House Price Prediction Methods Analysis Using Pattern Recognition Techniques

1st Mete DURLU

Computer Engineering Department
Baskent University
Ankara, TURKEY
metedurlu@gmail.com

2nd Turac SOGUTLU

Computer Engineering Department

Baskent University

Ankara, TURKEY

turac_000@hotmail.com

Abstract—There are a lot of pattern recognition techniques which help us analyze data and come up with valuable, refined information and predictions. The most common used techniques include Principle Component Analysis(PCA) [5]., Linear Discriminant Analysis(LDA) [6]., Naive Bayes, K-Nearest Neighbors(KNN) [1]. and Support Vector Machine(SVM) [3] [4]. . variants. These techniques help us bringing out hidden patterns and makes it easier for us to interpret our data. Experiments have been conducted with PCA [5], LDA [6], Naive Bayes, KNN [1], SVMs [3] [4] and Dense Neural Networks on House Price data-set [9] in order to better represent the effects of each method on a real world, time-series data with both categorical and numerical values as features. It has been concluded that a dense neural network with five hidden layers yields the best result with 85% among all methods mentioned. Source codes can be found at ... as notebooks.

Index Terms—pca, lda, svm, knn, regression, bayes, naivebayes, classification, neural, network, complex, data, data-set, analysis, pattern, recognition

I. INTRODUCTION

Before choosing a Pattern Recognition method data should be closely inspected. Pre-processing and interpretation methods should be chosen according to the data-set features and feature relations with problem specific target value. In this instance Ames Housing data-set [9] often cited as Boston Housing data-set has been used. Our data set has total of 81 features with multiple categorical, numerical labeled and continuous data. 81st feature is actually the target which is pursued to successfully predict. Other features include; overall quality of house, overall condition of house, linear feet of street connected to property, lot size in square feet, original construction date, remodel date, masonry veneer area in square feet, rating of basement finished area, total square feet of basement area, first floor area in square feet and many other explicit details about house samples. This relatively high feature amount and variance makes this data-set a tough estimation problem. Another field of difficulty arises when it is considered that some features for example "overall quality of house" scales from "0" to "10" but another feature "lot size in square feet" varies from 1300 to 215000 which causes logical and scalar imbalances. Yet another difficulty for pattern recognition is the high dimensionality of the data-set. Creates exponential complexity for some desired calculations such as KNN algorithm.

II. PRE-PROCESSING

While selecting the Pre-processing methods, nature of dataset at hand should always be in consideration. Not all data-sets are compatible with all Pre-processing methods. Inappropriate Pre-processing methods could have no effect or could even have negative effects.

A. Dealing With Missing Data

After careful inspection of data-set at hand it has been observed that missing values in data are actually correspond to missing features in houses. If an example need to be given; when a house has a garage, "garage size in square feet" has a positive numerical value, if there is no garage present then related feature gets a "Nan" value. Which has been solved by replacing "Nan" values with "0" values.

B. Normalization [7]

Normalization is a commonly used method before the pattern recognition applications. The goal of normalization is to transform features to be on a similar scale. This operation improves the performance and prevents the unbalanced contribution of feature values. It's simplicity makes it a go to Pre-processing method.

$$X' = (X - X_{min})/(X_{max} - X_{min})$$
 (1)

X being the feature/features selected to apply normalization, X_{max} is the maximum value for the selected feature and the X_{min} is the minimum value. By finding range $X_{max} - X_{min}$ then dividing the result of $X - X_{min}$, scaled values X' is produced.

C. Standardization [8]

Normalization is a commonly used method before the pattern recognition applications. With Standardization mean of values within features, will always be zero, and the standard deviation will always be one. The graph of standardized values will have exactly the same shape as the graph of raw data, but it may be a different size and have different coordinates. So basically feature space will be centered on origin point. These properties make standardization a meaningful choice of

Pre-processing. Rescaling formula is referred as the z-score formula.

$$Z = \frac{x - \mu}{\sigma} \tag{2}$$

Z = new feature values

x = original feature value

 μ = mean of the feature samples

 σ = standard deviation of the feature samples

D. Principle Component Analysis [5]

Principle component analysis is another popular Preprocessing method with a specific benefit of dimension reduction ability. PCA operation enables visualization of high dimensional data with meaningful 2D or 3D representations. Basically PCA constructs new variables using eigenvectors and eigenvalues. These values are chosen based on how well they can represent the data. If an eigenvalue of a principle component is high than that means it can be used to symbolize the change in data. PCA consists of four main steps.

