the example creates a histogram for each probability distribution, allowing the probabilities for each event to be directly compared.

We can see that indeed the distributions are different.

Start Machine Learning

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

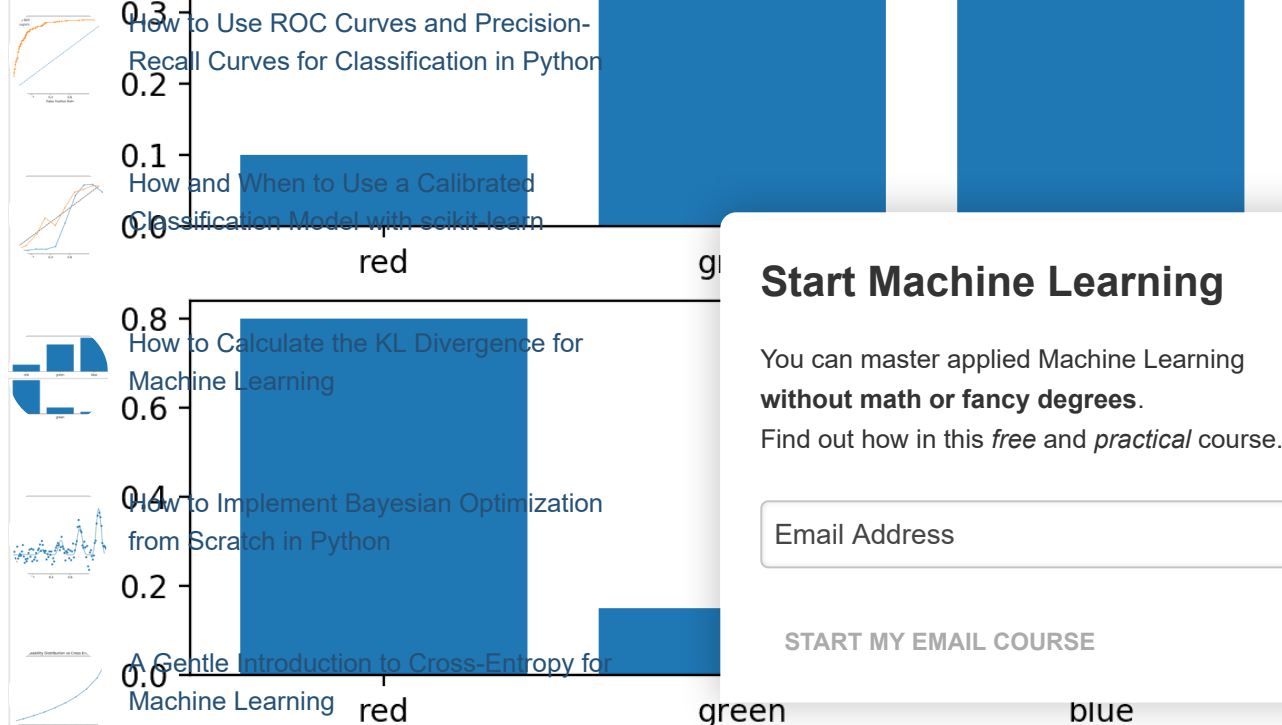
START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial:



Picked for you:



Histogram of Two Different Probability Distributions for the Same Random Variable

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

Calculate Cross-Entropy Between Distributions

Next, we can develop a function to calculate the cross-entropy between the two distributions.

The Probability for Machine Learning Ebook is

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

We will use log base-2 to ensure the result has units in bits.

```
>> SEE WHAT'S INSIDE
1 # calculate cross entropy
2 def cross_entropy(p, q):
3     return -sum([p[i]*log2(q[i]) for i in range(len(p))])
```

We can then use this function to calculate the cross-entropy of P from Q, as well as the reverse, Q from P.

```
1 ...
2 # calculate cross entropy H(P, Q)
3 ce_pq = cross_entropy(p, q)
4 print('H(P, Q): %.3f bits' % ce_pq)
5 # calculate cross entropy H(Q, P)
6 ce_qp = cross_entropy(q, p)
7 print('H(Q, P): %.3f bits' % ce_qp)
```

Tying this all together, the complete example is listed below.

```
1 # example of calculating cross entropy
```

Start Machine Learning

```

2 from math import log2
3 def cross_entropy(p, q):
4     return -sum([p[i]*log2(q[i]) for i in range(len(p))])
5
6 # define data
7 p = [0.10, 0.40, 0.50]
8 q = [0.80, 0.15, 0.05]
9
10 # calculate cross entropy H(P, Q)
11 ce_pq = cross_entropy(p, q)
12 print('H(P, Q): %.3f bits' % ce_pq)
13
14 # calculate cross entropy H(Q, P)
15 ce_qp = cross_entropy(q, p)
16 print('H(Q, P): %.3f bits' % ce_qp)

```

How and When to Use a Calibrated
 the example first calculates the cross-entropy of Q from P as just over 3 bits, then P from Q as just
 3 bits.

```

1 H(P, Q): 3.288 bits
2 H(Q, P): 2.906 bits

```

How to Implement Bayesian Optimization
 from Scratch in Python

If two probability distributions are the same, then the cross-entropy is the same as the entropy of the distribution.

A Gentle Introduction to Cross-Entropy for
 Machine Learning

demonstrate this by calculating the cross-entropy of P vs P and Q vs Q.

The complete example is listed below.

```

1 # example of calculating cross entropy for identical distributions
2 from math import log2
3
4 # calculate cross entropy
5 def cross_entropy(p, q):
6     return -sum([p[i]*log2(q[i]) for i in range(len(p))])
7
8 # define data
9 p = [0.10, 0.40, 0.50]
10 q = [0.80, 0.15, 0.05]
11
12 # calculate cross entropy H(P, P)
13 ce_pp = cross_entropy(p, p)
14 print('H(P, P): %.3f bits' % ce_pp)
15
16 # calculate cross entropy H(Q, Q)
17 ce_qq = cross_entropy(q, q)
18 print('H(Q, Q): %.3f bits' % ce_qq)

```

Running the example first calculates the cross-entropy of Q vs Q which is calculated as the entropy for Q, and P vs P which is calculated as the entropy for P.

```

1 H(P, P): 1.361 bits
2 H(Q, Q): 0.884 bits

```

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

Calculate Cross-Entropy Using KL Divergence

We can also calculate the cross-entropy using the KL divergence.

The cross-entropy calculated with KL divergence should be identical, and it may be interesting to calculate the KL divergence between the distributions as well to see the relative entropy or additional bits required instead of the total bits calculated by the cross-entropy.

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

Next, we can define a function to calculate the entropy

How to Calculate the KL Divergence for Machine Entropy

How to calculate the Bayes Cross-Entropy from Scratch in Python

A Gentle Introduction to Cross-Entropy for Machine Learning

the example is simple, we can compare the cross-entropy for $H(P, Q)$ to the KL divergence $KL(P || Q)$ and the entropy $H(P)$.

```
1 # calculate the kl divergence KL(P || Q)
2 def kl_divergence(p, q):
3     return sum(p[i] * log2(p[i]/q[i]) for i in range(len(p)))
```

Next, we can define a function to calculate the entropy

```
1 # calculate entropy H(P)
2 def entropy(p):
3     return -sum([p[i] * log2(p[i]) for i in range(len(p))])
```

How to calculate the Bayes Cross-Entropy from Scratch in Python

```
1 # calculate cross entropy H(P, Q)
2 def cross_entropy(p, q):
3     return entropy(p) + kl_divergence(p, q)
```

A Gentle Introduction to Cross-Entropy for Machine Learning

the example is simple, we can compare the cross-entropy for $H(P, Q)$ to the KL divergence $KL(P || Q)$ and the entropy $H(P)$.

