
CS 365 Project Final

Jae Hong Lee and Wyatt Napier
Boston University
jhonglee@bu.edu wnapier@bu.edu

Abstract

1 For our project milestone we're using Principal-Component Analysis (PCA) and
2 Support Vector Machines (SVM) to predict housing prices based off of this dataset
3 provided by Kaggle. This is our Google Colab notebook.

4 1 Introduction

5 1.1 Problem

6 As college students in Boston, housing prices are often on our minds, especially as we look for
7 off-campus housing. Housing price is based on a variety of factors, but we want to ensure that we are
8 getting the most value for our money. By creating a model to estimate housing prices based off of
9 the attributes of a specific house, we will more objectively be able to determine if the housing we're
10 looking at is truly worth its price. We will use PCA to do some preprocessing on our data and then
11 use support vector machines and radial basis functions for our regression model.

12 1.2 Dataset

13 The dataset that we're going to use to train our model comes from Kaggle. There are 1460 rows of
14 data which we will concatenate and do our own train/test split in order to optimize the ratio so as to
15 not over or undertrain. This dataset is incredibly verbose as it has roughly 80 different columns. We
16 can analyze the differences in results based on the selected attributes and diversify the combinations
17 of attributes chosen to increase accuracy.

18 1.3 Resources

19 In our project we used Google Colab to write and run our code. On Google Colab, we used roughly
20 28.2 GB of disk space as well as 1.4 GB of RAM. We were most frequently connected to the CPU
21 runtime since we could only use the T4 GPU for a limited amount of time. This was sufficient for
22 running all of the code besides a small portion at the end where we attempt to compare the results of
23 the model on the original pre-processed data with the results of the model on the PCA representation
24 of the original pre-processed data. In this case, the runtime disconnects before it completes. This
25 typically occurs after at least 1.5 hours of running. We will elaborate on why that occurred in the
26 results section of this paper. Some other tools that we relied on were Sci-Kit Learn documentation,
27 ChatGPT to help debug our code, and the various textbooks we cite in the references section.

28 2 Principal Component Analysis

29 Principal Component Analysis (PCA) is a method of compressing the dimension of data to make
30 learning easier and prediction more efficient. We want to go from representing the data in R^d to a
31 significantly smaller space R^p by using p principal components from the eigendecomposition.

32 We want to approximate an input vector $\mathbf{x}^{(n)} \in R^d$ with basis vectors $\mathbf{v}_1, \dots, \mathbf{v}_d \in R^d$ by fulfilling
33 the following equation where each α is a representation coefficient: $\mathbf{x}^{(n)} \approx \sum_{i=1}^p \alpha_i \mathbf{v}_i$. To do

so we need to find both the basis vectors $\{\mathbf{v}_i\}_{i=1}^p$ and coefficients $\{\alpha_i\}_{i=1}^p$. We do this using an eigendecomposition which is often derived from SVD or the definition of eigenvalues and eigenvectors: $\mathbf{A}\mathbf{u} = \lambda\mathbf{u} \rightarrow \mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1}$.

For just one row of data we would need to solve this optimization problem: $(\hat{\alpha}, \hat{\mathbf{v}}) = \underset{\|\mathbf{v}\|_2=1, \alpha}{\operatorname{argmin}} \|\mathbf{x} - \alpha\mathbf{v}\|^2$. Basically what this means is we want to find the coefficient and the basis vector such that we can minimize the distance of \mathbf{x} from the scaled vector $\alpha\mathbf{v}$. To generalize this to the entire matrix we have the following equations:

$$\hat{\mathbf{v}} = \underset{\|\mathbf{v}\|_2=1}{\operatorname{argmin}} E\|\mathbf{X} - \alpha\mathbf{v}\|^2 = \underset{\|\mathbf{v}\|_2=1}{\operatorname{argmax}} \mathbf{v}^T \mathbf{\Sigma} \mathbf{v}$$

We want to maximize the variance of the data, while also minimizing the distance of the data from the basis vectors. These two problems happen to actually be the same problem. Since in this case, $\mathbf{\Sigma}$ is just the covariance matrix from $\mathbf{\Sigma} = E[\mathbf{X}^T \mathbf{X}]$, we know that the $\hat{\mathbf{v}}$ has a solution \mathbf{u}_1 which is the first column of the eigenvector matrix \mathbf{U} assuming it is ordered by magnitude of corresponding eigenvalue.