- Compute the mean for every dimension of the whole data-set.
- 2) Compute the covariance matrix of the whole data-set.

$$cov(X_i, X_j) = \frac{1}{n-1} \sum_{k=1}^{n} (X_{ik} - \bar{x_i})(X_{jk} - \bar{x_j})$$
 (3)

3) Compute eigenvectors and the corresponding eigenvalues. (A = covariance matrix, I = identity matrix)

$$det(A - \lambda I) = 0 \tag{4}$$

4) Sort the eigenvectors by decreasing eigenvalues and choose k eigenvectors with the largest eigenvalues to form a d × k dimensional matrix W.

Acquired new matrix W, will be the new feature space for the data-set. PCA seems to have an effect of revealing linear distributions and clusters on simple data-sets.

E. Linear Discriminant Analysis [6]

Linear Discriminant Analysis (LDA) is a very common technique for dimensionality reduction problems as a preprocessing step for machine learning and pattern classification applications. LDA operation enables projection of high dimensional data with meaningful 2D or 3D representations. LDA consists of three main steps. The first step is to calculate the separability between different classes (i.e. the distance between the means of different classes), which is called the between-class variance or between-class matrix. The second step is to calculate the distance between the mean and the samples of each class, which is called the within-class variance or within-class matrix. The third step is to construct the lower dimensional space which maximizes the between-class variance and minimizes the withinclass variance.

F. Feature Correlation Matrix

Feature Correlation is basically the strength of relationship between features. Correlation can be calculated between each two feature independently. These relations makes possible to locate semantically identical features and most importantly gives information about which features affect the target the most.

$$cov(x,y) = \frac{1}{n-1} \sum_{k=1}^{n} (X_{ik} - \bar{x}_i)(X_{jk} - \bar{x}_j)$$
 (5)

$$r_{x,y} = \frac{cov(x,y)}{n-1} \tag{6}$$

- 1 indicates a strong positive relationship between x,y.
- -1 indicates a strong negative relationship between x,y.

G. Artificial Class Construction

When dealing with continuous data, options for pattern recognition methods are being reduced to regression. The subject data-set has continuous data as target values. This situation demands creative workarounds such as Artificial Class Construction as our terms. For this specific data-set 5 different classes were constructed by using standard deviation of target feature to compute class value intervals. Target feature got grouped into 5 classes with equal range. Equation below dishes out the class number.

$$r = \left(\frac{x}{2\sigma}\right) + 1\tag{7}$$

III. CLASSIFICATION METHODS

Classification is initially a method for predicting class values by putting feature values of a samples through certain processes or equations. It can be achieved through either supervised or unsupervised methods. The most common methods are Naive Bayes, K-Nearest Neighbors and Support Vector Machine. These methods are all supervised methods where the number of classes and class labels for each training sample is known as prior knowledge. All mentioned methods have been tried on data-set at hand with artificially crafted classes. Conclusions can be seen on results part with all evaluation applied.

A. Naive Bayes

Naive Bayes method is the most common and easy to implement supervised classification method. Although it is simple, it can produce relatively good results in some real-life cases. In this regard it has became the first method to try before anything else.

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$
(8)

$$P(c|X) = \prod_{i=1}^{n} P(x_i|c)$$
(9)

P(c|x) is the posterior probability of class given feature x is present. P(c) is the prior probability of class. P(x|c) is the

likelihood which is the features probability of occurrence on given class. P(x) is the prior probability of feature.

The main reason, it is not superior to the more complex methods, is the fact that Naive Bayes makes the assumption of features being independent from each other. Unfortunately this assumption is almost always wrong.

B. K-Nearest Neighbors [1]

K-Nearest Neighbors is a really powerful method with many applications. It is simple, easy to implement but has a tendency of getting too computationally expensive due to its nature. It is a supervised method relying on training sets class labels to classify new data. Basically it computes distance between each training sample and new input sample. Then uses the least distant "K" neighbors in order to decide on the class information. For example, if "K" equals three and then there are three votes to decide on new samples class. So three nearest neighbors could have class labels of one, one and zero. Then algorithm decides that new sample should have a class label of one. There is no need for hyper-parameters except "K" and there is no training but it becomes computationally expensive to calculate distances as the dimension and sample size increases. Most commonly euclidean distance is used to measure the distance between samples.