The complete example is listed below.

```
1 # example of calculating cross entropy with kl divergence
2 from math import log2
3
4 # calculate the kl divergence KL(P || Q)
5 def kl_divergence(p, q):
6     return sum(p[i] * log2(p[i]/q[i]) for i in range(len(p)))
7
8 # calculate entropy H(P)
9 def entropy(p):
10    return -sum([p[i] * log2(p[i]) for i in range(len(p))])
11
12 # calculate cross entropy H(P, Q)
13 def cross_entropy(p, q):
14    return entropy(p) + kl_divergence(p, q)
15
16 # define data
17 p = [0.10, 0.40, 0.50]
18 q = [0.80, 0.15, 0.05]
19 # calculate H(P)
20 en_p = entropy(p)
21 print('H(P): %.3f bits' % en_p)
22 # calculate kl divergence KL(P || Q)
23 kl_pq = kl_divergence(p, q)
24 print('KL(P || Q): %.3f bits' % kl_pq)
25 # calculate cross entropy H(P, Q)
26 ce_pq = cross_entropy(p, q)
```

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

```
27 print('H(P, Q): %.3f bits' % ce_pq)
```

Never miss a tutorial:

Running the example, we can see that the cross-entropy score of 3.288 bits is comprised of the entropy of 1.361 bits and the KL divergence of 1.927 bits calculated by the KL divergence.

This is a useful example that clearly illustrates the relationship between all three calculations.

Picked for you:

```
1 H(P): 1.361 bits
2 KL(P || Q): 1.927 bits
3 H(P, Q): 3.288 bits
```

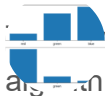
[Recall Curves for Classification in Python](#)



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

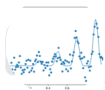
Cross-Entropy as a Loss Function

Cross-entropy is widely used as a loss function when



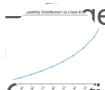
[How to Calculate the KL Divergence for Examples that you may encounter include the log](#)

and artificial neural networks that can be u



[How to Implement Bayesian Optimization ... using the cross-entropy error function instead from Scratch in Python](#)

problem leads to faster training as well as imp



[The 235 Pattern Recognition and Machine Learning](#)

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

Classification problems are those that involve one or more input variables and the prediction of a class label.

Classification tasks that have just two labels for the output variable are referred to as binary classification problems, whereas those problems with more than two labels are referred to as categorical or multi-class classification problems.

Loving the Tutorials?

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

- **Binary Classification:** Task of predicting one of two class labels for a given example.
- **Multi-Class Classification:** Task of predicting one of more than two class labels for a given example.

>> SEE WHAT'S INSIDE

We can see that the idea of cross-entropy may be useful for optimizing a classification model.

Each example has a known class label with a probability of 1.0, and a probability of 0.0 for all other labels. A model can estimate the probability of an example belonging to each class label. Cross-entropy can then be used to calculate the difference between the two probability distributions.

As such, we can map the classification of one example onto the idea of a random variable with a probability distribution as follows:

- **Random Variable:** The example for which we require a predicted class label.
- **Events:** Each class label that could be predicted.

Start Machine Learning

Never miss a tutorial:



Our model seeks to approximate the target probability distribution Q .

Picked for you:

In the language of classification, these are the actual and the predicted probabilities, or y and \hat{y} .



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

[Recall in Classification](#) in Python. The probability of each class label for an example in the dataset (P).

- **Predicted Probability (\hat{y}):** The probability of each class label an example predicted by the model (Q).



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

[Classification Model with cross-entropy](#) for a single

described above; for example.



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

Where each x in X is a class label that could be assigned to the class label and 0 for all other labels.



[How to Implement Bayesian Optimization](#)

[from Scratch in Python](#)

cross-entropy for a single example in a binary classification operation as follows:



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

You may see this form of calculating cross-entropy cited in textbooks.

If there are just two class labels, the probability is modeled as the [Bernoulli distribution](#) for the positive class label. This means that the probability for class 1 is predicted by the model directly, and the probability for class 0 is given as one minus the predicted probability, for example:

- **Predicted $P(\text{class } 0) = 1 - \hat{y}$**
- **Predicted $P(\text{class } 1) = \hat{y}$**

When calculating cross-entropy for classification tasks, the base-e or natural logarithm is used. This means that the units are in nats, not bits.

We are often interested in minimizing the cross-entropy for the model across the entire training dataset. This is calculated by calculating the average cross-entropy across all training examples.

Calculate Entropy for Class Labels

Recall that when two distributions are identical, the cross-entropy between them is equal to the entropy for the probability distribution.

Class labels are encoded using the values 0 and 1 when preparing data for classification tasks.

Start Machine Learning

For example, if a classification problem has three classes, and an example has a label for the first class, then the probability distribution will be $[1, 0, 0]$. If an example has a label for the second class, it will have a probability distribution over the events as $[0, 1, 0]$. This is called a **one hot encoding**.



This probability distribution has no information as the outcome is certain. We know the class. Therefore the **Picked for you** variable is zero.

 [How to Use ROC Curves and Precision-Recall Curves for Classification in Python](#)

Pretend we have a classification problem with 3 classes, and we have one example that belongs to each class. We can represent each example as a discrete probability distribution with a 1.0 probability for the class which the example belongs to and a 0.0 probability for all other classes.



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

We can calculate the entropy of the probability distribution.

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

```
1 # entropy of examples from a classification task
2 from math import log2
3 from numpy import asarray
4
5 # calculate entropy
6 def entropy(p):
7     return -sum([p[i] * log2(p[i]) for i in range(len(p))])
8
9 # class 1
10 p = asarray([1, 0, 0]) + 1e-15
11 print(entropy(p))
12 # class 2
13 p = asarray([0, 1, 0]) + 1e-15
14 print(entropy(p))
15 # class 3
16 p = asarray([0, 0, 1]) + 1e-15
17 print(entropy(p))
```

Running the example calculates the entropy for each random variable. The Probability for Machine Learning EBook is

where you'll find the **Really Good** stuff. We can see that in each case, the entropy is 0.0 (actually a number very close to zero).

Note that `>> SEE WHAT'S INSIDE` all value to the 0.0 values to avoid the `log()` from blowing up, as we cannot calculate the log of 0.0.

```
1 9.805612959471341e-14
2 9.805612959471341e-14
3 9.805612959471341e-14
```

As such, the entropy of a known class label is always 0.0.

This means that the cross entropy of two distributions (real and predicted) that have the same probability distribution for a class label, will also always be 0.0.

Recall that when evaluating a model using cross-entropy on a training dataset that we average the cross-entropy across all examples in the dataset.

Start 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

Therefore, a cross-entropy of 0.0 when training a model indicates that the predicted class probabilities are identical to the probabilities in the training dataset, e.g. zero loss.



You could just as easily minimize the KL divergence as a loss function instead of the cross-entropy.

