and the problem we are trying to solve. When someone tells you that "USA is the best country", the first question that you should ask is on what basis is this statement being made. Are we judging each country on the basis of their economic status, or their health facilities etc.? Similarly each machine learning model is trying to solve a problem with a different objective using a different dataset and hence, it is important to understand the context before choosing a metric.

discussing it and how it is even related to Data Science. Well, in this post, I will be discussing the usefulness of each error metric depending on the objective

Most Useful Metrics

Classification Regression MSPE Precision-Recall o ROC-AUC MSAE R Square Accuracy Adjusted R Square o Log-Loss

RMSE = $\sqrt{\frac{1}{n} \sum_{j=1}^{n} (y_j - \hat{y}_j)^2}$

whereas RMSE penalizes the higher difference more than MAE.

Let's understand the above statement with the two examples:

Case 1: Actual Values = [2,4,6,8] , Predicted Values = [4,6,8,10]

Case 2: Actual Values = [2,4,6,8] , Predicted Values = [4,6,8,12]

MAE for case 1 = 2.0, RMSE for case 1 = 2.0

MAE for case 2 = 2.5, RMSE for case 2 = 2.65

Let's take a simple linear model in one variable: y = mx+b

values (called residuals). Mathematically, it is calculated using this formula:

MAE MAE is the average of the absolute difference between the predicted values and observed value. The

calculated using this formula:

 $MAE = \frac{1}{n} \sum_{j=1}^{n} |y_j - \hat{y}_j|$ So which one should you choose and why? Well, it is easy to understand and interpret MAE because it directly takes the average of offsets

MAE is a linear score which means that all the individual differences are weighted equally in the

average. For example, the difference between 10 and 0 will be twice the difference between 5 and 0.

However, same is not true for RMSE which we will discuss more in details further. Mathematically, it is

If we define loss function (J) in terms of RMSE: then we can easily differentiate J wrt. to m and b and get

 $\frac{\partial}{\partial \mathbf{m}} = \frac{2}{N} \sum_{i=1}^{N} -x_i (y_i - (mx_i + b))$

From the above example, we can see that RMSE penalizes the last value prediction more heavily than

MAE. Generally, RMSE will be higher than or equal to MAE. The only case where it equals MAE is when

all the differences are equal or zero (true for case 1 where the difference between actual and

 $\frac{\partial}{\partial \mathbf{b}} = \frac{2}{N} \sum_{i=1}^{N} -(y_i - (mx_i + b))$

https://spin.atomicobject.com/wp-content/uploads/linear_regression_gradient1.png

R Squared (R²) and Adjusted R Squared

Mathematically, R_Squared is given by:

terms in a model. It is given by below formula:

Why should you choose Adjusted R² over R²?

Var1

x1

x2

their pros and cons.

Adjusted R²

be less than or equal to R²

Case 1

y1

y2

yhat = m.predict(X)

 $SS_Residual = sum((y-yhat)**2)$

return r_squared,adj_r_squared

model1 = linear_model.LinearRegression()

model1.fit(data.drop("x2", axis = 1),y)

metrics(model1,data.drop("x2", axis=1),y)

model2 = linear_model.LinearRegression()

model3 = linear_model.LinearRegression()

y = np.array([2.1, 4, 6.2, 8, 9])

y = np.array([2.1, 4, 6.2, 8, 9])

 $SS_Total = sum((y-np.mean(y))**2)$

r_squared = 1 - (float(SS_Residual))/SS_Total

print(yhat)

model2.fit(data,y)

model3.fit(data,y)

R_squared

possibility still exists.

Adj_R_squared

model 2 for more number of variables)

Comparison of Adjusted R² over RMSE

a model has adjusted R² equal to 0.05 then it is definitely poor.

metrics(model2,data,y)

8

9

10

11

12

13

14

16

17

19

20

22

23

25

26

with an example.

Var1

x1

x2

The above equations are simpler to solve and the same won't apply for MAE.

