
Predicting Global Sales of Video Games

Kevin Weng, Yi Lyu, Chen Li, Omar Hafez
MATH 156 Final Project: Group 5
University of California, Los Angeles (UCLA)
[Github Project](#)

Abstract

Our project aims to model global video game sales by applying a neural network, Random Forest, and k-Nearest Neighbors model to preprocessed data to successfully predict global sales and determine which model works best. We built the project on Python using sklearn, Keras, and Tensorflow for the models and SQL for data processing. We will evaluate the results by root mean squared error (RMSE). We hope to see whether or not global sales can be effectively modeled and also give an analysis on model strengths and weaknesses for this task.

1 Introduction

The expansion of the video game industry and eSports in the past decade has fueled gamers and game studios all over the world. With the backdrop of COVID-19, worldwide quarantines, and work-from-home structures, video games have garnered another boost in popularity and success in the first half of 2020. Looking forward, it is apparent that the landscape of video games and entertainment software is drastically changed with companies like Nintendo, Activision Blizzard, and many independent developers making an impact on modern culture. Creating a video game is an intensive project that requires a diverse array of resources and specialists that would be very costly for the studio if a release goes awry. As a video game producer, investor, or consultant, the ability to project global sales is very useful when it comes to considering possible translations, global releases, and general marketing investment in other parts of the world. In an effort to solve this problem, we will apply a neural network, Random Forest, and k-Nearest Neighbors to model global video game sales in a variety of different ways and determine which method is best and whether or not video game sales can be predicted from our data.

2 Related Works

There are other works related to our project in the video game industry, and countless more in the generalized sales and marketing prediction field. One of the projects involves internet search volume as a feature to predict global sale and is becoming increasingly relevant as social media dominates the information space in the majority of the video game industry's target audience^[5]. This data is likely heavily correlated to global sales as the consumer sentiment is captured even prior to release. Another paper that used neural networks predicted weekly game sales on PCA preprocessed data^[3]. The weekly timeframe is different from our cumulative global sales number and may be impacted seasonally. Additionally, one other paper used sexualized cover art content as a feature to predict sales for video games^[4]. In this case, another specific feature was analyzed that we were not able to consider for our project. As pertaining to sales, features regarding behavioral economics, consumer

psychology, and marketing can all be possible candidates to more successfully model video game sales.

3 Dataset and Features

The dataset we are going to train our model is **Video Game Sales with Ratings** from Kaggle. The dataset consists of 11,563 video game titles detailing release year, publisher, platform, genre, regional sales, global sales, critic and user scores, critic and user counts, and ESRB rating. Not all of the features are present for every title. The critic and user scores were obtained from Metacritic, a popular video game review site.

	Platform	Year_of_Release	Genre	Publisher	Global_Sales	Critic_Score	Critic_Count	User_Score	User_Count	Developer	Rating
0	Wii	2006.0	Sports	Nintendo	82.53	76.0	51.0	8	322.0	Nintendo	E
1	NES	1985.0	Platform	Nintendo	40.24	NaN	NaN	NaN	NaN	NaN	NaN
2	Wii	2008.0	Racing	Nintendo	35.52	82.0	73.0	8.3	709.0	Nintendo	E
3	Wii	2009.0	Sports	Nintendo	32.77	80.0	73.0	8	192.0	Nintendo	E
4	GB	1996.0	Role-Playing	Nintendo	31.37	NaN	NaN	NaN	NaN	NaN	NaN

Figure 1: First 5 entries of dataset

3.1 Features

- Name - The name of the video game.
- Platform - The console on which the game runs on. (Wii, PS4, PC, etc.)
- YearofRelease - The year the game was released.
- Genre - Category of the game. (Shooter, Racing, Puzzle, etc.)
- Publisher - Publisher of the game.
- NASales, EUSales, JPSales, OtherSales - Local region sales of video games in largest markets in millions of units.
- Global Sales - Total sales in the world in millions of units.
- Critic_score - Score by Metacritic's critics.
- Criticcount - Number of critics who contributed to Critic_score.
- User_score - Score by Metacritic's subscribers.
- Usercount - Number of users who contributed to User_score.
- Developer - Party who created the game.
- Rating - The Entertainment Software Rating Board (ESRB) rating.

3.2 Preprocessing

First, we removed the games containing missing Platform, Genre, Publisher, and YearofRelease because these variables could not be imputed effectively. Next, since many of the values for Critic_score, User_score, Criticcount, and Usercount were missing, we imputed the missing values with the median of the column. For the remaining categorical data in Genre and Publisher, we used one-hot-encoding to make the data suitable for fitting regression and dropped one column of each to decorrelate the columns. Our preprocessing also removed the local sales of games because they were obtained after global sales was calculated, hence would not be useful to show a correlation with prior regional video game sales and a global launch. Our preprocessing also removed the local sales of games because they were often obtained after global sales were calculated and thereby would not be useful for showing a correlation with prior regional video game sales and a global launch. Because of the unusual popularity of certain games, such as Grand Auto Theft V, we remove the top 10% and bottom 10% as outliers. After this preprocessing of the data, we split the

73 remaining into 80% training data and 20% test data. For k-Nearest Neighbors, we split
 74 the training data further into actual training data and cross validation data, which is used
 75 to determine the optimal k number of neighbors to be used that minimizes RMSE. The
 76 correlation matrix in Figure 2 was used to determine whether or not any features needed to
 77 be dropped.

