

ABSTRACT/PURPOSE

When working with a regression problem an issue that often arises is having an abundance of data that hinders our model. A dataset may include a large set of variables, but some of the variables do not contain much of the necessary information. One of the solutions to this problem is implementing an algorithm called Principal Component Analysis, this algorithm reduces the dimensions of our data set. Principal Component Analysis helps us in visualizing the main component in our data set and improves our modeling results. We are using data from a Kaggle competition that includes a training and testing data set, with 24 variables included in each data set. Our goal is to take an extensive look at the process of Principal Component Analysis and how it performs. We will do this by comparing the modelling results when using Principal Component Analysis and when not using Principal Component Analysis on the Kaggle dataset.

RESULTS

Multilinear Linear Regression

- We wanted to create a Multiple Linear Regression Model to later compare the results to our Principal Component Analysis Model.
- The image on the right shows our linear model using our 24 variables. This model was obtained using the testing dataset.
- We also calculated the R², Root Mean Square Error, and mean absolute error using our testing dataset.
- The higher the R² value means, that our model accounts for a higher percentage of our response variable variance.
- The lower the RMSE and MAE means that the model was performing better.

```
Residuals:
    Min       1Q   Median       3Q      Max
-8.1514 -1.8066 -0.3543  1.0614  9.4182

Coefficients:
(Intercept)      -2.184022    1.348481   -1.6200   0.10709
minutes_past     -0.005712    0.012995   -0.4400   0.66080
radardist_low    0.234396    0.077924    3.0080   0.00301
Ref              0.095947    0.084211    1.1390   0.23609
Ref_5x5_10th    0.032102    0.065657    0.4890   0.62549
Ref_5x5_30th    -0.234997    0.141548   -2.3607   0.01903
Ref_5x5_90th    0.212841    0.104995    2.0270   0.04415
RefComposite    -0.110030    0.111219   -1.3459   0.17911
RefComposite_5x5_10th -0.306883    0.094275   -3.2550   0.00136
RefComposite_5x5_30th  0.428391    0.214755    1.9950   0.04760
RefComposite_5x5_90th  0.106871    0.130885    0.8170   0.41530
RhoHV           -0.846304    1.483706   -0.5700   0.56965
RhoHV_5x5_10th -1.065291    1.339369   -0.7980   0.42573
RhoHV_5x5_30th  1.865122    1.689013    1.1040   0.27098
RhoHV_5x5_90th -0.528266    1.034144   -0.5110   0.61011
Zdr            -0.123846    0.310959   -0.3980   0.69094
Zdr_5x5_10th   0.159655    0.502030    0.3180   0.75085
Zdr_5x5_30th  0.276974    0.743168    0.3730   0.70983
Zdr_5x5_90th  0.515309    0.281191   -1.8330   0.06854
Kdp            0.022515    0.068759    0.3270   0.74372
Kdp_5x5_10th  0.125414    0.112032    1.1190   0.26446
Kdp_5x5_30th  0.110888    0.122942    0.9020   0.36831
Kdp_5x5_90th  -0.122505    0.105149   -1.1650   0.24556
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.936 on 177 degrees of freedom
Multiple R-squared:  0.3791, Adjusted R-squared:  0.302
F-statistic: 4.913 on 22 and 177 DF, p-value: 5.281e-10
```

R2	RMSE	MAE
0.2174564	2.356512	1.923406

CONCLUSIONS/IMPLICATIONS

- Why did our PCA model result in a lower R² and a higher RMSE and a higher MAE? Does this necessarily mean our PCA model was worse?

We were expecting the PCA model to perform slightly worse than the multilinear regression model because, models using PCA sacrifice some accuracy for a simpler data set. The simpler data set is easier to work with and ensures the data that is less significant to not be included, when creating the model.

- Why would you choose to implement PCA?

In today's world, there are huge datasets. Principal Component Analysis is able to simplify and compress those datasets into something simpler which still captures most of the information from the original dataset. This allows for faster computations within and upon those datasets.

METHODOLOGY

Problem Approach

- Our main goal was to reduce the amount of dimensions of our dataset using Principal Component Analysis. This allows for us to retain most of the datasets variance and information while making a large dataset smaller and therefore easier to work with. We then see how the Principal Component Analysis model holds up to a Multiple Linear Regression Model.