$$d(x,y) = \sum_{i=1}^{n} \sqrt{(x_i - y_i)^2}$$
 (10)

d = distance between sample x and sample y

 $x_i = i^{th}$ feature of sample x

 $y_i = i^{th}$ feature of sample y

n = number of features (dimension of feature space)

C. SVM Classifier(SVC) [3]

Support Vector Machine is a supervised machine learning method which is a powerful and efficient tool. It can be used for both classification and regression problems. SVMs main goal is finding a hyperplane that best divides a data-set to two different classes multiple times(as many times needed to match number of classes). Support vectors are the data points nearest to the hyperplane, these point help define the hyperplane so all computations are done through these points. This hyperplane creates an area of margin which divides two classes apart. Error function here is designed so that the margin becomes larger as error decreases. If there is no clearly dividing hyperplane then the whole feature space is transformed into a new higher dimension feature space. This is known as kernelling. SVMs produce accurate results on clean data-sets with small to medium sample size. When dealing with larger data-sets however computational costs can be too much to handle and also it is highly sensitive to the noisy nature of large data-sets.

$$\sum_{i \in SV} y_i \alpha_i K(x_i, x) + b \tag{11}$$

 $y = class based constant always \ge 1$

 α = Function coefficent

K = Kernel Function

x =Support Vector

b = Independent Term or Intercept

IV. REGRESSION METHODS

Regression analysis is a predictive technique which investigates the relationship between target feature and other given features in feature space. This technique is used for forecasting, time series modelling and finding the causal effect relationship between the variables.

A. SVM Regression(SVR) [4]

It has been already mentioned that SMVs are good at both regression and classification. Reminding, SVMs main goal is finding a hyperplane that best divides a data-set to two different classes for classification. On regression this hyperplane actually tries to mimic data-sets behaviour and error function tries to minimize points outside of the determined margin. Basically SVR is trying to decide a decision boundary at 'C' distance from the original hyperplane such that data points closest to the hyperplane or the support vectors are within that boundary line.

$$Prediction = \sum_{i \in SV} (\alpha_i - \alpha_i^*) K(x_i, x) + b$$
 (12)

 α, α^* = Function coefficients, bound to given distance 'C'

K = Kernel Function

x =Support Vector

b = Independent Term or Intercept

B. Dense Neural Network Regression

Over the last decade Neural Networks have became the primal solution to many problems including classification and regression in data-sets of all sizes. A Neural Network usually described as a blackbox taking input and producing output but actually they can be simply explained as long functions with adjusted weights. The smallest units in the Neural Networks are neurons. And the whole structure consists of layers made from neurons. Each neuron in a layer receives an input from all the neurons present in the previous layer. Before an input is used in a neuron it will be multiplied by the weight value adjusted by the optimizer in Neural Network. Optimizer as in its name optimizes the weight values so that the final output becomes similar to the desired value. On each iteration, output will be put in an error function and according to the error optimizer will adjust the weights on Neural Network. It could be imagined as a hyperplane of errors calculated by each output pushed out. The goal is to find the global minimum of the hyperplane by adjusting weights. There are a lot of equations to choose from in order to optimize the weights and calculate the error. The problem on hand should be considered before selecting these equations. In our situation:

Optimizer: ADAGRAD

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t \tag{13}$$

ADAGRAD is a gradient descend variant which is trying to find the global minimum of the error hyperplane. Error Function: Mean Absolute Error

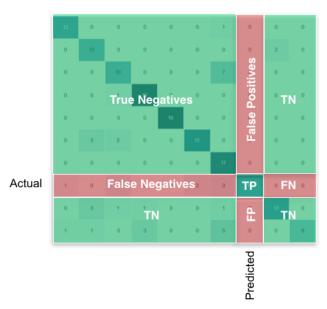
$$MAE = \frac{1}{n} \sum_{i=1}^{n} |x_i - x|$$
 (14)

V. EVALUATION METHODS AND METRICS

There are multiple evaluation methods and metrics used during this work. Primarily confusion matrix and metrics derived from it such as "Accuracy", "Sensitivity", "Precision", "Specificity" were used for classification evaluation. For regression evaluation \mathbb{R}^2 metric and for calculating error "Mean Absolute Error" has been used. These metrics were chosen according to the problem.