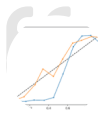
Picked for you:

Recall that the KL divergence is the extra bits required to transmit one variable compared to another. It is the cross-entropy without the entropy of the class label, which we know would be zero anyway.



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

[Recall Curves for Classification in Python](#), Minimizing the KL Divergence and the cross entropy for a classification task are identical.



[Minimizing this KL divergence corresponds exactly to minimizing the cross-entropy between the distributions](#)

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

— Page 132, [Deep Learning](#), 2016.



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

[Machine Learning](#)

ice, a cross-entropy loss of 0.0 often indicates that is another story.



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

Calculate Cross-Entropy Between Cla

The use of cross-entropy for classification often gives classes, mirroring the name of the classification task;

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

[Machine Learning](#)

Binary Cross-Entropy: Cross-entropy as a loss function for a binary classification task.

- **Categorical Cross-Entropy:** Cross-entropy as a loss function for a multi-class classification task.

We can make the use of cross-entropy as a loss function concrete with a worked example.

Loving the Tutorials?

Consider a two-class classification task with the following 10 actual class labels (P) and predicted class labels (Q).

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

```
1 .. where you'll find the Really Good stuff.
2 # define classification data
3 p = [1, 1, 1, 1, 1, 0, 0, 0, 0, 0]
4 q = [0.8, 0.9, 0.9, 0.6, 0.8, 0.1, 0.4, 0.2, 0.1, 0.3]
```

We can enumerate these probabilities and calculate the cross-entropy for each using the cross-entropy function developed in the previous section using `log()` (natural logarithm) instead of `log2()`.

```
1 # calculate cross entropy
2 def cross_entropy(p, q):
3     return -sum([p[i]*log(q[i]) for i in range(len(p))])
```

For each actual and predicted probability, we must convert the prediction into a distribution of probabilities across each event, in this case, the classes {0, 1} as 1 minus the probability for class 0 and probability for class 1.

We can then calculate the cross-entropy and repeat the process for all examples

[Start Machine Learning](#)


```

1 ...
2 # calculate cross entropy for each example
3 results = list()
4 for i in range(len(p)):
5     # create the distribution for each event {0, 1}
6     expected = [1.0 - p[i], p[i]]
7     predicted = [1.0 - q[i], q[i]]
8     # calculate cross entropy for the two events
9     ce = cross_entropy(expected, predicted)
10    print('>[y=%.1f, yhat=%.1f] ce: %.3f nats' % (p[i], q[i], ce))
11    results.append(ce)

```

Recall Curves for Classification in Python

we can calculate the average cross-entropy across the dataset and report it as the cross-entropy loss for the model on the dataset.

How and When to Use a Calibrated

```

1 ... Classification Model with scikit-learn
2 # calculate the average cross entropy
3 mean_ce = mean(results)
4 print('Average Cross Entropy: %.3f nats' % mean_ce)

```

How to Calculate the KL Divergence for Machine Learning

```

1 # calculate cross entropy for classification
2 from math import log
3 from numpy import mean
4
5 # calculate cross entropy
6 def cross_entropy(p, q):
7     return -sum([p[i]*log(q[i]) for i in range(len(p))])
8
9 # define classification data
10 p = [1, 1, 1, 1, 1, 0, 0, 0, 0, 0]
11 q = [0.8, 0.9, 0.9, 0.6, 0.8, 0.1, 0.4, 0.2, 0.1, 0.3]
12 # calculate cross entropy for each example
13 results = list()
14 for i in range(len(p)):
15     # create the distribution for each event {0, 1}
16     expected = [1.0 - p[i], p[i]]
17     predicted = [1.0 - q[i], q[i]]
18     # calculate cross entropy for the two events
19     ce = cross_entropy(expected, predicted)
20     print('>[y=%.1f, yhat=%.1f] ce: %.3f nats' % (p[i], q[i], ce))
21     results.append(ce)
22
23 # calculate the average cross entropy
24 mean_ce = mean(results)
25 print('Average Cross Entropy: %.3f nats' % mean_ce)

```

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 prints the actual and predicted probabilities for each example and the cross-entropy in nats.

The final average cross-entropy loss across all examples is reported, in this case, as 0.247 nats.

```

1 >[y=1.0, yhat=0.8] ce: 0.223 nats
2 >[y=1.0, yhat=0.9] ce: 0.105 nats
3 >[y=1.0, yhat=0.9] ce: 0.105 nats
4 >[y=1.0, yhat=0.6] ce: 0.511 nats
5 >[y=1.0, yhat=0.8] ce: 0.223 nats
6 >[y=0.0, yhat=0.1] ce: 0.105 nats
7 >[y=0.0, yhat=0.4] ce: 0.511 nats
8 >[y=0.0, yhat=0.2] ce: 0.223 nats

```

Start Machine Learning

```

9 >[y=0.0, yhat=0.1] ce: 0.105 nats
10 >[y=0.0, yhat=0.3] ce: 0.357 nats
11 Average Cross Entropy: 0.247 nats

```



is calculated when optimizing a logistic regression model or a neural network model under a cross-entropy loss function.

Picked for you:

Calculate Cross-Entropy Using Keras



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

confirm the same calculation by using the `binary_crossentropy()` function from the [Keras deep learning API](#) to calculate the cross-entropy loss for our small dataset.

Complete example is listed below

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

Notes: This example assumes that you have the [Keras](#) configured with a backend library such as [TensorFlow](#)

How to Calculate the KL Divergence for Machine Learning

```

1 # calculate cross entropy with keras
2 from numpy import asarray
3 from keras import backend
4 from keras.losses import binary_crossentropy
5 # prepare scratch Python data
6 p = asarray([1, 1, 1, 1, 1, 0, 0, 0, 0, 0])
7 q = asarray([0.8, 0.9, 0.9, 0.6, 0.8, 0.1, 0.1, 0.1, 0.1, 0.1])
8 # convert to keras variables
9 y_true = backend.variable(p)
10 y_pred = backend.variable(q)
11 # calculate the average cross-entropy
12 mean_ce = backend.eval(binary_crossentropy(y_true, y_pred))
13 print('Average Cross Entropy: %.3f nats' % mean_ce)

```

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, we can see that the same average cross-entropy loss of 0.247 nats is reported.

Loving the Tutorials?

This confirms the correct manual calculation of cross-entropy.

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

>> SEE WHAT'S INSIDE

Intuition for Cross-Entropy on Predicted Probabilities

We can further develop the intuition for the cross-entropy for predicted class probabilities.

For example, given that an average cross-entropy loss of 0.0 is a perfect model, what do average cross-entropy values greater than zero mean exactly?

We can explore this question no a binary classification problem where the class labels as 0 and 1. This is a discrete probability distribution with two events and a certain probability for one event and an impossible probability for the other event.

We can then calculate the cross entropy for different “predicted” probability distributions transitioning from a perfect match of the target distribution to the exact op

Start Machine Learning

We would expect that as the predicted probability distribution diverges further from the target distribution that the cross-entropy calculated will increase.