Stepping back for a moment, we want to maximize the variance along the basis vectors which will also minimize the distance of each point from that line. Since we're using the covariance matrix the diagonal will just be the variance. When we do the eigendecomposition, the eigenvalues are then a measure of the variance so the magnitude of the eigenvalues corresponds to the variance of each.

To clarify, the eigendecomposition of the $d \times d$ matrix $\mathbf{\Sigma}$ will give us the eigenvector matrix \mathbf{U} and the eigenvalue matrix \mathbf{S} . As part of PCA we want to use just p principle components (started with d components) so we just truncate \mathbf{U} and \mathbf{S} to use only the first p eigenvectors and eigenvalues. In doing so, we project the points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ onto the space spanned by the first p eigenvectors of \mathbf{U} . We use the following equation for a single point and repeat for the entire dataset to compress it:

$$\mathbf{U}_p^T \mathbf{x}^{(n)} = \alpha^{(n)}$$

Overall, the structure of PCA is to center the mean of the data on the origin and then find the eigendecomposition of the covariance matrix so that we can then do the reduction. In some applications the data is first scaled for each feature before applying SVD. PCA's foundations lay in simple linear algebraic techniques such as SVD and eigendecomposition and as such it is a true testament to the power of linear algebra.

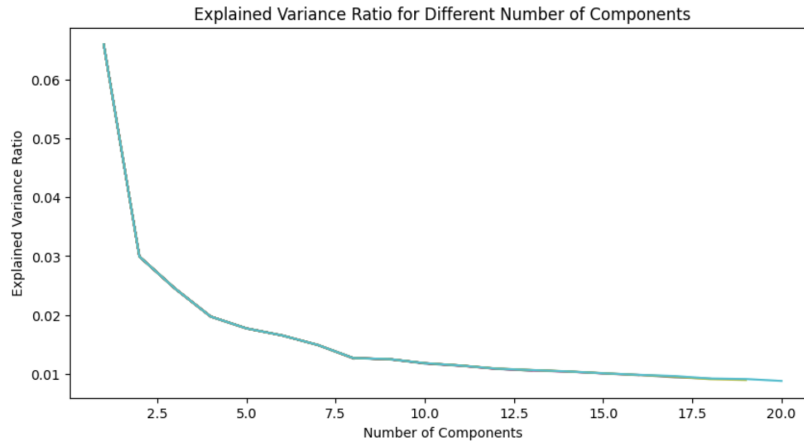


Figure 1: Explained Variance Ratio for Different Number of Components

When applying this to our dataset, we want to find the optimal number of principal components to maximize the explained variance ratio. The explained variance ratio equation is the percentage of the variance explained or captured by the selected component. In Figure 1 above, since the explained variance ratio has such a steep decay, only a small number of principal components are really needed. It is evident that the magnitude of the explained variance ratio after the 8th principal component seems to flatten out, so we chose to limit our PCA to 8 principal components to represent our preprocessed 271 columns.

