



New? Start here!

Machine Learning course

Subscribe to our newsletter

About

March 25, 2014 · MACHINE LEARNING

Simple guide to confusion matrix terminology

A confusion matrix is a table that is often used to **describe the performance of a classification model** (or "classifier") on a set of test data for which the true values are known. The confusion matrix itself is relatively simple to understand, but the related terminology can be confusing.

I wanted to create a "quick reference guide" for confusion matrix terminology because I couldn't find an existing resource that suited my requirements: compact in presentation, using numbers instead of arbitrary variables, and explained both in terms of formulas and sentences.

Let's start with an **example confusion matrix for a binary classifier** (though it can easily be extended to the case of more than two classes):

	Predicted:	Predicted:
n=165	NO	YES
Actual:		
NO	50	10
Actual:		
YES	5	100

What can we learn from this matrix?

There are two possible predicted classes:
 "yes" and "no". If we were predicting the
 presence of a disease, for example, "yes"
 would mean they have the disease, and





New? Start here!

Machine Learning course

Subscribe to our newsletter

About

- "no" would mean they don't have the disease.
- The classifier made a total of 165 predictions (e.g., 165 patients were being tested for the presence of that disease).
- Out of those 165 cases, the classifier predicted "yes" 110 times, and "no" 55 times.
- In reality, 105 patients in the sample have the disease, and 60 patients do not.

Let's now define the most basic terms, which are whole numbers (not rates):

- **true positives (TP):** These are cases in which we predicted yes (they have the disease), and they do have the disease.
- **true negatives (TN):** We predicted no, and they don't have the disease.
- **false positives (FP):** We predicted yes, but they don't actually have the disease. (Also known as a "Type I error.")
- false negatives (FN): We predicted no, but they actually do have the disease. (Also known as a "Type II error.")

I've added these terms to the confusion matrix, and also added the row and column totals:

n=165	Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60
Actual: YES	FN = 5	TP = 100	105
	55	110	

This is a list of rates that are often computed from a confusion matrix for a binary classifier:





New? Start here!

Machine Learning course

Subscribe to our newsletter

About

- **Accuracy:** Overall, how often is the classifier correct?
 - \circ (TP+TN)/total = (100+50)/165 = 0.91
- **Misclassification Rate:** Overall, how often is it wrong?
 - \circ (FP+FN)/total = (10+5)/165 = 0.09
 - equivalent to 1 minus Accuracy
 - o also known as "Error Rate"
- **True Positive Rate:** When it's actually yes, how often does it predict yes?
 - TP/actual yes = 100/105 = 0.95
 - also known as "Sensitivity" or "Recall"
- **False Positive Rate:** When it's actually no, how often does it predict yes?
 - \circ FP/actual no = 10/60 = 0.17
- **True Negative Rate:** When it's actually no, how often does it predict no?
 - TN/actual no = 50/60 = 0.83
 - equivalent to 1 minus False Positive
 Rate
 - also known as "Specificity"
- **Precision:** When it predicts yes, how often is it correct?
 - TP/predicted yes = 100/110 = 0.91
- **Prevalence:** How often does the yes condition actually occur in our sample?
 - o actual yes/total = 105/165 = 0.64

A couple other terms are also worth mentioning:

• Null Error Rate: This is how often you would be wrong if you always predicted the majority class. (In our example, the null error rate would be 60/165=0.36 because if you always predicted yes, you would only be wrong for the 60 "no" cases.) This can be a useful baseline metric to compare your classifier against. However, the best classifier for a particular application will sometimes have a higher error rate than





New? Start here!

Machine Learning course

Subscribe to our newsletter

About

- the null error rate, as demonstrated by the **Accuracy Paradox**.
- Cohen's Kappa: This is essentially a measure of how well the classifier performed as compared to how well it would have performed simply by chance. In other words, a model will have a high Kappa score if there is a big difference between the accuracy and the null error rate. (More details about Cohen's Kappa.)
- **F Score:** This is a weighted average of the true positive rate (recall) and precision.