The example below implements this and plots the cross-entropy result for the predicted probability distribution compared to the target of [0, 1] for two events as we would see for the cross-entropy in a binary classification task.

```
1 # how to use ROC Curves and Precision-Recall Curves
2 from math import log
3 from matplotlib import pyplot
4
5 # calculate cross-entropy
6 def cross_entropy(p, q, ets=1e-15):
7     return -sum([p[i]*log(q[i]+ets) for i in range(len(p))])
8
9 # define the target distribution for two events
10 target = [0.0, 1.0]
11 # define probabilities for the first event
12 probs = [1.0, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1]
13 # calculate probability distributions for the two events
14 dists = [[1.0 - p, p] for p in probs]
15 # calculate cross-entropy for each distribution
16 ents = [cross_entropy(target, d) for d in dists]
17 # plot probabilities vs cross-entropy
18 pyplot.plot([1-p for p in probs], ents, marker='o')
19 pyplot.title('Probability Distribution vs Cross-Entropy')
20 pyplot.xticks([1-p for p in probs], ['%.1f,%' % (1-p, p)])
21 pyplot.subplots_adjust(bottom=0.2)
22 pyplot.xlabel('Probability Distribution')
23 pyplot.ylabel('Cross-Entropy (nats)')
24 pyplot.show()
```

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 calculates the cross-entropy score for each probability distribution then plots the results as a line plot.

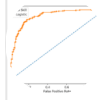
Loving the Tutorials?

We can see that as expected, cross-entropy starts at 0.0 (far left point) when the predicted probability distribution matches the target distribution, then steadily increases as the predicted probability distribution diverges.

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

We can also see a dramatic leap in cross-entropy when the predicted probability distribution is the exact opposite of the target, that is, [1, 0] compared to the target of [0, 1].

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 Calculate the KL Divergence for Machine Learning



How to Implement Bayesian Optimization from Scratch in Python



A Gentle Introduction to Cross-Entropy for Machine Learning

Line Plot of Probability Distribution vs Cross-Entropy for a Binary Classification Task

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?

We are not going to have a model that predicts the exact opposite probability distribution for all cases on a binary classification task.

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

The update [SEE WHAT'S INSIDE](#) listed below.

```
1 # cross-entropy for predicted probability distribution vs label
2 from math import log
3 from matplotlib import pyplot
4
5 # calculate cross-entropy
6 def cross_entropy(p, q, ets=1e-15):
7     return -sum([p[i]*log(q[i]+ets) for i in range(len(p))])
8
9 # define the target distribution for two events
10 target = [0.0, 1.0]
11 # define probabilities for the first event
12 probs = [1.0, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1]
13 # create probability distributions for the two events
14 dists = [[1.0 - p, p] for p in probs]
15 # calculate cross-entropy for each distribution
16 ents = [cross_entropy(target, d) for d in dists]
```

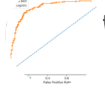
Start Machine Learning

```

17 # plot probability distribution vs cross-entropy
18 pyplot.plot([1-p for p in probs], ents, marker='.')
19 pyplot.title('Probability Distribution vs Cross-Entropy')
20 pyplot.xticks([1-p for p in probs], ['%.1f,%.1f'%(d[0],d[1]) for d in dists], rotation=70)
21 pyplot.subplots_adjust(bottom=0.2)
22 pyplot.xlabel('Probability Distribution')
23 pyplot.ylabel('Cross-Entropy (nats)')
24 pyplot.show()

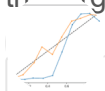
```

Running the example gives a much better idea of the relationship between the divergence in probability



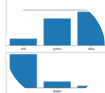
How to Use ROC Curves and Precision-Recall Curves for Classification in Python

We can see a super-linear relationship where the more the predicted probability distribution diverges from the target, the larger the increase in cross-entropy.

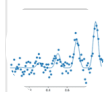


How and When to Use a Calibrated

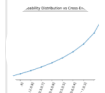
Classification Model with scikit-learn



How to Calculate the KL Divergence for Machine Learning



How to Implement Bayesian Optimization from Scratch in Python



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

Line Plot of Probability Distribution vs Cross-Entropy for a Binary Classification Task With Extreme Case Removed

A plot like this can be used as a guide for interpreting the average cross-entropy reported for a model for a binary classification dataset.

For example, you can use these cross-entropy values to interpret the mean cross-entropy reported by Keras for a neural network model on a binary classification task, or a binary classification model in scikit-learn evaluated using the logloss metric.

Start Machine Learning

You can use it to answer the general question:
Never miss a tutorial!



Picked for you:

If you are working in nats (and you usually are) and you are getting mean cross-entropy less than 0.2, you are off to a good start, and less than 0.1 or 0.05 is even better.



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

[Recall Curves if you are getting mean](#)

other hand, if you are getting mean cross-entropy greater than 0.2 or 0.3 you can probably improve, and if you are getting a mean cross-entropy greater than 1.0, then something is going on and you're making poor probability predictions on many examples in your dataset.



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

[Classification Model with scikit-learn](#)

[summarise these intuitions for the mean cross](#)

- **Cross-Entropy = 0.00:** Perfect probabilities.



[How to Calculate the K-L Divergence for](#)

[Machine Learning](#)

- **Cross-Entropy < 0.02:** Great probabilities.

- **Cross-Entropy < 0.05:** On the right track.

- **Cross-Entropy < 0.20:** Fine.

- **Cross-Entropy > 0.30:** Not great.



[How to Implement Bayesian Optimization](#)

[from Scratch in Python](#)

- **Cross-Entropy > 1.00:** Terrible.

- **Cross-Entropy > 2.00** Something is broken.

This listing will provide a useful guide when interpreting
[A Gentle Introduction to Cross-Entropy for](#)
[ion model, or your artificial neural network model](#)
[Machine Learning](#)



You can also calculate separate mean cross-entropy scores per-class and help tease out on which classes you're model has good probabilities, and which it might be messing up.

Cross-Entropy Versus Log Loss

Cross-Entropy is not Log Loss, but they calculate the same quantity when used as loss functions for classification problems. **Really Good** stuff.

Log Likelihood >> SEE WHAT'S INSIDE Log Likelihood

Logistic loss refers to the loss function commonly used to optimize a logistic regression model.

It may also be referred to as logarithmic loss (which is confusing) or simply log loss.

Many models are optimized under a probabilistic framework called the **maximum likelihood estimation**, or MLE, that involves finding a set of parameters that best explain the observed data.

This involves selecting a **likelihood function** that defines how likely a set of observations (data) are given model parameters. When a log likelihood function is used (which is common), it is often referred to as optimizing the log likelihood for the model. Because it is more common to minimize a function than to

Start Machine Learning

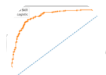
maximize it in practice, the log likelihood function is inverted by adding a negative sign to the front. This transforms it into a Negative Log Likelihood function or NLL for short.



In deriving the log likelihood function under a framework of maximum likelihood estimation for Bernoulli probability distribution functions (two classes), the calculation comes out to be:

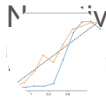
Picked for you:

- negative log-likelihood(P, Q) = $-(P(\text{class0}) * \log(Q(\text{class0})) + P(\text{class1}) * \log(Q(\text{class1})))$



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

[Recall Curves for Classification in Python](#)
Recall can be averaged over all training examples by calculating the average of the log of the likelihood function.



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

[Classification Model with scikit-learn](#)
Log-likelihood for binary classification problems is often shortened to simply “log loss” as the loss derived for logistic regression.