58 3 Support Vector Machine

59 3.1 Decision Rule

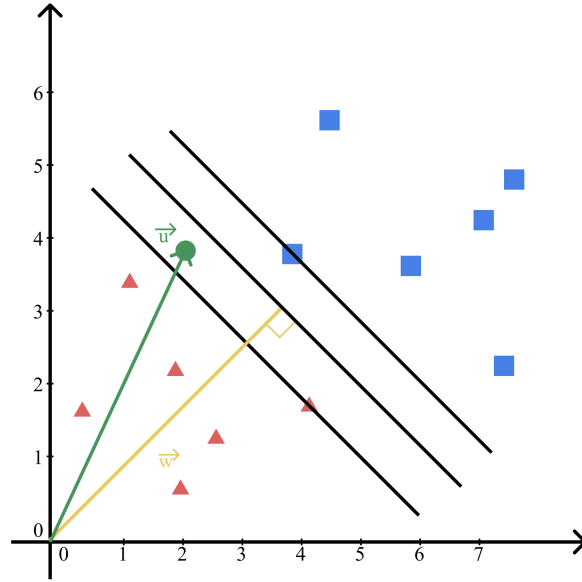


Figure 2: Decision Boundary

60 The space in the middle is drawn based on the widest path between the closest vectors (commonly
 61 called support vectors) of two different categories. Therefore, it's referred to as the widest street
 62 approach. In the method of dividing examples of two kinds: triangles and rectangles, this path is
 63 made as wide as possible. There is a vector \vec{w} of some length. This \vec{w} is perpendicular to the dividing
 64 center-line of the plane, and therefore the distance from the origin to the center-line is minimal. The
 65 distance of this vector is unknown.

66 And there is another unknown point. It is represented by \vec{u} in the figure. It is important to determine
 67 whether this point is closer to a triangle or a rectangle. Now, by projecting this point onto the \vec{w}
 68 vector, we can determine how many times it is proportional to the \vec{w} vector. If the length of the
 69 projection of \vec{u} onto \vec{w} is closer to the center-line, it is classified as a triangle; otherwise, if it is further
 70 than the center-line, it is classified as a rectangle. The figure above \vec{u} would be classified as a triangle.

71 The projection of \vec{u} to \vec{w} is the dot product of two vectors.

72 This is SVM's Decision Rule

$$\vec{w} \cdot \vec{u} \geq \text{some constant } c$$

73 if $c = -b$, $\vec{w} \cdot \vec{u} + b \geq 0$ then Rectangle

74 3.2 Margin Constraints

75 What conditions are needed for the constants b and vector \vec{w} ? We only know that vector \vec{w} is
 76 perpendicular to the central line. However, since \vec{w} can be perpendicular to the central line regardless
 77 of its length, there are no constraints on \vec{w} and b . To determine \vec{w} and b , some constraints are required.

78 Using Figure 2 as an example, to determine \vec{w} and b 's values, for rectangular vectors which is farther
 79 than the center-line, find two variables \vec{w} and b that satisfy the equation $\vec{w} \cdot \vec{x}_{rec} + b \geq 1$. The output
 80 will be greater than 1 for all rectangle vectors. In addition, the previously found \vec{w} and b should
 81 satisfy the equation $\vec{w} \cdot \vec{x}_{tri} + b \leq -1$ for triangle vectors.

82 However, finding the values of \vec{w} and b using these two equations is not an easy task. For mathematical
 83 convenience, let's define a variable y_i which returns a value of +1 for rectangle vectors and -1 for
 84 triangle vectors. Therefore, for the equation of the rectangle, $\vec{w} \cdot \vec{x}_{rec} + b \geq 1$ the expression involving

85 $y_i: y_i \cdot (\vec{w} \cdot x_{rec} + b) \geq 1$ and for the triangle: $y_i \cdot (\vec{w} \cdot x_{tri} + b) \geq 1$. They both have the same
86 equations $y_i \cdot (\vec{w} \cdot x_{rec \text{ or } tri} + b) - 1 \geq 0$. In a special case, for the vectors inside the widest street
87 area (gutter) the equation holds $y_i \cdot (\vec{w} \cdot x_{rec \text{ or } tri} + b) - 1 = 0$.