 (**More details about the F Score.**)
- ROC Curve: This is a commonly used graph that summarizes the performance of a classifier over all possible thresholds. It is generated by plotting the True Positive Rate (y-axis) against the False Positive Rate (x-axis) as you vary the threshold for assigning observations to a given class.

 (More details about ROC Curves.)

And finally, for those of you from the world of Bayesian statistics, here's a quick summary of these terms from **Applied Predictive Modeling**:

In relation to Bayesian statistics, the sensitivity and specificity are the conditional probabilities, the prevalence is the prior, and the positive/negative predicted values are the posterior probabilities.

Want to learn more?

In my new 35-minute video, **Making sense of the confusion matrix**, I explain these concepts in more depth and cover more **advanced topics**:





New? Start here!

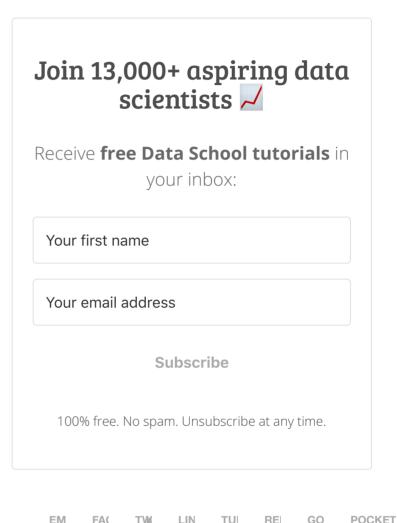
Machine Learning course

Subscribe to our newsletter

About

- How to calculate precision and recall for multi-class problems
- How to analyze a 10-class confusion matrix
- How to choose the right evaluation metric for your problem
- Why accuracy is often a misleading metric

Let me know if you have any questions!



Data School Comment Policy

All comments are moderated, and will usually be ap Kevin within a few hours. Thanks for your patience!



Join the discussion...





New? Start here!

Machine Learning course

Subscribe to our newsletter

About

LOG IN WITH

OR SIGN UP WITH DISQUS (?)

Name



Engr Ali Raza · 3 years ago

Dear Kevin.

can you please tell me the relationship betwe Misclassifcations and split value or split value explain that how we can calculate the split va decision algorithm?

SplitValue/Set Cases Misclassifications Class

22.65 5 0.8 1.6 Decision

N/A 1 0 1.6 Terminal

8.63 4 0.6 2.7 Decision

14.4 2 0.2 2.7 Decision

34.85 2 0.2 4 Decision

N/A 1 0 2.7 Terminal

N/A 1 0 5.2 Terminal

N/A 1 0 4 Terminal

N/A 1 0 7.5 Terminal

can anyone explain this table?how it ca be ge formulas used behind these values

78 ^ V · Reply · Share ›



Jenica J. Wilson → Engr Ali Raza • 2 years ε I know that this is 10 months ago but have not received an answer here it is

When developing a measure or test yo identify a unique group of people - for Rosenberg Self-Esteem Test purports people with low self-esteem. Good me particularly those used in research or users would like to see them have gre abilities. One way to do this is to cross splitting the data up is way to cross-va percentage of the data that is used to measure identifies people with, let say esteem, is then taken or split from the subset data is set buy you. So, if you well the measure identifies future exar test this on 30% of your data. So your would be set to .70, because 70% of be used as in the analysis to identify t low self-esteem. Based on the patterr from the 70% of the data set, the sens (what you have above) would use the





New? Start here!

Machine Learning course

Subscribe to our newsletter

About

to test how accurately the measure re participants back into their groups.

sensitivity = true positive/false positive (accuracy rates)

specificity = true negatives/ false negatives (accuracy rates)

Prevalence = proportion of group/tota Detection rates= measure ability to ac various groups

I am not sure what the program was w received the output above. So maybe would help people best identify how to findings and provide you on how the i be generated.



Kevin Markham Mod → Engr Ali Raza • 3 ye I'm sorry, I'm not familiar with split val table you included in your comment. (



ajay kumar • 5 years ago

Respected Sir,

I have two confusion matrices and I want to pure McNemar Test. It is hereby to requesting you how to generate the values of 2 by 2 matrix I find the values of f11, f12, f21 and f22 from the matrices.