Here, we are trying to find "m" and "b" and we are provided with data (x,y).

the updated m and b (this is how gradient descent works, I won't be explaining it here)

Edit: One important distinction between MAE & RMSE that I forgot to mention earlier is that minimizing the squared error over a set of numbers results in finding its mean, and minimizing the absolute error results in finding its median. This is the reason why MAE is robust to outliers whereas RMSE is not. This answer explains this concept in detail.

R Squared & Adjusted R Squared are often used for explanatory purposes and explains how well your

metrics are quite misunderstood and therefore I would like to clarify them first before going through

selected independent variable(s) explain the variability in your dependent variable(s). Both these

Just like R², adjusted R² also shows how well terms fit a curve or line but adjusts for the number of

 $R_{adj}^2 = 1 - \left[\frac{(1-R^2)(n-1)}{n-k-1} \right]$

where n is the total number of observations and k is the number of predictors. Adjusted R² will always

There are some problems with normal R² which are solved by Adjusted R². An adjusted R² will consider

the marginal improvement added by an additional term in your model. So it will increase if you add

increasing terms even though the model is not actually improving. It will be easier to understand this

у1

y2

Case 3

2*x1+0.1

y1

y2

Var2

2*x2

Var1

x1

x2

the useful terms and it will decrease if you add less useful predictors. However, R² increases with

Case 2

Var2

2*x1

2*x2

slight disturbance in var2 such that it is no longer perfectly correlated with var1.

as y values which is not true for R Square. The range of RMSE & MAE is from 0 to infinity.

However if you want a metric just to compare between two models from interpretation point of view,

then I think MAE is a better choice. It is important to note that the units of both RMSE & MAE are same

$$\hat{R}^2 = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2}$$
The numerator is MSE (average of the squares of the residuals) and the denominator is the variance in Y values. Higher the MSE, smaller the R_squared and poorer is the model.

хЗ y3 хЗ 2*x3 у3 хЗ y3 2*x3 + 0.12*x4 y4 y4 x4 2*x4 х4 х4 y4 2*x5 + 0.1x5 x5 2*x5 y5

Here, Case 1 is the simple case where we have 5 observations of (x,y). In case 2, we have one more

variable which is twice of variable 1 (perfectly correlated with var 1). In Case 3, we have produced a

So if we fit simple ordinary least squares (OLS) model for each case, logically we are not providing any

extra or useful information to case 2 and case 3 with respect to case 1. So our metric value should not

improve for these models. However, it is actually not true for R² which gives a higher value for model 2

and 3. But your adjusted R2 takes care of this problem and it is actually decreasing for both cases 2 &

3. Let's give some numbers to these variables (x,y) and look at the results obtained in Python. import numpy as np import pandas as pd from sklearn import datasets, linear_model def metrics(m,X,y):

 $adj_r_squared = 1 - (1-r_squared)*(len(y)-1)/(len(y)-X.shape[1]-1)$

data = $pd.DataFrame(\{"x1": [1,2,3,4,5], "x2": [2.1,4,6.1,8,10.1]\})$

data = $pd.DataFrame({"x1": [1,2,3,4,5], "x2": [2.1,4,6.1,8,10.1]})$

Case 1

metrics(model3,data,y) Linear_Regression_Python hosted with ♥ by GitHub **Note:** Predicted values will be same for both model 1 and model 2 and therefore, r_squared will also be same because it depends only on predicted and actual values.

0.985

0.981

case 1 to case 2, still R2 is increasing whereas adjusted R2 is showing the correct trend (penalizing

From the above table, we can see that even though we are not adding any additional information from

For the previous example, we will see that RMSE is same for case 1 and case 2 similar to R². This is the

values with actual values. Also, the absolute value of RMSE does not actually tell how bad a model is. It

can only be used to compare across two models whereas Adjusted R² easily does that. For example, if

However, if you care only about prediction accuracy then RMSE is best. It is computationally

Common Misconception: I have often seen on the web that the range of R² lies between 0 and 1

y_actual is positive. In this case, R² will be less than 0. This will be a highly unlikely scenario but the

Consider the case where model is predicting highly negative value for all the observations even though

which is not actually true. The maximum value of R² is 1 but minimum can be negative infinity.