88 3.3 Get the width of the street

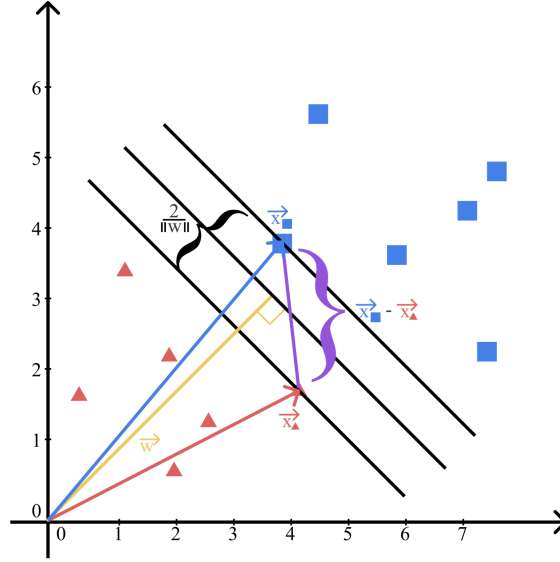


Figure 3: Width of the street

89 What is the widest path to divide triangles and rectangles? Despite not knowing the path widths, we
90 can calculate the difference between the two support vectors (the closest two vectors to the gutter) of
91 the triangles and rectangles. Using these two support vectors as x_{rec} and x_{tri} , respectively, we can
92 calculate the difference between the two vectors as $(x_{rec} - x_{tri})$. We can determine the width of the
93 street by taking the dot product of the difference between these two vectors and the unit normal to the
94 center-line. In other words, if we have a unit vector pointing in the direction of the center-line, we
95 can calculate the street width.

$$width = (x_{rec} - x_{tri}) \cdot \frac{\vec{w}}{\|\vec{w}\|}$$

96 From the previous margin constraints, we have the equation, $y_i \cdot (\vec{w} \cdot x_{rec \text{ or } tri} + b) - 1 = 0$. If we
97 have the rectangle support vector, y_{rec} will be 1 and rest of the equation will be $(\vec{w} \cdot x_{rec} + b) - 1 = 0$.
98 So $x_{rec} \cdot \vec{w}$ will be $1 - b$. Similarly, when we have the triangle support vector, y_{tri} will be -1 and
99 $x_{tri} \cdot \vec{w}$ will be $-1 - b$. With this equations we got $x_{rec} = 1 - b$ and $x_{tri} = -1 - b$. So the width is
100 $(x_{rec} - x_{tri}) \cdot \frac{\vec{w}}{\|\vec{w}\|} = \frac{x_{rec} \cdot \vec{w} - x_{tri} \cdot \vec{w}}{\|\vec{w}\|} = \frac{1 - b - (-1 - b)}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|}$.

101 We got the width: $\frac{2}{\|\vec{w}\|}$. We want to maximize the width. For mathematical convenience, it is okay
102 to maximize $\frac{1}{\|\vec{w}\|}$. Which also leads to minimize $\|\vec{w}\|$. Eventually it is minimizing $\frac{1}{2} \cdot \|\vec{w}\|^2$.

103 3.4 Maximize Width

104 From previous section, we learn in order to maximize the width of the street, need to minimize
105 $\frac{1}{2} \cdot \|\vec{w}\|^2$. To find the min $\|\vec{w}\|$ we use Lagrange multiplier.

106 Applying Lagrange Multiplier,

$$L = \frac{1}{2} \|\vec{w}\|^2 - \sum \alpha_i [y_i (\vec{w} \cdot \vec{x}_i + b) - 1]$$

107 In order to find the minimum, derivative the equation L and find w, b which makes 0.

$$\begin{aligned} \frac{dL}{dw} &= \vec{w} - \sum \alpha_i y_i \vec{x}_i \Rightarrow 0 \\ \Rightarrow \vec{w} &= \sum \alpha_i y_i \vec{x}_i \end{aligned}$$

\vec{w} is the linear sum of these vectors α_i, y_i, \vec{x}_i

$$\frac{dL}{db} = - \sum \alpha_i y_i \Rightarrow 0$$

Plug back in founded derivative to the original Lagrange Multiplier.

$$L = \frac{1}{2} \sum (\alpha_i y_i \vec{x}_i) \cdot \sum (\alpha_j y_j \vec{x}_j) - \sum (\alpha_i y_i \vec{x}_i) \sum (\alpha_j y_j \vec{x}_j) - \sum \alpha_i y_i b + \sum \alpha_i$$

we know b is constant so takes it outside of sum.

$$L = \frac{1}{2} \sum (\alpha_i y_i \vec{x}_i) \cdot \sum (\alpha_j y_j \vec{x}_j) - \sum (\alpha_i y_i \vec{x}_i) \sum (\alpha_j y_j \vec{x}_j) - b \sum \alpha_i y_i + \sum \alpha_i$$

From derivative of b , $-\sum \alpha_i y_i = 0$.