Thank You,



Peter → ajay kumar • 10 months ago

dear friend, try confusionMatrix functipackage



Kevin Markham Mod → ajay kumar • 5 year I'm not familiar with McNemar's test.



Todd de Quincey • 2 years ago

Great summary! Thanks for putting this toget going into my notebook of tips and tricks



Kevin Markham Mod → Todd de Quincey •

You're welcome!







New? Start here!

Machine Learning course

Subscribe to our newsletter

About



Davide • 3 years ago

Dear Kevin,

Thanks for the clear explanation.

I have a particular problem and I'm struggling hope your expertise can help me.

I want to compare the performance of a burn toward a reference one. My analysis is not pix not creating random points i the map and cre matrix. This also because the proportion of to vs not burned is really small. I'm therefore co detected by the product vs reference fires, ar matrix using the whole available data and not samples. Also the error matrix is missing True Concretely:

- reference dataset has 70 fires
- tested product has 36
- 21 fires are correctly detected by the tested
- 15 fires detected by the tested product are
- 49 reference fires are not detected by the te (FN)
- I didn't add TN

Several questions rise:

- can I use the available derived metrics to quesults? (F1, sensitivity and precision) even if RANDOM observations but the WHOLE AVAI OBSERVATIONS in the datasets? Or should I normal rates of detection and omission?
- should I add TN observations for the calcula ?(how many?) they will also be not random...
- if I add an other burned area product which fires, the error matrix will have a different nun observations... Are they comparable?
- Is there other statistic tools which would su case?

Thank you for your attention

5 ^ V · Reply · Share ›



Kevin Markham Mod → Davide • 3 years ag

Thanks for your detailed question! I'd help, but I'm not able to give you advi without having a lot more domain kno example, I don't know what a "burnec is, how you compare it against a "refe "pixel-based" analysis might look like, is the kind of question that I would impose the with a student for 20 minutes, and only





New? Start here!

Machine Learning course

Subscribe to our newsletter

About

give good advice! Without more inform be irresponisble of me to try to answe I'm sorry, and good luck!



Priyansh Agarwal · 3 years ago

Hi Kevin

I am confused with Type1 and TypeII errors

false positives (FP): We predicted yes, but the have the disease. (Also known as a "Type I er false negatives (FN): We predicted no, but the have the disease. (Also known as a "Type II e

Type 1 error is Worst than Type 2 error and bu Type 2 error looks worst than Type 1, becaus when we have a disease is more bad that the

Thank



Kevin Markham Mod → Priyansh Agarwal •

Thanks for your question. It appears to saying that in general, type I errors are type II errors. That's not actually the cois "worse" depends on the particular salso depends on your perspective. Ho



manan • 8 months ago

hi there.

i m working with election predictions. i have a model using 2018 election data-set and test 2008 election data-set. now my question is the mean of all confusion matrix for three electing single model.



Kevin Markham Mod → manan · 8 months : I don't think I understand your goals v provide any advice for this scenario -



NLR • 8 months ago

Hi Kevin, this is extremely helpful. Do you have number for your reference to Applied Prediction regard to getting Bayesian calculations from matrix?





New? Start here!

Machine Learning course

Subscribe to our newsletter

About



Kevin Markham Mod → NLR • 8 months age I don't have a page number, I'm sorry!

2 ^ | ∨ • Reply • Share •



Hritik Singh • 2 years ago

How to calculate the various metrics like pred with 3 classes (eg- the iris dataset)



Kevin Markham Mod → Hritik Singh • 2 years ago • edited

Great question! To calculate per-class iris dataset, for example, you are answ questions: (1) When it predicts setosa correct? (2) When it predicts versicolo correct? (3) When it predicts virginica, correct?

To answer question 1, for example, the "number of setosa predictions" and the "how many of those were correct". Yo that in order to answer questions 2 and

Similarly, to calculate per-class recall, answer questions like: (1) When the trasetosa, how often does it predict seto

You can see a simple example of 3-clarecall here: http://scikit-learn.org/sta..