A. Confusion Matrix

Confusion Matrix consists of True Positive(TP), True Negative(TN), False Positive(FP) and False Negative(FN) values. Which describe the models behavior on predicting classes. TPs represent predictions of actual positive cases as positive, TNs represent predictions of actual negative cases as negative, FPs represent predictions of actual negative cases as positive and FNs represent predictions of actual negative cases as negative.



Using these values, more elaborate metrics such as "Precision", "Sensitivity", "Specificity" and "Accuracy" can be calculated. These metrics can be calculated both classwise and modelwise.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{15}$$

$$Precision = \frac{TP}{TP + FP} \tag{16}$$

$$Sensitivity = \frac{TP}{TP + FN} \tag{17}$$

$$Specificity = \frac{TN}{TN + FP} \tag{18}$$

These metrics enables classwise semantic evaluation of model.

B. Mean Absolute Error

In regression problems there has to be an error function measuring models behaviour in order to optimize the model. For this purpose "Mean Average Error" function has been chosen. It reflects the absolute distance between sample points and regression function.

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |x_i - x|$$
 (19)

C. R² Method

 R^2 computes the coefficient of determination. It represents the proportion of variance (of target feature) that has been explained by the independent variables in the model. It provides an indication of goodness of fit and therefore a measure of how well unseen samples are likely to be predicted by the model, through the proportion of explained variance. Best possible score is 1.0 and it can be negative (because the model can be arbitrarily worse).

$$R^{2}(y,\hat{y}) = 1 - \frac{\sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}}$$
(20)

 \hat{y}_i = the predicted value of the $i^t h$ sample

 y_i = the corresponding true value for total n samples

VI. RESULTS

So far each method applied to data-set at hand has been explained briefly. From now on the empiric experimentation results will be discussed. On each step evaluation results were compared and for next step the process with best result has been carried on.

Pre-processing Evaluation

After a detailed inspection of data "Nan" values has been dropped initially but after further steps the definition of "Nan" values has became more clear as they represent when there is no instance of the feature on the sample. Regarding this information "Nan" values have been replaced by "0" and prevented data loss.

As scaling, Normalization and Standardization results have been compared. For Standardization changes the variance of the data Normalization has been chosen. By Normalization original shape of data has been stored.

Dimension reduction methods PCA and LDA has been applied seperately and consecutively but no improvement was observed on any case. Thus PCA and LDA was subtracted from pre-processing.

Classification Evaluation

Before applying classification methods, artificial classes created from target feature.

First of all Naive Bayes approach has been applied. Naive

Bayes has produced respectable accuracy values such as 81% but on a closer inspection it has been revealed that class imbalances were affecting sensitivity and specificity values negatively so much so that it made the model unreliable as it had 50% average sensitivity.

Second method was KNN. KNN had remarkable accuracy, precision and sensitivity results every measurement were over 80% even with using different "K" values. For instance; for "K" equals 1 accuracy was 82%. With "K" value at 5 best metric values (90% accuracy, 66% precision, 87% specificity) were observed and larger "K" values were producing worse results. But going into details, imbalances on class samples and nature of the data-set at hand revealed that KNN was also unreliable. Model had the tendency to ignoring classes with less sample size.

Lastly for classification methods SVM classifier was used. SVM classifier with a second degree polynomial kernel yielded subpar results. As the function degree get higher the produced results were also improving but SVM was also suffering greatly by imbalanced class samples. Although accuracy and specificity values were amazing, sensitivity values were suffering greatly.

To conclude classification evaluation, even if the SVM and KNN methods were quite successful, due to the nature of given problem and data-set results are not even close to the regression when converted from class labels back to continuous values. This proves that our "Artificial Class Creation" method was a bad solution to convert the problem from a regression case to classification case.

Regression Evaluation

SVM regression was the first method that has been applied. SVM regression was applied using polynomial kernel with various degree values. Observations showed that SVM regression had scored $0.84\ R^2$ score on training data but this did not translate to test data while predictions on test data could only achieve $0.51\ R^2$ at best. Next a dense neural network has been trained for classification. Neural network consisted of five hidden layers. The whole model can be summarized as:

TABLE I Model: "Sequential"

Layer (type)	Output Shape	Num of Params
dense (Dense)	(None, 256)	5888
dense_1 (Dense)	(None, 128)	32896
dropout (Dropout)	(None, 128)	0
dense_2 (Dense)	(None, 64)	8256
dropout_1 (Dropout)	(None, 64)	0
dense_3 (Dense)	(None, 32)	2080
dense_4 (Dense)	(None, 16)	528
dense_5 (Dense)	(None, 1)	17

