

MODEL COMPARISON ON MOBILE PRICE PREDICTION

Authors: Jingjie Ma; Ella Wang; Han Zhang

Data Science And Analyst Program Of Georgetown University

INTRODUCTION

In modern society, mobile phones are displaceable utilities that have been broadly used in recent decades. From the first handheld mobile was demonstrated in 1973, it has been well developed for 40 years in the market. The rapidly and continuously enhanced technologies and features in the competitive market determine the high variance of its price. The dataset gathered from Kaggle consists of different features of mobile phones. The goal is to compare different classification methods used to predict the mobile price by accuracy score and explain their algorithms.

Predictors

- battery_power**: Total energy a battery can store (quantitative)
- clock_speed**: Speed at which microprocessor executes instructions (quantitative)
- ram**: Random Access Memory in Megabytes (quantitative)
- dual_sim**: Has dual sim support or not (0 for no, 1 for yes)
- int_memory**: Internal Memory in Gigabytes (quantitative)
- n_cores**: Number of cores of processor (quantitative)
- four_g**: Has 4G or not (0 for no, 1 for yes)
- three_g**: Has 3G or not (0 for no, 1 for yes)
- blue**: Has bluetooth or not (0 for no, 1 for yes)
- touch_screen**: Has touch screen or not (0 for no, 1 for yes)
- wifi**: Has wifi or not (0 for no, 1 for yes)
- fc**: Front Camera megapixels (quantitative)
- pc**: Primary Camera megapixels (quantitative)
- m_dep**: Mobile Depth in cm (quantitative)
- mobile_wt**: Weight of mobile phone (quantitative)
- px_height**: Pixel Resolution Height (quantitative)
- px_width**: Pixel Resolution Width (quantitative)
- sc_h**: Screen Height of mobile in cm (quantitative)
- sc_w**: Screen Width of mobile in cm (quantitative)
- talk_time**: Longest time that a single battery charge will last in hours (quantitative)

Response Variable

- Price Range**: Value of 0 (low cost), 1 (medium cost), 2 (high cost) and 3 (very high cost)

The dataset has 2000 rows and 21 columns in total, with 20 predictors and 1 response variable.

EXPLORING DATA

The response variable of this dataset is price range, which is split into 4 classes. To verify the classification accuracy, multi-class classification is needed to be done rather than binary classification. Also, this dataset has 20 features and 800 samples, so $n_samples$ (number of rows) > $n_features$ (number of features). This might be a crucial factor for the accuracy of the classification methods below.