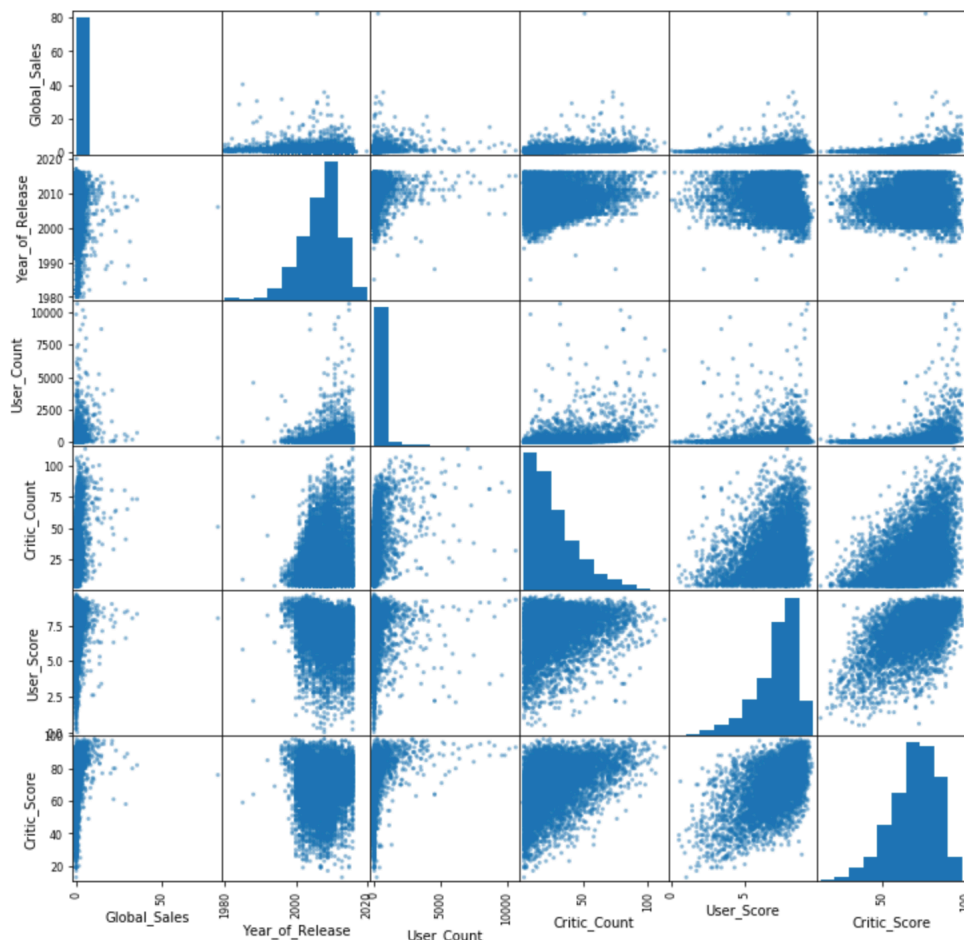


Figure 2: Correlation Matrix of Features

78 4 Methods

79 4.1 Models

80 Neural Network - The Neural Network as studied in class is a model comprised of hidden
 81 layers of nodes known as neurons that takes in an input and, by feeding it through
 82 the hidden layers, produces a result in the output layer. The neural network is built
 83 to handle large volumes of data, but since our dataset does not have relatively as
 84 many examples and some were removed during preprocessing, the neural network
 85 may not take advantage of its processing ability to return the strongest predictor.
 86

87 For our experiment, we will build an Artificial Neural Network (ANN) from the
 88 Sequential class in Keras. We have two hidden layers of 6 and 4 nodes respectively,
 89 and use the Exponential Linear Unit (ELU) activation function. We chose ELU
 90 because of its fast learning rate. Additionally, we used the 'Adam' optimizer offered
 91 by Keras, which is a stochastic gradient descent method and trained for 100 epochs.

92 Random Forest - The Random Forest model is in ensemble method that trains multiple
 93 decision trees and outputs a class by majority vote^[1]. For regression tasks, instead
 94 of mode, the mean prediction of the individual trees is returned. For reference, a
 95 decision tree is a popular machine learning algorithm that uses many input variables
 96 to traverse down a tree, and returns a prediction from a leaf. The benefit of using
 97 multiple trees is a reduced variance as a single tree can easily overfit the data.
 98 Random Forest differs from simply bagging multiple decision trees by selecting a
 99 random subset of the features for each tree so that if some features are stronger pre-
 100 dictors than others, such trees would be correlated known as the 'Random Subspace
 101 method^[6]'. Fewer correlated trees would improve the prediction while also making
 102 feature evaluation more accessible (We do not conduct feature evaluation in this
 103 project but it would be a very important area of exploration for practical purposes).
 104

105 For our experiment, we used the built-in RandomForestRegressor in sklearn with
 106 mse criterion. This model was the simplest to implement.

107 k-Nearest Neighbors - The k-Nearest Neighbors algorithm takes an element and looks for
 108 its closest neighbors to take a majority vote in the classification case, and the mean
 109 value of the nearest neighbors for regression^[2]. For our purposes, we used the mean.
 110 The basic structure of the algorithm is as follows:

- 111 1. We choose a k from all possible k for our model.
- 112 2. We divide the training dataset into p equal parts.
- 113 3. We randomly choose one part for cross validation and the remaining $p - 1$ parts
 114 to for training, yielding us an error, which we will repeat for p times so that
 115 each part is used once as cross validation set. We then compute the average
 116 error of this model given k over the p parts.
- 117 4. We find the k that minimizes the RMSE and return the model.