Data

- First we used a dataset that included 24 variables in order to predict "the actual gauge observation in millimeters" for rain (Kaggle).
- This dataset had upwards of 1,000,000 entries. A lot of these rows had a vast amount of missing data so we ended up using only 252 rows to do our models upon.

Multiple Linear Regression

- First, we split our revised dataset into a training and test dataset.
- Then we ran a multiple linear regression which did not incorporate any Principal Component Analysis and used all variables in our dataset.
- We used our model to get Predicted values for the expected rainfall and compared it to the actual values of the rainfall.
- We observed the R² and RMSE for this model in order to see how it compared to the Principal Component Analysis Model.

Principal Component Analysis

- Like in Multiple Linear Regression, we split our revised dataset into a training and test dataset using the same seed we set in Multiple Linear Regression.
- We then ran Principal Component Analysis on our 24 variables in order to see how many principal components we want in our new model.
- When we ran a summary of after running Principal Component Analysis, it gave us a table that included Validation:RMSEP and TRAINING: % variance explained.
- From Validation:RMSEP it calculates the RMSE as you add more components. The Lowest RMSE happens at 15 components.
- From Training Percent Variance, it tells us the variance in our response variable (expected) which is explained by the amount of principal components. At all 22 components you reach 100 percent, but at 15 you already reached 98.68 percent. Therefore based off of Validation:RMSEP and TRAINING: %, we deduced it was best to run 15 components onto our new model.
- We then ran the model based off of 15 principal components and recorded R², RMSE and MAE.