CLASSIFICATION METHODS AND ACCURACY

- K-Nearest Neighbors
- Random Forest
- Support Vector Machine
- Neural Network

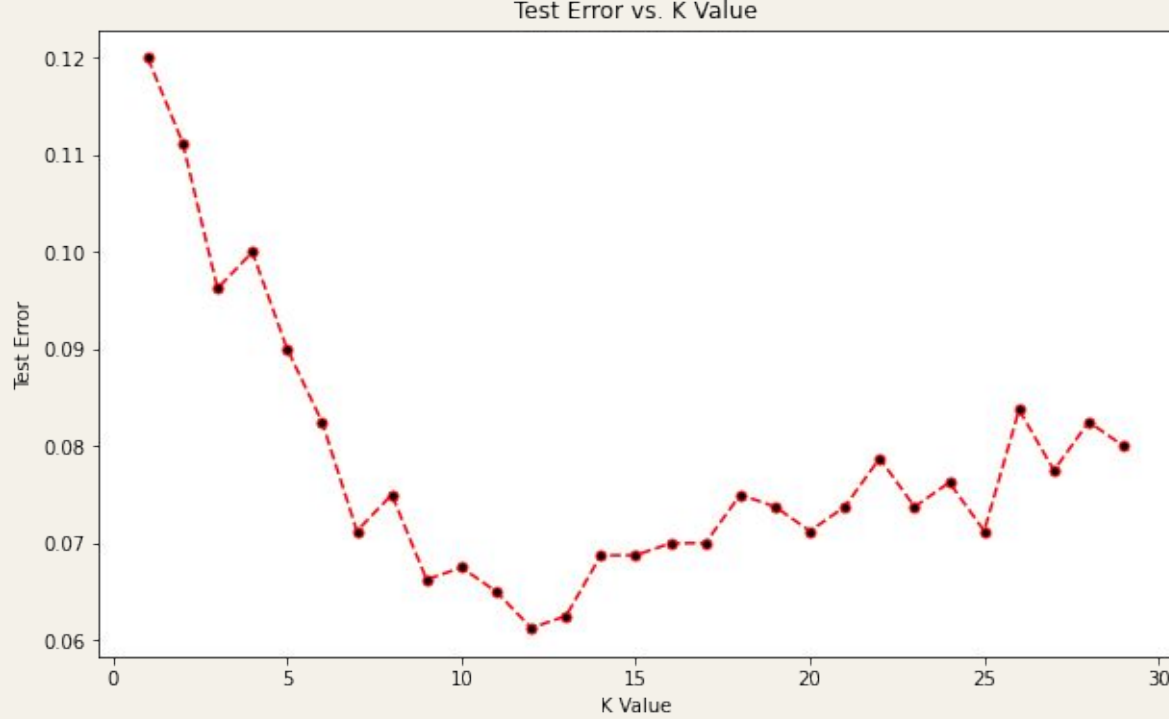
The criteria to compare each classification method is accuracy, which is the “score” output of all methods above. It measures the mean percent of train vectors are correctly classified.

K-NEAREST NEIGHBORS

The k-nearest neighbors algorithm for classification in Sklearn Package measures the Euclidean distance between the target vectors and the train vectors and uses the majority vote of the nearest k neighbors of each point to decide the classification result. This is a **convex** program. The Euclidean distance is calculated as:

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

where x and y are data points in n-space. The optimal k value was chosen by plotting different k values against their test error. In this project, the accuracy is **93.9%** with k = 12.



RANDOM FOREST

Random Forest is built on Decision Trees and predictions are the averaged prediction of the individual classifiers. For Random Forest Classifier in Sklearn package, each tree in the ensemble is built from a sample drawn with replacement (bootstrapping sample) from the training set. $n_estimators$ and $max_features$ are two main parameters to adjust. The best parameter values should always be cross-validated by the accuracy score. So this program is optimized when the maximum of its accuracy score is reached. In this project, the accuracy of random forest classification is **92.1%** with $n_estimators = 100$ and $max_features = 14$.

SUPPORT VECTOR MACHINE

The Support Vector Machine algorithm is embedded in the package “scikit-learn”. In the dataset, there are 4 classes in “price range” but SVM is designed for binary classification problems. In this package, multi-class classification problem can be split into multiple binary classification problems and train a binary classification model each (OvO). Hence, this solves the **convex** primal program for each pair of classes and for 6 classifications ($4 \times 3/2$), for class i and j :

$$\min \frac{1}{2} (w)^T w + C \sum_{i=1}^n (\zeta_i)$$

subject to $((w)^T \phi(x_i) + b) \geq 1 - \zeta_i$,
 $((w)^T \phi(x_i) + b) \leq -1 + \zeta_i, \zeta_i \geq 0$

where $x_i \in R^p$, $\phi(x_i)$ maps x_i into a higher dimensional space and $y_i = \{-1, 1\}^n$. And $C > 0$ is the regularization parameter. Here, let $C = 1$. The kernel function of radial SVM is:

$$K(x_i, x_j) = \exp(-\gamma ||x_i - x_j||^2) = \phi(x_i)^T \phi(x_j)$$

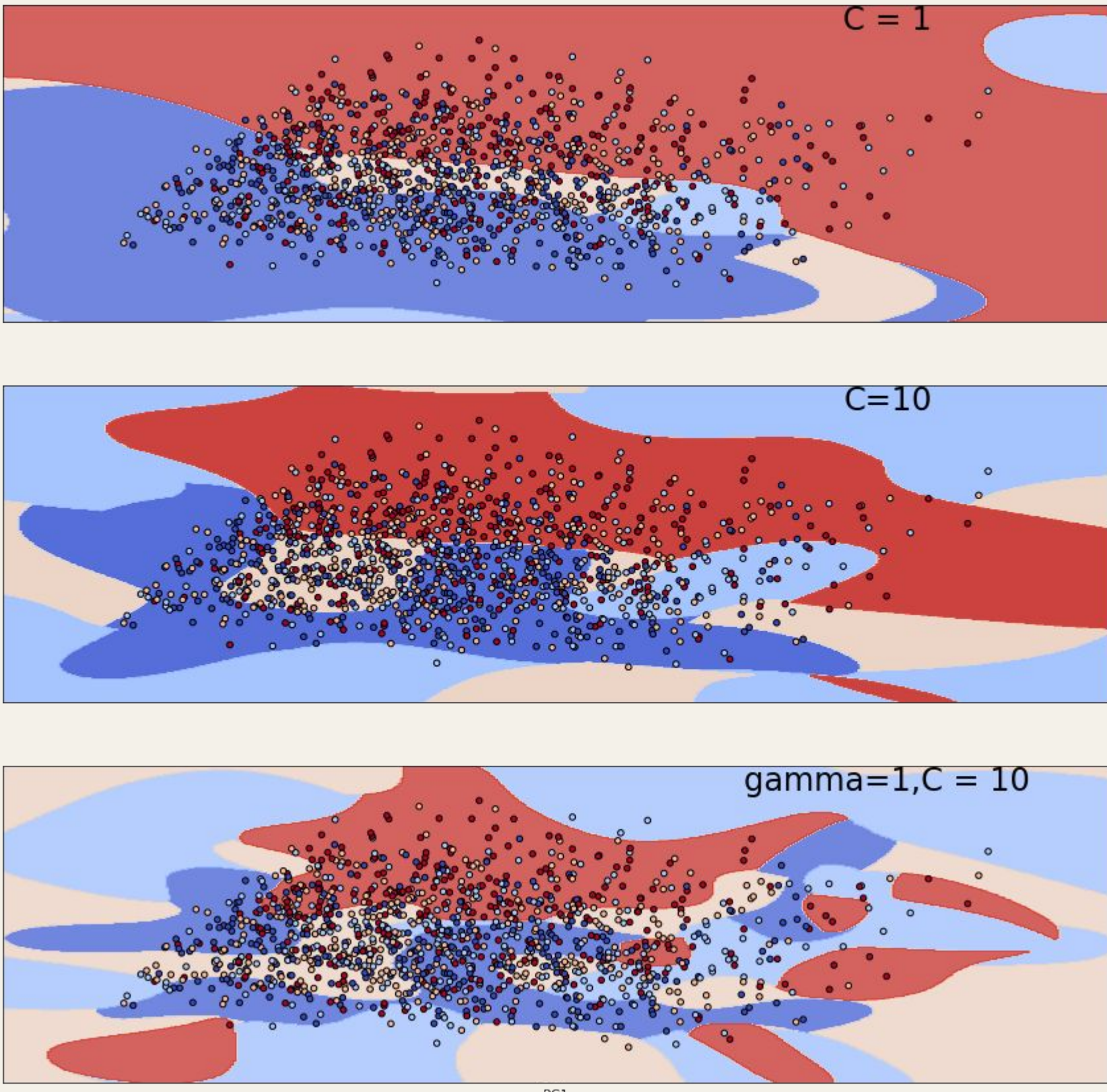
where $\gamma = \frac{1}{20} = 0.05$ here since there are 20 features. Then, solve the corresponding dual program to get the optimal outputs ω and b . And by this decision function: $\text{sgn}(\omega^T \phi(x_i) + b)$ The predicted labels are gathered and recorded. The final predicted labels are determined by majority voting of predictions from 6 such dual programs.