For our experiment, we used the KNeighborsRegressor from sklearn with a calcu-
 lated optimal $k = 22$ and weighed the points by the inverse of their distance, making
 closer neighbors have greater influence:

$$\hat{y} = \sum_{i=1}^k \frac{d(x, x_i) y_i}{\sum_{i=1}^k d(x, x_i)}$$

118 5 Experiments/Results/Discussion

119 5.1 Methodology

After training the three models, we will evaluate the results with the root mean squared
 error (RMSE) metric to evaluate efficacy and accuracy.

$$RMSE = \sqrt{\left(\frac{1}{n}\right) \sum_{i=1}^n (y_i - x_i)^2}$$

120 RMSE takes into account negative values and is a commonly used metric in determining
 121 the performance of regression models.

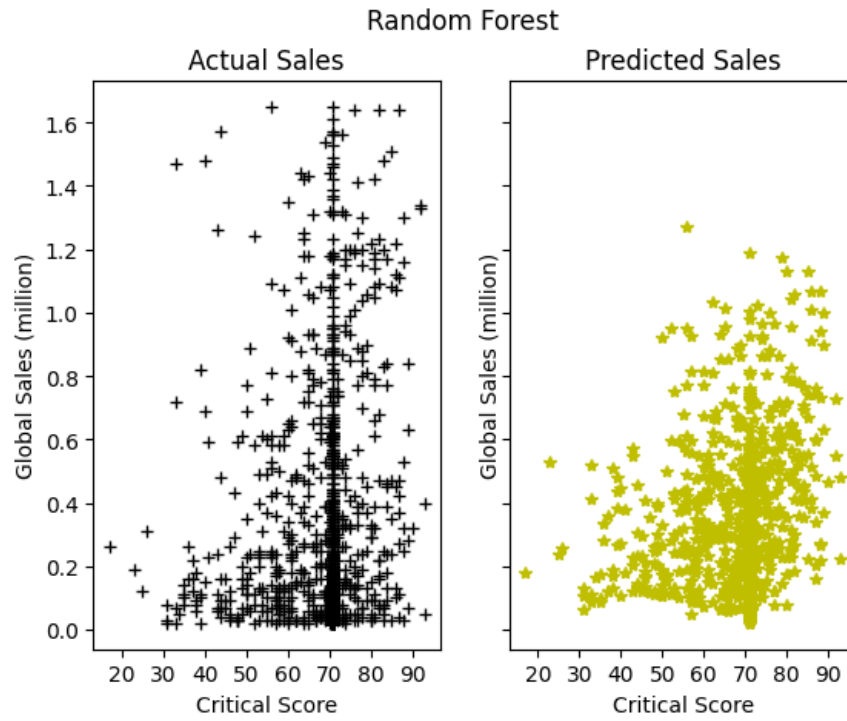


Figure 3: Random Forest Predictions

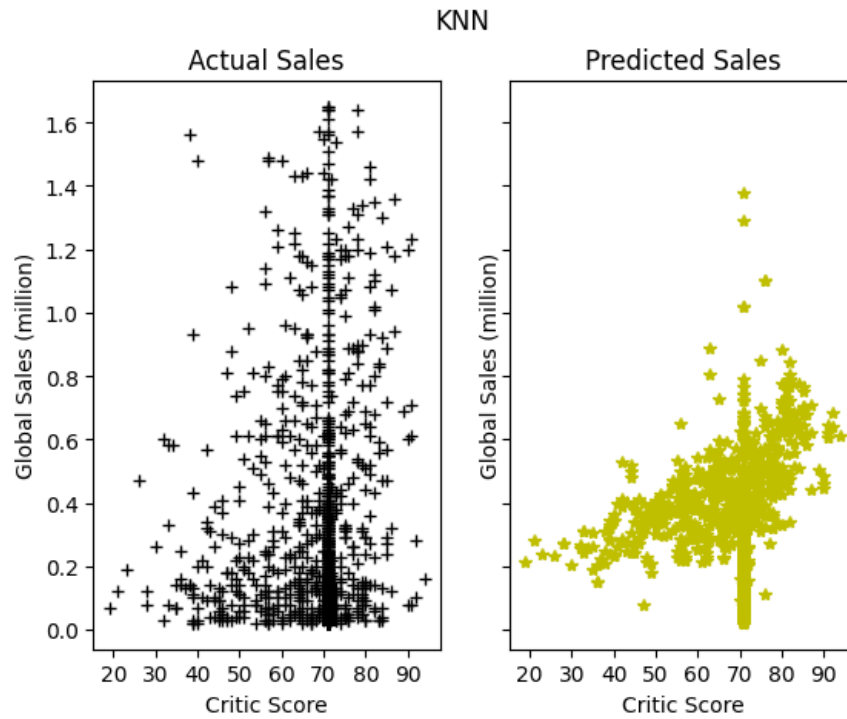


Figure 4: k-Nearest Neighbors Predictions

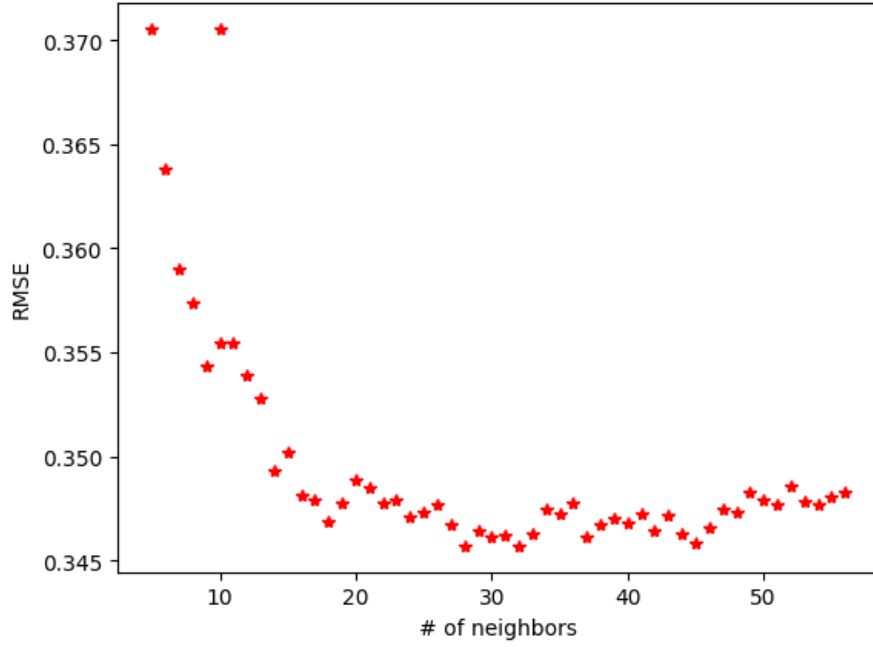


Figure 5: Choice of k neighbors for kNN

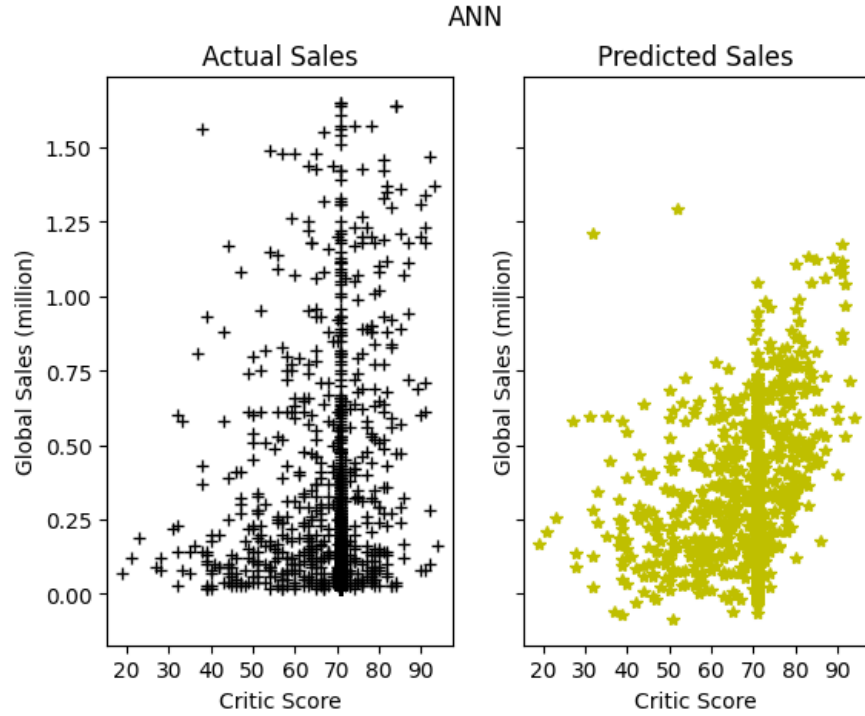


Figure 6: Neural Network Predictions

$$\begin{aligned}
 RMSE_{RF} &: 0.3242 \\
 RMSE_{kNN} &: 0.3289 \\
 RMSE_{ANN} &: 0.3363
 \end{aligned}$$

As we can see in Figure 3 and 4, the predicted sales roughly capture the general distribution of the global sales, but the one produced by k nearest neighbors regression is more localized than the one by Random Forest, which can be ascribed to the fact that the k -Nearest Neighbors regression only looks at the neighbors of the input for prediction. We plotted sales against Critic score to help visualize our predictions.

5.3 Discussion

From our RMSE values, we see that the strongest performer was Random Forest and the weakest was the ANN. Naturally, we will discuss the strengths and weaknesses of our models. For Random Forest, an ensemble of regression trees, the individual trees are easy to understand, specifically for mirroring human behavior purposes which is relevant to our task of suggesting that features impact decisions to purchase video games. A weakness of the neural network is its data volume demand with the relatively small dataset we are working with. In contrast with Random Forest, the model that the neural network produces is incomprehensible and a "black box" which makes it possibly less attractive for practical uses of sales prediction. However, as mentioned previously, there have been neural networks applied to PCA processed data for sales prediction. k -Nearest Neighbors falls in between the other two models, and we can see that the RMSE of k NN converges as we increase the number of k neighbors used.