RESULTS

PCA Model

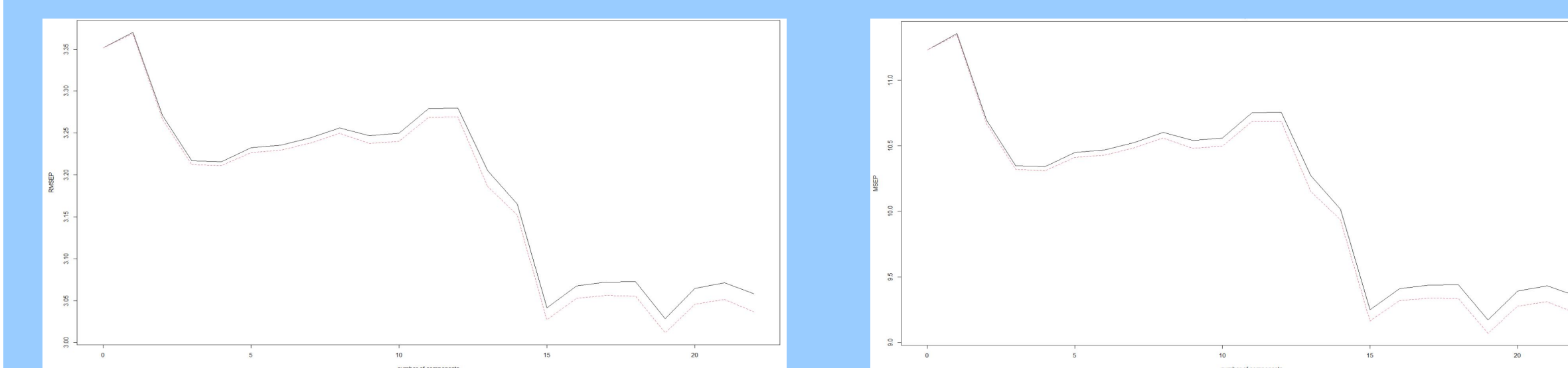
- To create our model using PCA, we first needed to calculate our principle components.
- The image on the right shows the principle components calculated by PCA and this was information we used to decide, which components to include in our model.

```
Data: X dimension: 24 22
Fit method: svd
Number of components considered: 22

VALIDATION: RMSEP
Cross-validated using 10 random segments.
CV (Intercept) 1 comps 2 comps 3 comps 4 comps 5 comps 6 comps 7 comps 8 comps 9 comps 10 comps 11 comps 12 comps 13 comps 14 comps 15 comps 16 comps 17 comps 18 comps 19 comps 20 comps 21 comps 22 comps
RMSECV  3.151  3.130  3.125  3.121  3.118  3.115  3.112  3.110  3.109  3.108  3.107  3.106  3.105  3.104  3.103  3.102  3.101  3.100  3.099  3.098  3.097  3.096  3.095
MAECV   3.250  3.238  3.234  3.229  3.226  3.223  3.220  3.218  3.216  3.215  3.214  3.213  3.212  3.211  3.210  3.209  3.208  3.207  3.206  3.205  3.204  3.203  3.202

TRAINING: % variance explained
% var  1 comps 2 comps 3 comps 4 comps 5 comps 6 comps 7 comps 8 comps 9 comps 10 comps 11 comps 12 comps 13 comps 14 comps 15 comps 16 comps 17 comps 18 comps 19 comps 20 comps 21 comps 22 comps
Expected 36.822  54.450  65.14  72.09  76.81  81.22  84.66  87.92  90.08  91.09  91.66  92.00  92.22  92.36  92.44  92.50  92.54  92.57  92.59  92.60  92.61  92.62  92.63  92.64
X      30.45  82.67  94.61  96.81  98.06  99.10  99.85  99.98  99.99  99.99  99.99  99.99  99.99  99.99  99.99  99.99  99.99  99.99  99.99  99.99  99.99  99.99  99.99
MAE    30.48  82.44  94.40  96.66  97.82  98.75  99.37  99.57  99.67  99.73  99.76  99.78  99.79  99.80  99.81  99.81  99.82  99.82  99.82  99.82  99.82  99.82  99.82
RMSE   30.48  82.44  94.40  96.66  97.82  98.75  99.37  99.57  99.67  99.73  99.76  99.78  99.79  99.80  99.81  99.81  99.82  99.82  99.82  99.82  99.82  99.82  99.82
MAE    30.48  82.44  94.40  96.66  97.82  98.75  99.37  99.57  99.67  99.73  99.76  99.78  99.79  99.80  99.81  99.81  99.82  99.82  99.82  99.82  99.82  99.82  99.82
RMSE   30.48  82.44  94.40  96.66  97.82  98.75  99.37  99.57  99.67  99.73  99.76  99.78  99.79  99.80  99.81  99.81  99.82  99.82  99.82  99.82  99.82  99.82  99.82
MAE    30.48  82.44  94.40  96.66  97.82  98.75  99.37  99.57  99.67  99.73  99.76  99.78  99.79  99.80  99.81  99.81  99.82  99.82  99.82  99.82  99.82  99.82  99.82
```

To decide which components to include in our model, we looked at the RMSEP and Training % Variance explained in our model when using the training data. We wanted to find the number of components which we can include in our model that would not give us diminishing returns,



Below are the metrics in which we calculate how well our model performed, when using 15 components in our model.

R2	RMSE	MAE
0.1966326	2.453501	1.992071

Our PCA model resulted in a lower R² value, higher RMSE, and higher MAE value.

BIBLIOGRAPHY

- Ph.D., Benjamin Obi Tayo. "Machine Learning: Dimensionality Reduction via Principal Component Analysis." Medium, Towards AI, 11 June 2019.
<https://pub.towardsai.net/machine-learning-dimensionality-reduction-via-principal-component-analysis-1bdc77462831>.
- Jaadi, Zakaria. "A Step-by-Step Explanation of Principal Component Analysis (PCA)." Built In, 1 Apr. 2021.
<https://builtin.com/data-science/step-by-step-explanation-principal-component-analysis>.
- Zach. "Principal Components Analysis in R: Step-by-Step Example." Statology, 1 Dec. 2020.
<https://www.statology.org/principal-components-analysis-in-r/>
- "How Much Did It Rain? II." Kaggle, 7 Dec. 2015.
<https://www.kaggle.com/competitions/how-much-did-it-rain-ii/data>.

ACKNOWLEDGEMENTS

Thank you Dr. Ivona Grzegorzcyk for the guidance and encouragement that you have given us throughout this project.
Thank you Kaggle for providing the dataset.