simple, easily differentiable and present as default metric for most of the models.

case where Adjusted R² does a better job than RMSE whose scope is limited to comparing predicted

Case 2

0.985

0.971

Here is an interesting metric to know about if you are interested in NLP and you deserve it for reaching the end. I recently got to know about it in Andrew Ng Deep Learning course and found it worth sharing.
BLEU (Bilingual Evaluation Understudy)
It is mostly used to measure the quality of machine translation with respect to the human translation. It uses a modified form of precision metric.
Steps to compute BLEU score:
1. Convert the sentence into unigrams, bigrams, trigrams, and 4-grams
2. Compute precision for n-grams of size 1 to 4

 $BLEU = BP * \exp\left(\sum w_n \log\left(P_n\right)\right)$

3. Take the exponential of the weighted average of all those precision values

4. Multiply it with brevity penalty (will explain later)

Let's compare the above two translations with BLEU score.

1 on the

1 mat is

1 is a

0.83

0.25

0 a cat

P2

1 the mat

Unigram Match | Bigram

Ion

!the

mat

is

ıa

!cat

Weights

Ρ1

Reference: The cat is sitting on the mat Machine Translation 1: On the mat is a cat Machine Translation 2: There is cat sitting cat

Match Trigram

Here I am using nltk.translate.bleu_score package from nltk.translate.bleu_score import sentence_bleu reference = [['the', 'cat',"is","sitting","on","the","mat"]] candidate = ["on",'the',"mat","is","a","cat"] score = sentence_bleu(reference, candidate) print(score)

0.40

0.25

Why do we add brevity penalty? Brevity Penalty penalizes candidates shorter than their reference translations. For example, if the candidate for the above mentioned reference is "The cat", then it will have a high precision for unigram and bigram because, both the words are present in the reference in the same order. However, the length is too short and does not actually reflect the meaning of reference. With this brevity penalty in place, a high-scoring candidate translation must now match the reference in terms of length, same words and order of words. Hopefully, most of the useful measures in regression are covered in this blog. Please comment if

Related: CatBoost vs. Light GBM vs. XGBoost • Learning Curves for Machine Learning

Bio: Alvira Swalin (Medium) is currently pursuing Master's in Data Science at USF, I am particularly

interested in Machine Learning & Predictive Modeling. She is a Data Science Intern at Price (Fx).

Join the discussion... **LOG IN WITH** OR SIGN UP WITH DISQUS (?) DfyG Name

not the best indicator of what the model will do. I almost always tie model accuracy to \$\$ or counts of

trades off errors in the way that best matches the business objective? :)

DhananJay Shembekar → Dean Abbott • 2 years ago wow.true that, but the author aims something different, its when one identifies the problem. **Subscribe** Add Disgus to your site DISQUS **⚠** Do Not Sell My Data <= Previous post

1. Data Science Minimum: 10 Essential Skills You Need to Know to Start Skills You Need to Know to Start **Doing Data Science Doing Data Science** 2. Introduction to Time Series Analysis 2. Introduction to Time Series Analysis in Python in Python

5. Autograd: The Best Machine **Learning Library Youre Not Using?** 2020 Update 6. How I Consistently Improve My 6. Deep Learning's Most Important

in 2020

KDnuggets Home » News » 2018 » Apr » Tutorials, Overviews » Choosing the Right Metric for Evaluating Machine Learning Models – Part 1 (18:n18)

WHAT DATA SCIENTISTS SPEND THE MOST TIME DOING



- **Analysis in Python**
- **Learning Course From Amazon** 3. How LinkedIn Uses Machine **Learning in its Recruiter Recommendation Systems**

4. A step-by-step guide for

- for Android 7. A step-by-step guide for creating an authentic data science portfolio project
- Omission: Business Leadership fastcore: An Underrated Python Library

Top tweets, Oct 7-13: Every

DataFrame Manipulation,

Free From MIT: Intro to

Deep Learning Design Patterns

Computational Thinking and Data

Explain...

Science Goodhart's Law for Data Science and what happens when a meas... Getting Started with PyTorch

Prod...