- log loss = negative log-likelihood, under a Bernoulli

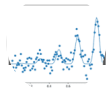


[How to Calculate the KL Divergence for](#)

[Machine Learning](#)
The negative log-likelihood is the same

for Bernoulli probability distribution functions (two events)

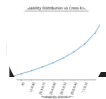
Multinoulli distributions (multi-class classification) also



[How to Implement Bayesian Optimization](#)

[from Scratch in Python](#)
The log loss and the negative log likelihood, s

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



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

[Machine Learning](#)
Log Loss and Cross Entropy Calculate the Same Thing

For classification problems, “log loss”, “cross-entropy” and “negative log-likelihood” are used interchangeably.

Loving the Tutorials?

More generally, the terms “cross-entropy” and “negative log-likelihood” are used interchangeably in the context of loss functions for classification models.

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

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



The negative log-likelihood for logistic regression is given by [...] This is also called the cross-entropy

>> SEE WHAT'S INSIDE

— Page 246, [Machine Learning: A Probabilistic Perspective](#), 2012.

Therefore, calculating log loss will give the same quantity as calculating the cross-entropy for Bernoulli probability distribution. We can confirm this by calculating the log loss using the `log_loss()` function from the scikit-learn API.

Calculating the average log loss on the same set of actual and predicted probabilities from the previous section should give the same result as calculating the average cross-entropy.

The complete example is listed below.

```
1 # calculate log loss for classification problem
```

Start Machine Learning

```

2 from sklearn.metrics import log_loss
3 from numpy import asarray
4 # define classification data
5 p = [1, 1, 1, 1, 1, 0, 0, 0, 0, 0]
6 q = [0.8, 0.9, 0.9, 0.6, 0.8, 0.1, 0.4, 0.2, 0.1, 0.3]
7 # define data as expected, e.g. probability for each event {0, 1}
8 y_true = asarray([[1-v, v] for v in p])
9 y_pred = asarray([[1-v, v] for v in q])
10 # calculate the average log loss
11 ll = log_loss(y_true, y_pred)
12 print('Average Log Loss: %.3f' % ll)

```

Recall Curves for Classification in Python

Running the example gives the expected result of 0.247 log loss, which matches 0.247 nats when calculated using the average cross-entropy.

How and When to Use a Calibrated
1 Average Log Loss: 0.247
Classification Model with scikit-learn

This does not mean that log loss calculates cross-entropy

How to Calculate the KL Divergence for
Machine Learning
they are different quantities, arrived at from different perspectives. When using a loss function for a classification task, recall that a cross-entropy loss function is equivalent to a maximum likelihood estimation of the probability distribution.

How to Implement Bayesian Optimization
from Scratch in Python
It demonstrates a connection between the study of probability theory for discrete probability distributions.

A Gentle Introduction to Cross-Entropy for
Machine Learning
limited to discrete probability distributions, and it's the first time.

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

Specifically, a linear regression optimized under the maximum likelihood estimation framework assumes a Gaussian continuous probability distribution for the target variable and involves minimizing the mean squared error function. This is equivalent to the cross-entropy for a random variable with a Gaussian probability distribution.

Loving the Tutorials?

The Probability for Machine Learning EBook is

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

Any loss consisting of a negative log-likelihood is a cross-entropy between the empirical distribution defined by the training set and the probability distribution defined by model. For example, the cross-entropy between the empirical distribution and a Gaussian model.

— Page 132, Deep Learning, 2016.

This is a little mind blowing, and comes from the field of differential entropy for continuous random variables.

It means that if you calculate the mean squared error between two Gaussian random variables that cover the same events (have the same mean and standard deviation), then you are calculating the cross-entropy between the variables.

Start Machine Learning

It also means that if you are using mean squared error loss to optimize your neural network model for a regression problem, you are in effect using a cross entropy loss.



Further Reading

Picked for you:

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



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

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



How and When to Use a Calibrated Logistic Regression With Maximum Likelihood Estimation

How to Choose Loss Functions When Training Deep Learning Models

- Loss and Loss Functions for Training Deep Learning Models



How to Calculate the KL Divergence for Machine Learning

- Machine Learning: A Probabilistic Perspective, 2006
- Pattern Recognition and Machine Learning, 2006



How to Implement Bayesian Optimization from Scratch in Python

API



A Gentle Introduction to Cross-Entropy for Machine Learning

learn.metrics.log_loss API.

Articles

- Kullback-Leibler divergence, Wikipedia.
- Cross entropy, Wikipedia.
- Joint Entropy, Wikipedia.
- The Probability for Machine Learning EBook is where you'll find the **Really Good** stuff.
- Loss functions for classification, Wikipedia.
- Differential entropy, Wikipedia.

>> SEE WHAT'S INSIDE

Summary

In this tutorial, you discovered cross-entropy for machine learning.

Specifically, you learned:

- How to calculate cross-entropy from scratch and using standard machine learning libraries.
- Cross-entropy can be used as a loss function when optimizing classification models like logistic regression and artificial neural networks.
- Cross-entropy is different from KL divergence but can be calculated using KL divergence, and is different from log loss but calculates the same quantity when used as a loss function.

Start Machine Learning

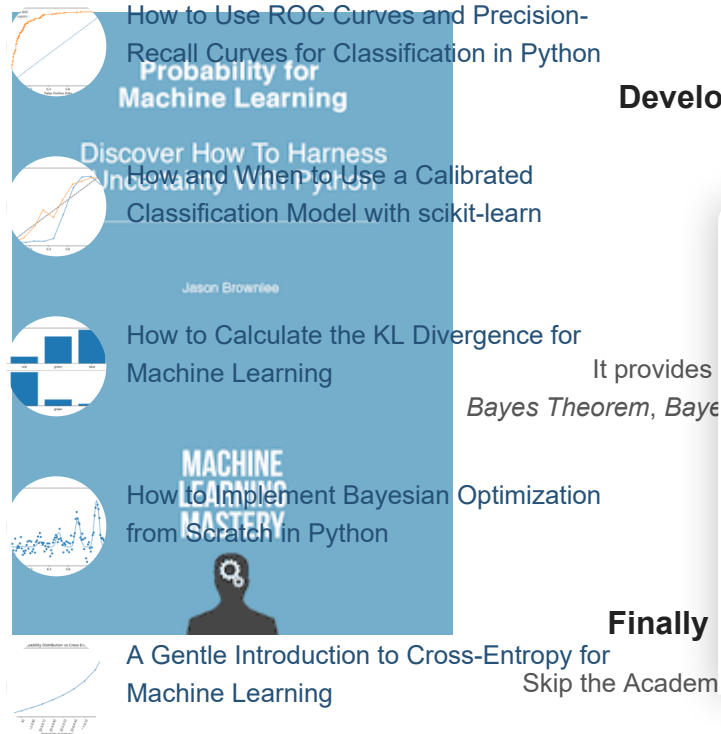
Do you have any questions?

Never miss a tutorial:

Ask your questions in the comments below and I will do my best to answer.



Picked for you: **Get a Handle on Probability for Machine Learning!**



Develop Your Understanding of Probability

...with just a few lines of python code

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

SEE WHAT'S INSIDE

Loving the Tutorials?

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

>> SEE WHAT'S INSIDE **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 →](#)

< Naive Bayes Classifier From Scratch in Python

A Gentle Introduction to Maximum Likelihood Estimation for Machine Learning >

49 Responses to *A Gentle Introduction to*

Start Machine Learning

Never miss a tutorial:**Markus** October 22, 2019 at 6:57 am #

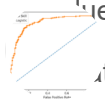
REPLY ↩



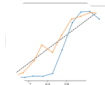
Thanks for this blog post.