$$\begin{aligned} L &= \frac{1}{2} \sum (\alpha_i y_i \vec{x}_i) \cdot \sum (\alpha_j y_j \vec{x}_j) - \sum (\alpha_i y_i \vec{x}_i) \sum (\alpha_j y_j \vec{x}_j) + \sum \alpha_i \\ &= \sum \alpha_i - \frac{1}{2} \sum (\alpha_i y_i \vec{x}_i) \cdot \sum (\alpha_j y_j \vec{x}_j) \\ &= \sum \alpha_i - \frac{1}{2} \sum \sum a_i a_j y_i y_j \cdot x_i \cdot x_j \end{aligned}$$

We see that the optimization is only depends on the pair of samples.

3.5 Apply back to Decision Rule

Now plug w back in decision rule from the figure 1.

$$\sum (a_i y_i \vec{x}_i) \cdot \vec{u} + b \geq 0, \text{ then Rectangle}$$

4 Kernel Trick

In SVM, sometimes the data isn't linearly separable. In this case, a non-linear transformation must be applied to the data so that it can then be categorized. Then, SVM doesn't need to know the function for transformation, it just needs to know how the data points compare to each other after they're all transformed by using the kernel function which is centered around the dot product.

In the linear kernel, $f(x) = x$ and the kernel function is $k(x, x') = x^T x'$. This attempts to create a linear boundary between data points with minimal transformation, so it often isn't enough to create a proper boundary.

In the polynomial kernel, the equation of the transformation would be $f(x) = (x_1, x_2, x_1 x_2, x_1^2, x_2^2)$ and the polynomial kernel is $k(x, x') = (1 + x^T x')^2$. In this case, it is particularly useful because you don't need to know the transformation function itself. This will create a boundary that appears to be much more dynamic because the transformation that is applied is non-linear.

In the RBF (Radial Basis Function) kernel, the transformation function captures infinite dimensions so the abstraction of $f(x)$ in the kernel function is necessary to properly fit the data. Instead we just use the simplified kernel function: $k(x, x') = e^{-r \|x - x'\|^2}$.

133 5 Experiments

134 5.1 PCA Result

135 Based on our research on PCA and SVM methods, we applied Support Vector Regression (SVR) to
136 our housing price dataset, aiming to predict house prices using 80 features. The dataset contains a
137 large number of features, so we used PCA to reduce the number of feature columns to 8.

138 To apply PCA, we first scaled our data so that the eigendecomposition wasn't skewed and then we
139 used sklearn's PCA method to reduce the dimensionality of our dataset to 8 columns by projecting
140 the data onto the space spanned by these 8 eigenvectors. We then centered these results around the
141 origin by transforming the data.

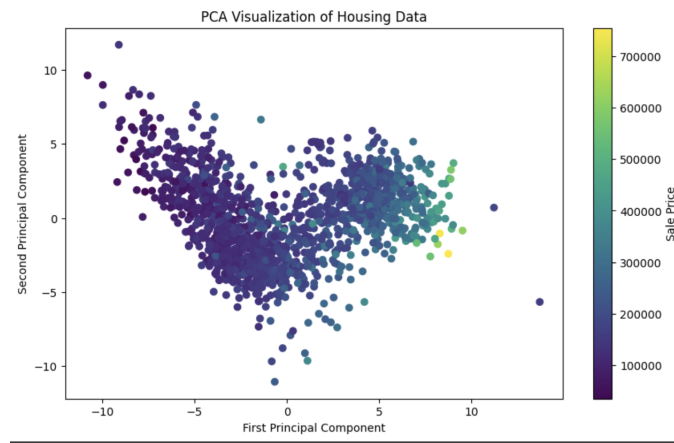


Figure 4: PCA Plot of YearBuilt and Foundation_CBlock

142 These 8 features that captured the greatest variance were, in order from greatest magnitude to
143 least: YearBuilt, Foundation_CBlock, GrLivArea, 2ndFlrSF, MSSubClass, BldgType_1Fam, Bldg-
144 Type_Duplex, and SaleCondition_Normal.

145 One thing to keep in mind is that the principal components that we choose to retain are those that
146 capture the greatest variance, not necessarily those that have the greatest impact on the dataset. These
147 components are not directly interpretable in terms of the original features but represent patterns of
148 variation in the data. Thus, we can draw very limited conclusions from the results of the PCA.