6 Conclusion and Future Work

With a RMSE of around 0.32, our objective of modeling global sales of video games is decently accurate and can provide some insight for future video game releases. The range of global sales value is 0 to 60, so our error metric is actually very small and so our model can be considered very strong. With more video game data we can verify that our models are not overfitted, but ensemble methods and cross validation are innately preventative measures against overfitting. Out of the regressors that we explored, we conclude that the Random Forest model produced the best prediction of global video game sales by measure of RMSE, albeit by a small margin. To continue our exploration in the realm of video games, we see that some of the most successful and popular games today are free to play, offering paid in-game content that users can elect to pay for or not. Our project does not strongly consider this model of games and exploration of different payment structures may be interesting to consider in the future. Furthermore, as seen in related work, there are many esoteric features that were not taken into account. As with any consumer product, sentiment is a major factor in how a product is received and could also be a viable direction of exploration. Separate analyses of games and marketing strategies can be helpful to building a more holistic and complete model of global sales.

References

- [1] Anava, Oren, and Kfir Levy. "k*-nearest neighbors: From global to local." *Advances in neural information processing systems*. 2016.
- [2] Breiman, Leo. "Random forests." *Machine learning* 45.1 (2001): 5-32.
- [3] Marcoux, Julie, and Sid-Ahmed Selouani. "A hybrid subspace-connectionist data mining approach for sales forecasting in the video game industry." *2009 WRI World Congress on Computer Science and Information Engineering*. Vol. 5. IEEE, 2009.
- [4] Near, Christopher E. "Selling gender: Associations of box art representation of female characters with sales for teen-and mature-rated video games." *Sex roles* 68.3 (2013): 252-269.
- [5] Ruohonen, Jukka, and Sami Hyrynsalmi. "Evaluating the use of internet search volumes for time series modeling of sales in the video game industry." *Electronic Markets* 27.4 (2017): 351-370.
- [6] Wikipedia contributors. "Random subspace method." *Wikipedia, The Free Encyclopedia*. Wikipedia, The Free Encyclopedia, 3 Nov. 2019. Web. 27 Jul. 2020.

```

173 """
174
175 == Filename: helper.py
176 == Author: Yi Lyu
177 == Status: Complete
178
179 """
180
181 import numpy as np
182 import pandas as pd
183 import pickle
184 import sqlite3
185 import re
186 import os
187 from sklearn.impute import SimpleImputer as Imputer
188 from sklearn.impute import KNNImputer
189 from sklearn.preprocessing import LabelEncoder
190
191 __all__ = [ 'Videogames' ]
192
193 def get_dir(path):
194     return os.path.join(getWorkDir(), path)
195
196 def getWorkDir():
197     pathlist = os.path.abspath(os.curdir).split('/')
198     path = ''
199     for p in pathlist:
200         path = os.path.join(path, p)
201         if p == 'video-game-sales-predictor' or p == 'video-game-sales-
202 predictor-master':
203             break
204     return path
205
206 class Videogames(object):
207     def __init__(self, database_dir, data_dir='data/', storage='data'):
208         self.database_dir = database_dir
209         self.table = ''
210         self.data_dir = data_dir
211         self.storage = '{0}.pickle'.format(storage)
212
213         self._has_data = False
214         self._headers = []
215         self._dtypes = []
216         self._connection = None
217         try:
218             with open(get_dir(data_dir + self.storage), "rb") as f:
219                 self.table, self._headers, self._dtypes, self._has_data =
220 pickle.load(f)
221         except:
222             pass
223
224     @property
225     def table_name(self):
226         return self.table
227
228     @property
229     def status(self):
230         return self.get_status()
231
232     @property
233     def headers(self):
234         return self._headers
235
236     @property
237     def dtypes(self):

```



```

238         return self._dtypes
239
240     def get_status(self):
241         return self._has_data
242
243     def get_headers(self):
244         return self._headers
245
246     def get_dtypes(self):
247         return self._dtypes
248
249     def read_data_in(self, filepath, table, write_headers=False):
250         conn = sqlite3.connect(database=self.database_dir)
251         cur = conn.cursor()
252         data = pd.read_csv(filepath, delimiter=";", encoding="
253 unicode_escape")
254
255         self.table = table
256         headers = self._get_headers(data)
257         dtypes = self._get_dtypes(data)
258         self._create_table(headers, dtypes, cur)
259
260         if write_headers:
261             with open(get_dir(self.data_dir + 'headers.csv'), "w+") as f:
262                 f.write(";\n".join(headers))
263
264         if not self._has_data:
265             data = self._remove_missing(data)
266             with open(get_dir(self.data_dir + self.storage), "wb+") as f:
267                 self._insert_data(data, headers, dtypes, cur)
268                 self._has_data = True
269                 pickle.dump((self.table, headers, dtypes, True), f,
270 pickle.HIGHEST_PROTOCOL)
271
272         del data
273         conn.commit()
274         conn.close()
275
276     def _remove_missing(self, data):
277         data = data.replace(r'tbd', np.nan, regex=True)
278         data['User_Score'] = data['User_Score'].astype(np.float64)
279         condition = (data['Platform'].notnull() & data['Genre'].notnull()
280 & data['Publisher'].notnull() & data['Year_of_Release'].notnull())
281         data = self._imputation(data[condition])
282         return data
283
284     def _imputation(self, data):
285         imp = Imputer(strategy='median')
286         attributes = ['Critic_Score', 'User_Score', 'Critic_Count', '
287 User_Count']
288         for item in attributes:
289             data[item] = imp.fit_transform(data[[item]]).ravel()
290         return data
291
292     def get_col(self, *header):
293         if not self._connection:
294             self._connection = sqlite3.connect(self.database_dir)
295             cur = self._connection.cursor()
296
297             command = "SELECT {0} FROM {1};".format(self._list2str(header),
298 self.table)
299             return self._col2list(cur.execute(command).fetchall())
300
301     def execute(self, command):
302         if not self._connection:

```

```

303         self._connection = sqlite3.connect(self.database_dir)
304         cur = self._connection
305
306         if bool(re.match("[\t\n]*SELECT", command, re.I)):
307             return list(self._col2list(cur.execute(command).fetchall()))
308         else:
309             print("ILLEGAL COMMAND")
310
311
312     ## Helper Functions ##
313     def _get_headers(self, data):
314         """Return the headers of the data
315
316         Args:
317             data DataFrame: the data we read from csv.
318
319         Returns:
320             list: the headers of the data
321         """
322         if not self._headers:
323             self._headers = list(map(lambda col: col.lower(), data.
324 columns))
325         return self._headers
326
327     def _get_dtypes(self, data):
328         if not self._dtypes:
329             self._dtypes = [self._process_dtype(data[col][0]) for col in
330 data.columns]
331         return self._dtypes
332
333     def _create_table(self, headers, dtypes, cur):
334         """Execute the following SQL command
335
336         CREATE TABLE IF NOT EXISTS {table} (
337             name VARCHAR(80),
338             ...
339         );
340
341         Args:
342             headers (list): the list of columns where each header is
343 lowercase.
344             dtypes (list): the list of types where each type is either
345 NUMBER or VARCHAR(80) based on this data set.
346             cur (sqlite3.connection.cursor): a connection cursor of
347 sqlite3 database
348         """
349         command = "CREATE TABLE IF NOT EXISTS {0} ("
350         template = "{0} {1}"
351
352         n = len(headers)
353
354         ## Convert the data to suitable form for _list2str function
355         data = [template.format(headers[i], dtypes[i]) for i in range(n)]
356
357         command += self._list2str(data)
358         command += ");"
359         cur.execute(command)
360
361     def _insert_data(self, data, headers, dtypes, cur):
362         command_template = "INSERT INTO {0} ({1}) VALUES ({2});"
363         for i, itr in data.iterrows():
364             res = list(map(self._str_classifier, list(itr)))
365             command = command_template.format(self.table, ", ".join(
366 headers),

```

```

367                                     self._list2str(res ,
368 classify=self._row_classifier(res , dtypes)))
369                                     cur.execute(command)
370
371 def _list2str(self , data , delimiter=",", classify=lambda x: x):
372     """Convert the list to a string
373
374     I have not found such a function in Python and therefore
375     wrote one.
376
377     Args:
378         data (list): the row of the table
379         delimiter (str, optional): the delimiter.
380         classify (function, optional): a function that classifies the
381         data in the row.
382
383     Returns:
384         str: a string representing the data converted to a string.
385     """
386     res = ""
387     for i in range(len(data)):
388         res += classify(data[i])
389         if i != len(data) - 1:
390             res += delimiter + " "
391     return res
392
393 def _row_classifier(self , data , dtypes):
394     """classifies the data in a row in the table
395     def classifier(x):
396         i = data.index(x)
397         if dtypes[i] == "NUMBER":
398             if x == "NULL" or x == 'tbd':
399                 return "-1"
400             else:
401                 return str(x)
402         else:
403             return "\">{0}\{}".format(x)
404     return classifier
405
406 def _str_classifier(self , x):
407     """classifies the data so that it does not contain nan
408     if type(x) == float and np.isnan(x):
409         return -1
410     return x
411
412 def _process_dtype(self , var):
413     dtype = type(var)
414     if dtype == str and var.isnumeric():
415         return "NUMBER"
416     type_converter = {type(',') : "VARCHAR(80)", np.float64: "NUMBER"
417 , np.int64: "NUMBER"}
418     return type_converter[dtype]
419
420 def _col2list(self , col):
421     n = len(col[0])
422     return list(map(lambda x: list(x)[:n] , col))
423
424
425 == Filename: main.py
426 == Author: Yi Lyu
427 == Status: Complete
428
429
430

```

```

431 import numpy as np
432 import matplotlib.pyplot as plt
433 import pandas as pd
434 import seaborn as sns
435 import pickle
436 import os
437 from sklearn.decomposition import PCA
438 from keras.models import load_model
439 from sklearn.model_selection import train_test_split
440
441 from helper import Videogames, getWorkDir, get_dir
442 from models import *
443 from plotting import *
444
445 def read_data():
446     videogames = Videogames(get_dir("data/math156.db"))
447     videogames.read_data_in(get_dir("data/videogames.csv"), "VIDEOGAMES",
448                             True)
449     res = np.array(videogames.execute('''
450         SELECT name, g_total, cscore, uscore, genre, publisher FROM (
451             SELECT name AS name,
452                   SUM(global_sales) AS g_total,
453                   critic_score AS cscore,
454                   user_score AS uscore,
455                   genre AS genre,
456                   publisher AS publisher
457             FROM VIDEOGAMES
458             WHERE year_of_release >= 2004 and uscore != 0 and cscore != 0
459             GROUP BY name) AS VideogameSummary
460         WHERE g_total != 0
461         ORDER BY g_total DESC;
462     '''))
463     return res
464
465 if __name__ == "__main__":
466     ## the critic scores and user scores
467     columns = ['name', 'gtotal', 'cscore', 'uscore', 'genre', 'publisher',
468               ]
469     res = pd.DataFrame(read_data(), columns=columns)
470
471     n = len(res)
472     factor = 0.1
473     quantile1 = round(n * factor)
474     quantile2 = n - round(n * factor)
475     res = res.loc[quantile1:quantile2 + 1, :]
476
477     ## Transform data into appropriate form for regression
478     scores = res[['cscore', 'uscore']]
479     genre = pd.get_dummies(res['genre'], drop_first=True)
480     publisher = pd.get_dummies(res['publisher'], drop_first=True)
481
482     X = pd.concat((scores, genre, publisher), axis=1).astype('float64')
483     Y = res['gtotal'].astype('float64')
484
485
486     X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=
487     .20, train_size=.80, random_state = 40)
488
489     try:
490         with open(get_dir('data/models.pickle'), 'rb') as f:
491             rfreg, knnreg = pickle.load(f)
492             annreg = load_model(get_dir('data/ann'))
493     except:
494         annreg = ann(X_train, Y_train.ravel()) ## ANN