Picked for you:

What confuses me a bit is the fact that we interpret the labels 0 and 1 in the example as the probability values for calculating the cross entropy between the target distribution and the predicted distribution!

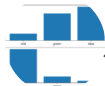
[How to Use ROC Curves and Precision](#)[Recall Curves for Classification in Python](#)

What if the labels were 4 and 7 instead of 0 and 1?!

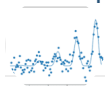
[How and When to Use a Calibrated](#)[Classification Model with scikit-learn](#)**Jason Brownlee**

October 22, 2019 at 1:44 pm #

Thanks.

[How to Calculate the KL Divergence for](#)[Scikit-Learn](#) For binary classification we map the

I recommend reading about the Bernoulli distribution

<https://machinelearningmastery.com/discrete-probability-distributions/>[How to Implement Bayesian Optimization](#)[from Scratch in Python](#)**sharpblade4**

October 22, 2019 at 5:35 pm #

A Gentle Introduction to Cross-Entropy for Machine Learning

Hi Jason,

A small fix suggestion: in the beginning of the article in section "What Is Cross-Entropy?" you've mentioned that "The result will be a positive number measured in bits and 0 if the two probability distributions are identical."

However, if the two probability distributions $H(P, P)$ is the entropy for the probability-distribution $H(P)$, opposed to KL divergence of the same probability-distribution which would indeed outcome zero.

The [Probability for Machine Learning](#) EBook iswhere you'll find the **Really Good** stuff.Thanks,
Ron U

>> 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

Jason Brownlee

October 23, 2019 at 6:34 am #

REPLY ↩

Thanks Ron!

I'll schedule time to update the post and give an example of exactly what you're referring to. E.g.:

```
1 # example of calculating cross entropy
2 from math import log2
3
4 # calculate cross entropy
5 def cross_entropy(p, q):
6     return -sum([p[i]*log2(q[i]) for i
```

Start Machine Learning

Never



Picked

```

7
8 # define data
9 p = [0.10, 0.40, 0.50]
10 q = [0.80, 0.15, 0.05]
11 # calculate cross entropy H(P, P)
12 ce_pp = cross_entropy(p, p)
13 print('H(P, P): %.3f bits' % ce_pp)
14 # calculate cross entropy H(Q, Q)
15 ce_qq = cross_entropy(q, q)
16 print('H(Q, Q): %.3f bits' % ce_qq)

```



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

H(P, P): 1.361 bits

H(Q, Q): 0.884 bits



How and When to Use a Calibrated

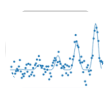
Classification Model with scikit-learn

Thanks Again!



How to Calculate the KL Divergence for Machine Learning

Hugh Brown October 26, 2019 at 4:39 pm #



How to Implement Bayesian Optimization
Why not use zip?
from Scratch in Python

```

1 # calculate cross entropy
2 def cross_entropy(p, q):
3     return -sum(pp*log2(qq) for pp, qq in zip(p, q))

```

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?

Jason Brownlee October 27, 2019 at 5:43 am #

REPLY ↩

The Probability for Machine Learning EBook is
Thanks for the tip Hugh, that is a much cleaner approach!
where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Anthony The Koala October 28, 2019 at 3:48 am #

REPLY ↩

Dear Dr Jason,

In the last few lines under the subheading "How to Calculate Cross-Entropy", you had the simple example with the following outputs:

```

1 H(P): 1.361 bits          #entropy
2 KL(P || Q): 1.927 bits   #kl divergence
3 H(P, Q): 3.288 bits      #cross entropy = entropy + kl divergence

```

What is the interpretation of these figures in 'plain English' please.

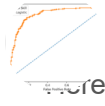
For example if the above example produced the follow

Start Machine Learning

Ne 1 $H(P)$: 0.361 bits #entropy
 2 $KL(P || Q)$: 2.927 bits #kl divergence
 3 $H(P, Q)$: 3.288 bits #cross entropy = entropy + kl divergence



Pic 1 $H(P)$: 1.927 bits #entropy
 2 $KL(P || Q)$: 0.361 bits #kl divergence
 3 $H(P, Q)$: 3.288 bits #cross entropy = entropy + kl divergence

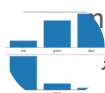


How to Use ROC Curves and Precision-

Recall Curves for Classification in Python

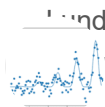
There is another example of made up figures.

1 $H(P) = 0.05$ bits #entropy
 2 $KL(P || Q) = 0.2$ bits #kl divergence
 3 $H(P, Q) = 0.25$ bits #cross entropy = entropy + kl divergence



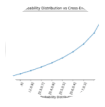
Comparing the first set of 'made up figures' does not make sense. The KL divergence is a measure of the difference between two probability distributions, not a measure of the quality of a model's predictions.

Also:



Understand that a bit is a base 2 number. Eg $1 = 1$ (base 2) number of bits in a base 2 system is an integer. We can't have a fraction of a bit.

Anthony of Sydney



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

Jason Brownlee October 28, 2019 at 6:10 am #

REPLY ↩

Typically we use cross-entropy to evaluate a model, e.g. true classes vs probability predictions.

Loving the Tutorials?

In that case would compare the average cross-entropy calculated across all examples and a lower value would represent a better fit. Interpreting the specific figures is often not useful.

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

>> SEE WHAT'S INSIDE [Jala](#) October 28, 2019 at 6:18 am #

REPLY ↩

Dear Dr Jason,
 Thank you for response.
 The other question please:
 How can you have a fraction of a bit. For example entropy = 3.2285 bits. What is 0.2285 bits.
 Thank you,
 Anthony of Sydney

Jason Brownlee October 28, 2019 at 7:00 am #

REPLY ↩

Recall, it is an average over a dis

Start Machine Learning

Think of it more of a measure and less like the crisp bits in a computer.
Never miss a tutorial:



Farhan October 28, 2019 at 6:32 pm #

REPLY ↩

Picked for you:

Great Article, Hope to see more more content on machine learning and AI.



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

Jason Brownlee October 29, 2019 at 5:21 am #

REPLY ↩



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

Thanks!



How to Calculate the KL Divergence for Machine Learning

Hi,



How to Implement Bayesian Optimization from Scratch in Python

Could you provide an example of this sentence "The errors of these values will equal 0.0."?

Does this mean a distribution with a mixture of these values?



A Gentle Introduction to Cross-Entropy for Machine Learning

Thank you!

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?

Jason Brownlee October 31, 2019 at 1:46 pm #

REPLY ↩

The Probability for Machine Learning Ebook is where you'll find the **Really Good** stuff. Sorry that is confusing.

I'm >> SEE WHAT'S INSIDE distribution for a class label will always be zero.

I have updated the tutorial to be clearer and given a worked example.

Xin November 13, 2019 at 9:35 am #

REPLY ↩

Hi Jason! This is the best article I've ever seen on cross entropy and KL-divergence! Finally I can understand them 😊 Thank you so much for the comprehensive article.