So far, the major algorithm the package is using is **coordinate descent**. By python, the accuracy score of this method is: **86.5%**.

SVM IN SKLEARN

- Linear SVM:
Kernel SVM is a great method with higher accuracy than linear SVM but sometimes, for example, when the number of samples is small, there is no need to use kernel trick. Instead, the linear SVM will be more suitable. In sklearn package, there is also algorithm to solve this optimization program. With the same primal program as Kernel SVM, linear SVM only need to solve 4 such programs since it trains the dataset with OvR, not OvO.
- Parameter C, gamma:
C and gamma are two parameters of SVM in sklearn (only C if it is linear). Changing C and gamma will give a accurate classification result. The greater C is, the smoother and curlier the decision function will be. And the greater gamma is, the better to capture the shape of train vectors. Thus, increasing C and gamma leads to overfitting.

Compare different C and gamma



Neural Network

We built a sequential neural network to predict the price. We used three layers of neurons with dense 16, 12, and 4, respectively. We are using Relu as activation function for the first two layers and softmax as the output function. The Relu function preserves the non-negative part of the variable. The softmax function takes a vector of values and transforms it into a vector of values that sum up to 1. We used Adam as the optimizer. Adam is a popular method that is used by stochastic gradient descent. Unlike regular gradient descent, stochastic gradient descent uses only a stochastic or random batch of training samples to update parameters. Thus, stochastic gradient descent is suitable when the dataset is very large. The Adam optimizer does not only get stochastic gradient at each iteration, it also updates the first moment (the mean) and the second raw moment (the uncentered variance) accordingly. The algorithm of the Adam optimizer is shown below:

Require: α : Stepsize
Require: $\beta_1, \beta_2 \in [0, 1]$: Exponential decay rates for the moment estimates
Require: $f(\theta)$: Stochastic objective function with parameters θ
Require: θ_0 : Initial parameter vector
 $m_0 \leftarrow 0$ (Initialize 1st moment vector)
 $v_0 \leftarrow 0$ (Initialize 2nd moment vector)
 $t \leftarrow 0$ (Initialize timestep)
while θ_t not converged **do**
 $t \leftarrow t + 1$
 $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)
 $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)
 $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)
 $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)
 $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)
 $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)
end while
return θ_t (Resulting parameters)

We also used dropout technique to avoid potential overfitting. Dropout basically ignores some units and they will not be included in both forward and backward path. The accuracy of the neural network model on the testing set is **93.5%**.

CONCLUSION

Method	Accuracy
K-Nearest Neighbors	93.9%
Random Forest	92.1%
SVM	94.8%
Neural Network	93.5%

REFERENCES

Platt JC. Probabilistic outputs for SVMs and comparisons to regularized likelihood methods. <http://research.microsoft.com/~jplatt>. 1999 Mar 26. <https://www.cs.colorado.edu/~mozer/Teaching/syllabi/6622/papers/Platt1999.pdf>

Chang C-C, Lin C-J. LIBSVM: A Library for Support Vector Machines. <https://www.csie.ntu.edu.tw>. 2019 Nov 20.

Hsu C-W, Chang C-C, Lin C-J. A practical guide to support vector classification. Edu.tw. [accessed 2020 Dec 2]. <https://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>

sklearn.tree.DecisionTreeClassifier — scikit-learn 0.23.2 . [accessed 2020 Dec 2]. <http://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

Nearest Neighbors Classification | Statistical Classification . [accessed 2020 Dec 2]. <https://www.scribd.com/document/421227117/Nearest-Neighbors-Classification>

Kingma DP, Ba J. Adam: A Method for Stochastic Optimization. arXiv.org. 2017 Jan 30 [accessed 2020 Dec 5]. <https://arxiv.org/abs/1412.6980>

Contact Information

Ella Wang (yw701@georgetown.edu)
Jingjie Ma (jm3292@georgetown.edu)
Han Zhang (hz274@georgetown.edu)