```

```

495         rfregr = random_forest(X_train, Y_train.ravel())    ## Random
496 Forest
497         knnregr = knn(X_train, Y_train.ravel())             ## KNN
498         with open(get_dir('data/models.pickle'), 'wb+') as f:
499             pickle.dump((rfregr, knnregr), f, pickle.HIGHEST_PROTOCOL)
500             annregr.save(get_dir('data/ann'))
501
502         print("The mean is", np.mean(Y))
503     ## RMSE
504     rmse_template='RMSE\t{name:25}{value:18}'
505     print(rmse_template.format(name='random forest', value=rmse(X_test,
506 Y_test, rfregr)))
507     print(rmse_template.format(name='Knn', value=rmse(X_test, Y_test,
508 knnregr)))
509     print(rmse_template.format(name='ANN', value=rmse(X_test, Y_test,
510 annregr)))
511
512     plot_predictions(X_test, Y_test, rfregr, knnregr, annregr)
513
514     """
515
516     print(=====)
517
518     ## R2
519     r2_template = 'R^2\t{name:25}{value:18}'
520     print(r2_template.format(name='random forest', value=rfregr.score(
521 X_test, Y_test)))
522     print(r2_template.format(name='Knn', value=knnregr.score(X_test,
523 Y_test)))
524     print(r2_template.format(name='Aan', value=annregr.score(X_test,
525 Y_test)))
526
527     """
528
529 """
530 =====
531 == Filename: models.py
532 == Author: Yi Lyu
533 == Status: Complete
534 =====
535 """
536
537 import numpy as np
538 import matplotlib.pyplot as plt
539 from sklearn.linear_model import Ridge
540 from sklearn.linear_model import GammaRegressor
541 from sklearn.preprocessing import PolynomialFeatures, StandardScaler
542 from sklearn.pipeline import make_pipeline
543 from sklearn.ensemble import RandomForestRegressor
544 from sklearn.neighbors import KNeighborsRegressor
545 from sklearn.model_selection import train_test_split, cross_val_score
546 from sklearn.metrics import mean_squared_error
547 from keras.models import Sequential
548 from keras.layers import Dense
549
550 ## just in case someone wants to implement them instead of using sklearn
551
552 def rmse(X_test, Y_test, model):
553     Y_pred = model.predict(X_test)
554     return mean_squared_error(Y_test, Y_pred, squared=False)
555
556 def plot_knn(ns, rmses):
557     plt.plot(ns, rmses, 'r*')
558     plt.xlabel('# of neighbors')

```

```

559 plt.ylabel('RMSE')
560
561 plt.savefig('graphs/knn_choice_n.png', bbox_inches='tight')
562 plt.clf()
563
564 def knn(xs, ys, n=10):
565     X_train, X_test, Y_train, Y_test = train_test_split(xs, ys, test_size
566     = .1, random_state = 40)
567     num_cols = len(X_train.columns)
568     i = 5
569
570     best_index = 4
571     best_score = 10000
572     nums = [i for i in range(5, int(np.sqrt(num_cols)) + 5)]
573     cvs = []
574
575     for num in nums:
576         model = KNeighborsRegressor(n_neighbors=num, algorithm='kd_tree',
577         weights='distance')
578         temp = cross_val_score(model, xs, ys, cv=5).mean()
579         temp = np.sqrt(1 - temp)
580         if temp < best_score:
581             best_score = temp
582             best_index = num
583     print(best_index)
584     return KNeighborsRegressor(n_neighbors=best_index, algorithm='kd_tree'
585     ', weights='distance').fit(xs, ys)
586     """
587
588     best_model = KNeighborsRegressor(n_neighbors=i, algorithm='kd_tree',
589     weights='distance').fit(X_train, Y_train)
590     best_rmse = rmse(X_test, Y_test, best_model)
591
592     ### Cross Validation
593     ns = [n]
594     rmses = [best_rmse]
595     cvs = []
596     ### You can change 5 to * 2 or * 3 here for a better result , but
597     slower.
598     for n in range(i, int(np.sqrt(num_cols)) + 5):
599         model = KNeighborsRegressor(n_neighbors=n, algorithm='kd_tree',
600         weights='distance').fit(X_train, Y_train)
601         temp = rmse(X_test, Y_test, model)
602         ns.append(n)
603         rmses.append(temp)
604         if temp < best_rmse:
605             best_model = model
606             best_rmse = temp
607     plot_knn(ns, rmses)
608
609     """
610
611     return best_model
612
613 def ann(xs, ys):
614     n = len(xs.columns)
615     ANN = Sequential()
616     ANN.add(Dense(units = 6, activation = "elu", input_dim = n))
617     ANN.add(Dense(units = 4, activation = "elu"))
618     ANN.add(Dense(units = 1))
619
620     ANN.compile(optimizer = "adam", loss = "mean_squared_error")
621     ANN.fit(xs, ys, batch_size = 2, epochs = 100)
622     return ANN
623