I have one small question: in the section "Intuition for Cross-Entropy on Predicted Probabilities", in the first code block to plot the visualization, the code is as follows:

Start Machine Learning

```
# define the target distribution for two events
Never miss a tutorial:
target = [0.0, 0.1]
```

```
# define probabilities for the first event
probs = [0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0]
```

```
# create probability distributions for the two events
```

Picked for you: `p] for p in probs]`

```
# calculate cross-entropy for each distribution
```

 = **How to Use ROC Curves and Precision**

Recall Curves for Classification in Python

...not quite understand why the target probability for the two events are $[0.0, 0.1]$? Could you explain a bit more? Thank you!

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

Jason Brownlee November 13, 2019 at 1:45

How to Calculate the KL Divergence for Machine Learning

Yes, looks like a typo. I'll fix it ASAP. It should be [C]

Update: I have updated the code and re-generated

A Gentle Introduction to Cross-Entropy for Machine Learning

Jason, I so appreciate all your various posts on ML topics. If I may add one comment regarding the found helpful in the past:

One point that I didn't see really emphasized here that I've seen in other treatments (e.g.,

<https://tdhopper.com/blog/cross-entropy-and-kl-divergence/>) is that cross-entropy and KL difference “differ by a constant”, i.e. in your expression

$H(P, Q) = H(P) + KL(P \parallel Q)$

the H(P) is constant with respect to Q. In most ML tasks, P is usually fixed as the “true” distribution” and Q is the dist “ ” trying to refine until it matches P.

>> SEE WHAT'S INSIDE

"In many of these situations, p is treated as the 'true' distribution, and q as the model that we're trying to optimize.... Because p is fixed, $H(p)$ doesn't change with the parameters of the model, and can be disregarded in the loss function." (<https://stats.stackexchange.com/questions/265966/why-do-we-use-kullback-leibler-divergence-rather-than-cross-entropy-in-the-t-sne/265989>)

You do get to this when you say “As such, minimizing the KL divergence and the cross entropy for a classification task are identical.”

And yet for me at least, knowing that the two “differ by a constant” makes it intuitively obvious why minimizing one is the same as minimizing the other, even if they’re actually intended to measure different things.

...Thus I think this “differ by a constant” is another reason that people get mixed up about cross-entropy vs KL divergence, and why guides like yours are so helpful

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?

$H(P, Q) = H(P) + KL(P \parallel Q)$

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

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

the $H(P)$ is constant with respect to Q. In most ML tasks, P is usually fixed as the “true” distribution” and Q is the dist “ ” “ ” “ ” trying to refine until it matches P.

>> SEE WHAT'S INSIDE

"In many of these situations, p is treated as the 'true' distribution, and q as the model that we're trying to optimize.... Because p is fixed, $H(p)$ doesn't change with the parameters of the model, and can be disregarded in the loss function." (<https://stats.stackexchange.com/questions/265966/why-do-we-use-kullback-leibler-divergence-rather-than-cross-entropy-in-the-t-sne/265989>)

You do get to this when you say “As such, minimizing the KL divergence and the cross entropy for a classification task are identical.”

And yet for me at least, knowing that the two “differ by a constant” makes it intuitively obvious why minimizing one is the same as minimizing the other, even if they’re actually intended to measure different things.

...Thus I think this “differ by a constant” is another reason that people get mixed up about cross-entropy vs KL divergence, and why guides like yours are so helpful

Never miss a tutorial:



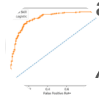
December 15, 2019 at 6:01 am #

REPLY ↩

Thanks for your note Scott.

Picked for you:

Yes, $H(P)$ is the entropy of the distribution. This becomes 0 when class labels are 0 and 1. I outline this at the end of the post when we talk about class labels.



[How to Use ROC Curves and Precision-Recall Curves for Classification in Python](#)

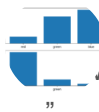
[Recall Curves for Classification in Python](#)

A constant of 0 in that case means using KL divergence and cross entropy result in the same numbers, e.g. the KL divergence.



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

Arij February 11, 2020 at 6:52 am #



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

thanks for a grate article!
what confused me that in your article you have me

The result will be a positive number measured in bits a



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

probability distributions are identical."

then again mentioned



"If two probability distributions are the same, then the c

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

then

"This means that the cross entropy of two distributions (real and predicted) that have the same probability distribution for a class label, will also always be 0.0."

Loving the Tutorials?

"Therefore, a cross-entropy of 0.0 when training a model indicates that the predicted class probabilities are identical to the probabilities in the training dataset, e.g. zero loss."

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

Am I missing something?

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

>> SEE WHAT'S INSIDE

Jason Brownlee February 11, 2020 at 1:42 pm #

REPLY ↩

The statements are correct.

The cross-entropy will be the entropy between the distributions if the distributions are identical.

In this case, if we are working with class labels like 0 and 1, then the entropy for two identical distributions will be zero.

We demonstrate this with a worked example in the above tutorial.

Does that help?

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!

April 11, 2020 at 5:36 pm #

REPLY ↩



It was not a bad example, they were understandable, thanks.

Just I could not imagine and understand them numerically.

Reading them again I understand that when the values of any distribution are only one or zero then entropy,

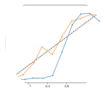
Picked for you:

cross-entropy, KL all will be zero.



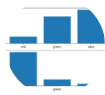
How to Use ROC Curves and Precision-Recall Curves for Classification in Python

Recall Curves for Classification in Python



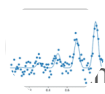
How and Why to Use a Calibrated Classification Model with scikit-learn

Exactly!



How to Calculate the KL Divergence for Machine Learning

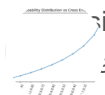
Grzegorz Kępisty April 17, 2020 at 4:23 pm #



How to Implement Bayesian Optimization from Scratch in Python

Thank you for great post!

Question on KL Divergence: In its definition we have log



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?

Jason Brownlee April 18, 2020 at 5:41 am #

REPLY ↩

Good question, no problem as probabilities are always greater than zero, so log never blows

The Probability for Machine Learning EBook is

where you'll find the **Really Good** stuff.
More on KL divergence here too:

<http://www.jbrownlee.com/divergence-between-probability-distributions/>

>> SEE WHAT'S INSIDE

Allan April 27, 2020 at 2:25 am #

REPLY ↩

Is it possible to use KL divergence as a classification criterion?

Jason Brownlee April 27, 2020 at 5:37 am #

REPLY ↩

Probably, it would be the same as log loss and cross entropy when using class labels instead of probabilities.

Start Machine Learning

Never miss a tutorial:



All April 28, 2020 at 12:46 pm #

REPLY ↩

Thanks for your reply. So let say the final calculation result is "Average Log Loss", what

Picked for you:



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

Jason Brownlee April 28, 2020 at 6:39 am #

REPLY ↩



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



How to Calculate the KL Divergence for Machine Learning

Gledson April 29, 2020 at 6:21 am #

Hello Jason, Congratulations on the explanation on my crossentropy. Should I replace it with some



How to implement Bayesian Optimization from Scratch in Python

Gledson.



A Gentle Introduction to Cross-Entropy for Machine Learning

Jason Brownlee April 29, 2020 at 6:36 am #

REPLY ↩

You cannot log a zero. It is a good idea to always add a tiny value to anything to log, e.g.

$\log(\text{value} + 1e-8)$

Loving the Tutorials?

The **Probability for Machine Learning** EBook is

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

Zeinab May 13, 2020 at 3:47 am #

REPLY ↩

>> SEE WHAT'S INSIDE

Thanks for all your great post, I've read some of them.