2. Data Science Minimum: 10 **Essential Skills You Need to Know to Start Doing Data** Science 3. 10 Best Machine Learning Courses in 2020

4. Strategies of Docker Images

- **Optimization** 5. The Best Free Data Science eBooks: 2020 Update 6. How LinkedIn Uses Machine **Learning in its Recruiter Recommendation Systems**
- 7. Introduction to Time Series
- **Most Shared** 1. 10 Best Machine Learning Courses in 2020
- 2. Free Introductory Machine
- - coding challenge Text Mining with R: The Free eBook

KDnuggets 20:n39, Oct 14: A stepby-step guide for creating...

How to be a 10x data scientist Top Stories, Oct 5-11: A step-bystep guide for creating an au...

creating an authentic data science portfolio project 5. Annotated Machine Learning **Research Papers** 6. A Guide to Preparing OpenCV **More Recent Stories** Machine Learning's Greatest How to ace the data science

in

+

98

SHARES

Top September Stories: Free From MIT: Intro to Computer Scienc...

SIAM launches activity group,

The Future of Fake News

publications for data scientists

Software Engineering Tips and

Best Practices for Data Science

Uber Open Sources the Third

5 Best Practices for Putting

Machine Learning Models Into

M...

Release of Ludwig, its Code-Free

Exploring The Brute Force K-Nearest Neighbors Algorithm **Annotated Machine Learning** Research Papers

view raw

0.987

0.975

Case 3

Bonus!

Here BP is the brevity penalty, r & c is the number of words in reference & candidate respectively, w weights, P—Precision values **Example:**

1 on the mat

1 the mat is

0 mat is a

0 is a cat

Match 4-gram

0

0.25

0.25

Р3

1 on the mat is

0 the mat is a

0 mat is a cat

Match

0.00

0.25

Ρ4

anything is missed, will really appreciate ideas for the next blog. I will discuss the important metrics of classification and unsupervised learning in the next blog. Stay tuned! In the meantime, check out my other blogs here! LinkedIn: www.linkedin.com/in/alvira-swalin References 1. http://thestatsgeek.com/2013/10/28/r-squared-and-adjusted-r-squared/ 2. https://www.aclweb.org/anthology/P02-1040.pdf

3. https://www.nltk.org/_modules/nltk/translate/bleu_score.html

Original. Reposted with permission.

Machine Learning Model Metrics

2 ^ | V · Reply · Share ›

3. Automating Every Aspect of Your

Python Project

KDnuggets ■ Disqus' Privacy Policy **1** Login 2 Comments Sort by Best -Commend 2 f Share **Tweet**

Dean Abbott • 2 years ago Choosing the right metric is a HUGE issue. I've never used any of the Regression metrics above to describe to stakeholders how good a regression model is (these were for customer analytics and fraud detection problems). Why not? Because (1) business stakeholders don't have intuition about what they actually mean and (2) most of these kinds of problems include a "selection" part where one contacts or audits a subset of all records scored. Since R-squared, MSE, etc. are all metrics that give an average influence for all records, and if we aren't contacting most of the population anyway, global metrics are

individuals/tax returns/cases the model will correctly identify or product a positive \$\$ return. In fact, I've

seen cases where the Spearman rank correlation for say 1000 models is close to 0 when I compare MSE

and something like Lift at the 3rd decile. Models tradeoff errors differently. Why not pick the one that

Next post => **Top Stories Past 30 Days Most Popular Most Shared** 1. Data Science Minimum: 10 Essential

4. Machine Learning from Scratch: 4. Machine Learning from Scratch: **Free Online Textbook Free Online Textbook** 5. The Best Free Data Science eBooks: **Machine Learning Models From 80% Ideas**

3. 10 Best Machine Learning Courses

to Over 90% Accuracy 7. Online Certificates/Courses in Al, **Data Science, Machine Learning** 7. The Best Free Data Science eBooks: 2020 Update from Top Universities

Subscribe to KDnuggets News

© 2020 KDnuggets. | About KDnuggets | Contact | Privacy policy | Terms of Service