```

```

624 def gamma_model(xs, ys):
625     model = GammaRegressor().fit(xs, ys)
626     return model
627
628 def linear_model(xs, ys, m):
629     model = make_pipeline(PolynomialFeatures(m), Ridge(normalize=True)).
630     fit(xs, ys)
631     return model
632
633 def random_forest(xs, ys):
634     model = RandomForestRegressor(criterion='mse').fit(xs, ys)
635     return model
636
637 """
638 == Filename: plotting.py
639 == Author: Yi Lyu
640 == Status: Complete
641 """
642
643
644 import numpy as np
645 import matplotlib.pyplot as plt
646 import pandas as pd
647 import seaborn as sns
648 import os
649
650 from helper import getWorkDir, get_dir
651
652 def predict(X_test, Y_test, model):
653     """Predict the sales based on the dataset
654
655     Args:
656         X_test (DataFrame): Data
657         Y_test (Series): Actual Sales
658         model (object): Model we are using
659
660     Returns:
661         DataFrame: predicted scales
662     """
663     return pd.DataFrame(model.predict(X_test))
664
665 def plot_helper(xs, data_ys, predict_ys, model_name='Unknown'):
666     """Plot the predicted sales
667
668     Args:
669         xs (Series): the x values
670         data_ys (Series): the actual sales
671         predict_ys (Series): the predicted sales
672         model_name (str, optional): the name of the model. Defaults to '
673         Unknown'.
674     """
675     xs = xs.astype(np.float64)
676     fig, (ax1, ax2) = plt.subplots(1, 2, sharex=True, sharey=True)
677     plt.xticks(np.linspace(0, 100, 11))
678
679     fig.suptitle(model_name)
680
681     ax1.plot(xs, data_ys, 'k+', label='data')
682     ax1.set_title('Actual Sales')
683     ax1.set_xlabel='Critic Score', ylabel='Global Sales (million)'
684
685     ax2.plot(xs, predict_ys, 'y*', label='prediction')
686     ax2.set_title('Predicted Sales')
687     ax2.set_xlabel='Critic Score', ylabel='Global Sales (million)'

```

```

688     pic_path = 'graphs/{0}.png'.format(model_name.replace(' ', '_').lower())
689     ()
690
691     pic_dir = get_dir(pic_path)
692
693     plt.savefig(pic_dir, bbox_inches='tight')
694     plt.clf()
695
696
697 def plot_helper2(data_ys, predicted_ys, model_name='Unknown'):
698     fig, (ax1, ax2) = plt.subplots(1, 2, sharex=True, sharey=True)
699     fig.suptitle(model_name)
700
701     bins = np.arange(0, 6, 0.1)
702     sns.distplot(data_ys, bins=bins, hist=True, kde=True, ax=ax1, color='r',
703                 axlabel='Sales')
704     ax1.set_title('Actual Sales -- Density Plot')
705     ax1.set_xlim(0, 2)
706     sns.distplot(predicted_ys, bins=bins, hist=True, kde=True, ax=ax2,
707                 color='b', axlabel='Sales')
708     ax2.set_title('Predicted Sales -- Density Plot')
709     ax2.set_xlim(0, 2)
710
711     pic_path = 'graphs/{0}_hist.png'.format(model_name.replace(' ', '_').lower())
712     pic_dir = get_dir(pic_path)
713
714     plt.savefig(pic_dir, bbox_inches='tight')
715     plt.clf()
716
717
718 def plot_predictions(X_test, Y_test, rfregr, knnregr, annregr):
719     """Plot the Predicted sales of each model
720
721     Args:
722         X_test (DataFrame): data
723         Y_test (Series): actual sales
724         rfregr (RandomForestRegressor): Random Forest Regressor
725         knnregr (KNNRegressor): KNN Regressor
726         annregr (ANNRegressor): Artificial Neural Network Regressor
727
728     """
729     cscores = X_test['cscore']
730     ## Get predicted sales
731     rfres = predict(X_test, Y_test, rfregr)
732     knnres = predict(X_test, Y_test, knnregr)
733     annres = predict(X_test, Y_test, annregr)
734
735     ## Correct the indices in case
736     temp = pd.DataFrame(pd.concat([cscores, Y_test], axis=1).to_numpy(),
737                         columns=['cscore', 'gtotal'],
738                         index=np.arange(0, len(cscores), 1))
739
740     ## Create a pandas DataFrame sorted by Critic Score
741     df = pd.concat([temp, rfres, knnres, annres], axis=1)
742     df = pd.DataFrame(df.sort_values(by='cscore', ascending=True).to_numpy(),
743                     columns=['cscore', 'gtotal', 'rfres', 'knnres', 'annres'])
744
745     plot_helper(df['cscore'], df['gtotal'], df['rfres'], 'Random Forest')
746     plot_helper(df['cscore'], df['gtotal'], df['knnres'], 'KNN')
747     plot_helper(df['cscore'], df['gtotal'], df['annres'], 'ANN')
748
749     plot_helper2(df['gtotal'], df['rfres'], 'Random Forest')
750     plot_helper2(df['gtotal'], df['knnres'], 'KNN')
751     plot_helper2(df['gtotal'], df['annres'], 'ANN')
752

```