Total params: 49,665 Trainable params: 49,665 Non-trainable params: 0 As error function Mean Absolute Error(MAE) has been used and Adamgrad was chosen as optimizer. After initializing model several training sessions with different hyperparameters were constructed. With a high learning rate model would overfit with up to $0.95\ R^2$ score on training data and run worse on test data predictions. Even with higher regularization rates overfit situation could not be averted. Best results on test data were observed with low learning rate and also low regularization values. Epochs were taken as constants as higher epoch values would result in overfit models. After hyperparameters were tuned, best test data results were computed with DNN. Latest model achieved $0.85\ R^2$ score on training data and $0.86\ R^2$ score on test data.

To conclude regression evaluation. It has been proven that House Price data-set was most suitable for regression methods. Although results on first glance could favor classification methods, closer inspection revealed that accuracy could not be used alone for evaluating models. And Dense Neural Networks has surpassed SVM on the task of regression.

TABLE II PATTERN RECOGNITION RESULTS

Methods	Metric Results (0-1)			
Classification	Hyper-params.	Accuracy	Sensitivity	
Naive Bayes	-	0.81	0.45	
KNN	K = 1 distance = "Euclidian"	0.8	0.54	
KNN	K = 3 distance = "Euclidian"	0.82	0.56	
KNN	K = 5 distance = "Euclidian"	0.84	0.58	
KNN	K = 7 distance = "Euclidian"	0.83	0.52	
SVC	f = Polynomial d = 1 c = 1.0	0.8	0.49	
SVC	f = Polynomial d = 2 c = 1.0	0.85	0.61	
SVC	f = Polynomial d = 3 c = 1.0	0.86	0.61	
SVC	f = Polynomial d = 5 c = 1.0	0.87	0.6	
SVC	f = Polynomial d = 7 c = 1.0	0.86	0.59	
Regression	Hyper-params.	Error Func.	R^2	
SVR	f = Polynomial d = 1 c = 1.0	-	0.56	
SVR	f = Polynomial d = 2 c = 0.5	-	0.42	
SVR	f = Polynomial d = 3 c = 0.9	-	0.31	
SVR	f = Polynomial d = 5 c = 1.0	-	0.34	
SVR	f = Polynomial d = 7 c = 1.0	-	0.23	
DNN	$\lambda = 0.01 \text{ E} = 250 \text{ D} = 0.2$	MAE	0.62	
DNN	$\lambda = 0.001 \text{ E} = 250 \text{ D} = 0.15$	MAE	0.67	
DNN	$\lambda = 0.0001 \text{ E} = 250 \text{ D} = 0.1$	MAE	0.86	

K = Number of Neighbors, f = Kernel Function, λ = Learning Rate, E = Epochs, D = Regularization Rate

ACKNOWLEDGMENT

We thank our Lecturer Dr. Duygu DEDE ŞENER from Computer Engineering Department on Baskent University (ddede@baskent.edu.tr) for all of her efforts and support.

REFERENCES

- [1] Guo, Gongde and Wang, Hui and Bell, David and Bi, Yaxin. KNN Model-Based Approach in Classification. 08/2004.
- [2] Harris Drucker, Chris J.C. Burges, Linda Kaufman, Alex Smola, Vladimir Vapoik. Support Vector Regression Machines. Bell Labs and Monmouth University Department of Electronic Engineering
- [3] Sklearn: SVM-Classification, https://scikit-learn.org/stable/modules/svm.html#classificat

[4] Sklearn: SVM-Regression,

https://scikit-learn.org/stable/modules/svm.html#svr

[5] Sklearn: PCA,

https://scikit-learn.org/stable/modules/decomposition.html#principal-component-analysis-pca

[6] Sklearn: LDA,

https://scikit-learn.org/stable/modules/lda_qda.html#linear-and-quadratic-discriminant-analysis

[7] Sklearn: Normalization,

https://scikit-learn.org/stable/modules/preprocessing.html#standardization-or-mean-removal-and-variance-scaling

[8] Sklearn: Standardization,

https://scikit-learn.org/stable/modules/preprocessing.html#standardization-or-mean-removal-and-variance-scaling

[9] Dean De Cock Ames, Iowa: Alternative to the Boston Housing Data as an End of Semester Regression Project. Journal of Statistics Education Volume 19, Number 3(2011)