149 As mentioned in the resources section above, when we tried to run the Support Vector Regression
150 model on the original data, the runtime disconnected after roughly 1.5 hours rather than completing
151 the computation. This means the computation was incredibly demanding and had to be abandoned.
152 This demonstrates the value of PCA in increasing efficiency of training a model. Without PCA, we
153 couldn't even train the model because it was too inefficient due to the size of the data, but in the
154 reduced form, it took well under a minute to fit the model for the PCA adjusted data.

155 5.2 Support Vector Regression Model

156 The experiment dataset underwent PCA, and the eight columns with the highest variance were
157 selected. These datasets were then used to split the data into a training set and a test set in a 7:3 ratio.
158 This ensured that we could accurately train our data, but wouldn't overfit so testing would still be
159 accurate. The models were trained on this split data, and accuracy was evaluated.

160 We conducted experiments using three models: linear, RBF, and polynomial, based on the kernel
161 trick used before calculating the support vector.

162 5.2.1 Fine Tuning

163 SVR models, contain several parameters, but I will only mention the parameters we touched.

- 164 • kernel : { 'linear', 'rbf', 'poly' } The types of kernel trick we will use in our SVR Model
- 165 • C : { float, default = 1.0 (Positive) } Regularization parameter.

- degree : { default = 3 (Non-Negative) } Degree of the polynomial kernel function.
- gamma : { 'auto', 'scale' or float default = 'scale' } Kernel Coefficient for RBF, Poly, and Sigmoid SVR
- coef0 : { default = 0.0 (Non-Negative float) } Independent term in kernel function. Only applicable for Poly and Sigmoid

We used GridSearchCV to calculate the optimal parameters for 3 models, which search over specified arrays of parameter values. In order to see the significant difference in the model score, we set the input parameter array as follows.

- kernel : { 'linear', 'rbf', 'poly' }
- C : [0.1, 1, 10, 100, 1000]
- degree : { 3, 5, 8 }
- gamma : { 'auto', 'scale', 0.1, 0.01, 0.001, 0.0001 }
- coef0 : { 10, 25, 50, 100 }

5.3 Result

We applied GridSearchCV through three different kernels to SVR: the Linear Kernel, the Poly Kernel, and the RBF Kernel. Based on the results of the experiments, the R-squared scores are as follows:

Grid Search Parameter result			
Parameters	Linear Kernel Model	RBF Kernel Model	Poly Kernel Model
Kernel	linear	rbf	poly
C	1350	1000	210
degree	1	Not-applicable	3
gamma	Not-applicable	0.01	scale
coef0	Not-applicable	Not-applicable	25

Linear SVR R-squared score was: 0.7769931267399243

RBF SVR R-squared score was: 0.5085046369569125

Poly SVR R-squared score was: 0.8689104943523224

As observed, using Poly kernel generated better results. We used the R-squared score as the scoring function. $R^2 = 1 - \frac{u}{v}$, where $u = \sum (y_{true} - y_{predict})^2$ and $v = \sum (y_{true} - y_{true\ mean})^2$. A value closer to 1 indicates a better model fit.

5.3.1 Hypothesis to the result:

Based on comprehensive analysis, the Poly kernel performed better than the other two kernels, RBF and Linear kernel applied models. According to our hypothesis, polynomial kernel has better estimation than others due to the trend that more features lead to higher house prices, but no clear separation exists between the data. Separating the datasets with a liner kernel would not be enough, and RBF has a too complex way of setting support vectors. So the polynomial has some terms to distinguish the dataset, but they are not too complex like RBF.

Regularization strength is determined by the inverse proportion of C values. The C values are very high after parameter tuning. In our assumption, there is no clear separation in the dataset, so less regularization helps models perform better. If there was a lot more variance between data, regularization would help us generalize the model and counteract overfitting.

The coef0 values are only applicable to Polynomial kernel. Considering that it determines the number of independent terms in kernel functions, the polynomial model had a better result.

There is one aspect of the experiment that is unexpected: the testing score is higher than the training score. Here is the table of training and testing accuracy. (In the method of calculating scores, we used the R2 Score function).