I'm working on traffic classification and I've converted my data to string of bits, I want to use cross-entropy on bytes.

Assume below lists:

$p = [1, 0, 1, 1, 0, 0, 1, 0]$

$q = [1, 1, 1, 0, 1, 0, 0, 1]$

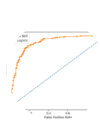
When I use `-sum([p[i] * log2(q[i]) for i in range(len(p))])`, I encounter this error :ValueError: math domain error. S

Start Machine Learning

Would you please tell me what I'm doing wrong here and how can I implement cross-entropy on a list of bits?



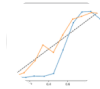
Picked for you:



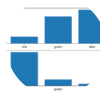
Jason Brownlee May 13, 2020 at 6:43 am #
How to Use ROC Curves and Precision-Recall Curves for Classification in Python

REPLY ↩

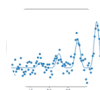
Can't calculate log of 0.0. Try adding a tiny value to the equation, e.g. $1e-8$ or $1e-15$



How and When to Use a Calibrated Classification Model with scikit-learn
Zeinab May 13, 2020 at 11:03 pm #

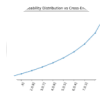


I've converted the traffic to string of bits any value.
How to Calculate the KL Divergence for Machine Learning



How to Implement Bayesian Optimization from Scratch in Python
TuanVu May 13, 2020 at 7:58 pm #

Why we use log function for cross entropy?



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

Jason Brownlee May 14, 2020 at 5:46 am #

REPLY ↩

Good question, perhaps start here:
<https://machinelearningmastery.com/what-is-information-entropy/>

Loving the Tutorials?

The Probability for Machine Learning EBook is where you'll find the **Really Good** stuff.
TuanVu May 14, 2020 at 9:22 am #

REPLY ↩

>> SEE WHAT'S INSIDE

Jason Brownlee May 14, 2020 at 1:25 pm #

REPLY ↩

You're welcome.

Zahir May 29, 2020 at 3:36 pm #

REPLY ↩

Hello Jason,
Thank you so much for all your great posts.

Start Machine Learning

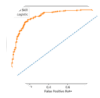
I have a question, if we have conditional entropy $H(y|x) = -\sum P(x,y) \log(P(y|x))$ could we say that it is equal to cross-entropy $H(x,y) = -\sum y \log y^x$?



And if that correct when we could say that? i.e., under what assumptions.

Thank you in advance.

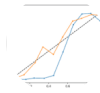
Picked for you:



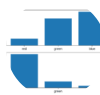
How to Use ROC Curves and Precision-Recall Curves for Classification in Python
Jason Brownlee May 30, 2020 at 5:52 am #

REPLY ↩

I think you're asking me if the conditional entropy is the same as the cross entropy. I don't think it is off the cuff, but perhaps confirm with a good textbook.



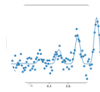
How and When to Use a Calibrated Classification Model with scikit-learn



How to Calculate the KL Divergence for Machine Learning
Zahir May 30, 2020 at 10:23 am #

Thank you so much for your replay,

I found it in "Privacy-Preserving Adversarial Neural Networks" as a cost function, but when they implement the why? How to Implement Bayesian Optimization from Scratch in Python



A Gentle Introduction to Cross-Entropy for Machine Learning
Jason Brownlee May 31, 2020 at 6



Fascinating.

Perhaps email the authors directly.

Loving the Tutorials?

The Probability of Machine Learning
BARURI SAI AVINASH November 1, 2020 at 4:22 pm #

REPLY ↩

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

How can be Number of bits per charecter in text generation is equal to loss ???

>> SEE WHAT'S INSIDE

Jason Brownlee November 2, 2020 at 6:38 am #

REPLY ↩

Perhaps try re-reading the above tutorial that lays it all out.

Eric Orr December 22, 2020 at 11:05 am #

REPLY ↩

This is excellent Introduction to Cross-Entropy. It seems that one of the following sentences may have a typo in the stated notion of "surprise". My first impression is that the second sentence should have said "are less surprising". Is that true?

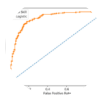
Start Machine Learning

Never miss a tutorial: “Low probability events are more surprising therefore have a larger amount of information. Whereas probability distributions where the events are equally likely are more surprising and have larger entropy.”



Picked for you: [Jason Brownlee](#) December 22, 2020 at 1:34 pm #

REPLY ↩



Thanks
[How to Use ROC Curves and Precision-Recall Curves for Classification in Python](#)

yes it could be clearer. Information is about events, entropy is about distributions, cross-entropy is about comparing distributions.



Surprise means something different when talking about information/events as compared to [Classification Mode](#) with scikit-learn

I mixed the discussion of the two at the start of the



Also see this: <https://machinelearningmastery.com/what-is-information-entropy/>

Does that help?



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

Eric Orr December 22, 2020 at 11:48 am #



[A Gentle Introduction to Cross-Entropy for Machine Learning](#)
After additional consideration, it appears that

OWS.

“In probability distributions where the events are equally likely, no events have larger or smaller likelihood (smaller or larger surprise, respectively), and the distribution has larger entropy.”

Loving the Tutorials?

Sorry for belaboring this. It is a good point but sometimes confusing.

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

stefani January 6, 2021 at 11:19 pm #

REPLY ↩

>> SEE WHAT'S INSIDE

I understood a concept when you can describe it with very simple words and I feel that is the case here. Wonderful job!

Jason Brownlee January 7, 2021 at 6:18 am #

REPLY ↩

Thanks!

Aaron February 20, 2021 at 8:52 am #

REPLY ↩

Start Machine Learning

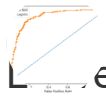
“Relative Entropy (KL Divergence): Average number of extra bits to represent an event from Q instead of P.”
Never miss a tutorial:

Shouldn't it rather say: Relative Entropy (KL Divergence): Average number of extra bits to represent an event
 coming from Q instead of P.

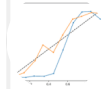


Since the expectation is over $P(x)[\dots]$

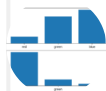
Picked for you:



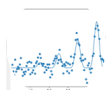
How to Use ROC Curves and Precision-Recall Curves for Classification in Python



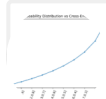
How and When to Use a Calibrated Classification Model with scikit-learn



How to Calculate the KL Divergence for Machine Learning



How to Implement Bayesian Optimization from Scratch in Python



A Gentle Introduction to Cross-Entropy for Machine Learning

Website

Start Machine Learning



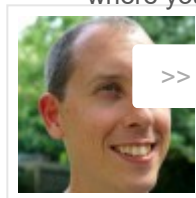
You can master applied Machine Learning **without math or fancy degrees**.
 Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Loving the Tutorials?

SUBMIT COMMENT

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



Welcome!

>> SEE WHAT'S INSIDE

Read more

hD
 s get results with **machine learning**.

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 Calculate the KL Divergence for Machine Learning



How to Implement Bayesian Optimization from Scratch in Python



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?

© 2021 Machine Learning Mastery Pty Ltd. All Rights Reserved.
The Probability for Machine Learning EBook is
When you'll find the **Really Good** stuff
LinkedIn | Twitter | Facebook | Newsletter | RSS

Privacy | D >> SEE WHAT'S INSIDE itemap | Search

Start Machine Learning