Hope that helps!



Jenica J. Wilson → Hritik Singh • 2 years ag if you are using R here is the formula: This will give you the actual accuracy classifier, here it is LDA, reclassified th

p1 <- predict(lda,dataset)\$class
tab <- table(Predict=p1,Actual=datase
of accurate predictions/confusion mat
accuracy <- sum(diag(tab))/sum(tab)#a
model
tab

accuracy

If you want to cross-validate in R see Packages that you would need install.packages("caret") library("caret") library(klaR)





New? Start here!

Machine Learning course

Subscribe to our newsletter

About

You can choose your split

see more



hana • 2 years ago

I work on credit risk project the confusion ma TN=6 FP 10 AND FN=3 I calculate some Met , error sensitivity and other also I suggest the money from confusion matrix my QU is I wan for the calculation of earn and lose money on thanks



Kevin Markham Mod → hana • 2 years ago I'm sorry, I don't completely understar Good luck!



Tamil Selvan · 3 years ago · edited

can you please clarify me what is the differen misclassification error and misclassification random misclassification ran



Kevin Markham Mod → Tamil Selvan • 3 year I wouldn't recommend using the term "misclassification error". A "classification incorrect, and a "misclassification" is a whereas "misclassification error" is a or

"Misclassification rate", on the other had percentage of classifications that were 3 ^ | V · Reply · Share >



Raj Kandala • 4 years ago

Dear Kevin

Sub: In multi class confusion matrix, Finding each class.

According to accuracy definition, Acc= (TP+TN)/(TP+TN+FP+FN). accuracy is high e⁻ TP=0,. How can we judge the performance ir the same time precsion= TP/(TP+FP) is zero. explain this situation.



Kevin Markham Mod → Raj Kandala • 4 yea





New? Start here!

Machine Learning course

Subscribe to our newsletter

About



In a multi-class problem (meaning mo classes), the accuracy is simply "corredivided by "total predictions".

If you want to evaluate each class ind option is to calculate the per-class preclass recall. scikit-learn's classification this, for example: http://scikit-learn.or



Raj Kandala → Kevin Markham • 4 y Thank you kevin



Raj Kandala • 4 years ago

Hi, I want to calculate ROC plots using multic matrix. Is it possible Mr. Kevin?



Kevin Markham Mod → Raj Kandala • 4 yea ROC curves can only be drawn for bir problems, meaning problems with two classes. (However, you can turn a mul into a binary problem using a "one ver approach.)

Also, drawing an ROC curve requires the predicted probability of class men observation, rather than just the class Therefore, a confusion matrix alone (e case) does not provide enough data in draw the ROC curve.

This post provides more information a curves: http://www.dataschool.io/ro...

Hope that helps!



Raj Kandala → Kevin Markham • 4 y
Thank You Kevin



Abdullah Nazzal · 4 years ago · edited

what if we have more than two classes ... say how to calculate this .. ? should i convert it to if so then how ?



Kevin Markham Mod → Abdullah Nazzal • 4



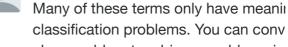


New? Start here!

Machine Learning course

Subscribe to our newsletter

About



classification problems. You can conv class problem to a binary problem sir output classes together, though I wou doing that just so that you can use a r to evaluate your model. Rather, you sh convert it from multi-class to binary if sense in the context of your problem.



Hassan Aftab Mughal → Kevin Mai 4 years ago

I think if you have more than tv your transfusion matrix will be n=165 Class1 Class2 Class3 C

Class1 10 3 3 0

Class2 5 15 1 2

Class3 0 5 20 1

Class4 0 2 1 9

2 ^ | V · Reply · Share ›



Kevin Markham Mod A H. Mughal • 4 years ago

That's correct, you can confusion matrix for a n problem. (My apologies otherwise!) However, m terminology outlined in applies to binary classif problems.

Here's an example of a matrix: http://scikit-lear

As a side note, the "n=" upper left corner of the refers to the number of so in the case of your 4 matrix, it would be n=77

© 2020 Data School. All rights reserved. Privacy policy.