Calculate training score	
Linear SVR R-Squared Score	0.7688902758375236
RBF SVR R-Squared Score	0.5168741555025075
POLY SVR R-Squared Score	0.8697297388733204

Calculate Testing Score	
Linear SVR R-Squared Score	0.8548107563261176
RBF SVR R-Squared Score	0.5498264522835576
POLY SVR R-Squared Score	0.8937439503876586

Based on both training and testing accuracy results, we found that the testing scores were slightly higher than the training scores. Using the SVR models, the testing prices of houses are predicted based on the training inputs. Our assumption was that the results might be misleading. In order to ensure that the result is correct, we trained and tested a linear regression model with the same PCA applied dataset.

The linear regression result is as follows.

Linear Regression model results	
Linear regression training score	0.9295112045275672
Linear regression testing score	0.7599034228795364

Using linear regression, it clearly shows that the training score is about 17% higher than the testing score. Consequently, the results would suggest that the interpretation is more likely to be accurate. Additionally, we can assume that the unexpectedly high accuracy of the testing results of the POLY SVR can be attributed to very fortunate tuning results, so in this case the POLY SVR model is still valid and the most accurate.

6 Conclusion

One of our greatest challenges throughout this process was managing the sheer magnitude of this data, especially due to the frequency of categorical variables. The original dataset had 80 features, but after one-hot encoding all of the categorical variables and preprocessing, it grew to 271 columns. This limited the amount of direct testing we could do due to runtime issues, especially as we did more intensive fitting for our SVM model. In the future, we would like to expand this project by trying some different models such as XG Boost. Since it is a gradient boosting tree model, based on each feature it makes separate trees and combines it all at once so it is better than the basic tree method and could provide better insights on the categorical data than what all of our one-hot encoding can do. Alternatively, we would like to explore other options outside of one-hot encoding for capturing the categorical variables. I believe that better preprocessing could greatly improve the accuracy of our results.

231 References

- 232 [1] Stanley H Chan. “Introduction to probability for data science”. In: *(No Title)* (2021).
- 233 [2] Mark Crovella. *Applications of the SVD - Linear Algebra, Geometry, and Computation*. Ac-
234 cessed 21 Apr. 2024. n.d. URL: [https://mcrovella.github.io/CS132-Geometric-
235 Algorithms/L26ApplicationsOfSVD.html](https://mcrovella.github.io/CS132-Geometric-Algorithms/L26ApplicationsOfSVD.html).
- 236 [3] IntuitiveMachineLearning. *Support Vector Machines: All You Need to Know!* Accessed 23 Apr.
237 2024. 2020. URL: [https://www.youtube.com/watch?v=ny1iZ5A8i1A&ab_channel=
238 IntuitiveMachineLearning](https://www.youtube.com/watch?v=ny1iZ5A8i1A&ab_channel=IntuitiveMachineLearning).
- 239 [4] MIT OpenCourseWare. *16. Learning: Support Vector Machines*. Accessed 21 Apr. 2024.
240 2014. URL: [https://www.youtube.com/watch?v=_PwhiWxHK8o&ab_channel=
241 MITOpenCourseWare](https://www.youtube.com/watch?v=_PwhiWxHK8o&ab_channel=MITOpenCourseWare).
- 242 [5] *Regularization*. URL: <https://www.ibm.com/topics/regularization>.
- 243 [6] *SciKitLearn Documentation*. Accessed 21 Apr. 2024. URL: [https://scikit-learn.org/
244 stable/index.html](https://scikit-learn.org/stable/index.html).
- 245 [7] *SciKitLearn Support Vector Regression Documentation*. Accessed 23 Apr. 2024. URL: [https:
246 //scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html](https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html).
- 247 [8] VisuallyExplained. *The Kernel Trick in Support Vector Machine (SVM)*. Accessed 30 Apr.
248 2024. 2022. URL: [https://www.youtube.com/watch?v=Q7vT0--5VII&ab_channel=
249 VisuallyExplained](https://www.youtube.com/watch?v=Q7vT0--5VII&ab_channel=VisuallyExplained).