

Proyecto IA

Inteligencia Artificial

Departamento de Electrónica

Profesor: Francisco Carlos Calderón Bocanegra



Video Games Sales Prediction

Facultad de Ingeniería, Pontificia Universidad Javeriana, Bogotá D.C, Colombia.

Sergio Enrique González. sergioegonzalez@javeriana.edu.co

En este proyecto se encontrará la predicción de ventas de videojuegos en Estados Unidos basado en el dataset obtenido en la página de Kaggle. Este dataset se puede descargar y encontrar en el siguiente enlace: <https://www.kaggle.com/datasets/migeruj/videogames-predictive-model>

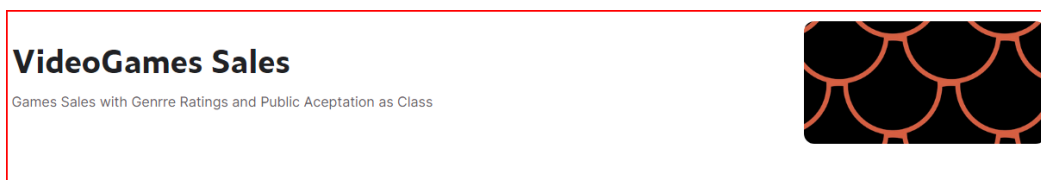


Ilustración 1. Dataset obtenido en Kaggle

Este proyecto utilizo la herramienta Visual Studio Code con un lenguaje de programación en Python, en donde se realizaron diferentes tipos de pruebas y modelos para obtener el mejor método de análisis el cual este problema de aprendizaje correspondía a uno de no supervisado y de regresión.

1. El archivo fue creado en .ipynb
2. Se procedió hacer la limpieza respectiva del dataset.
3. Se normalizaron los datos donde posteriormente se realizó PCA
4. Al comprobar los datos obtenidos en PCA fue necesario utilizar los 5 pliegues ya que cada uno de los componentes de la variable X aportan información importante en la predicción que se va a desarrollar durante el código.

PCA

```
PCA = decomposition.PCA(n_components=5)
PCA.fit(X_train)
X_train = PCA.transform(X_train)

print("Pesos n de PCA:", PCA.explained_variance_ratio_)
print("Suma total de pesos n de PCA", sum(PCA.explained_variance_ratio_))
```

```
Pesos n de PCA: [0.24404026 0.21613138 0.20139132 0.17850932 0.15992771]
Suma total de pesos n de PCA 0.9999999999999999
```

Ilustración 2. PCA de la variable X

- Después se realizó un gridsearch y cross-validación utilizando GridSearchCV para optimizar los hiperparámetros de los diferentes modelos para R2, arrojando los mejores parámetros del modelo seleccionado para obtener el mejor Score.
- Los modelos utilizados fueron KNN donde se varió la métrica y peso de la distancia, Máquina de vectores de soporte "SMV" y para finalizar se utilizó Decision Tree se varió el criterio y profundidad de los árboles

SVM

Optimizar los hiperparámetros utilizando Grid-search

```
param_grid=[{'kernel':['linear','poly','rbf'],'gamma':['scale','auto'],'C':[0.1,1,2.5,5,7.5,10,30,50,75,100]]
gridsearchCV=GridSearchCV(estimator=SVR(),param_grid=param_grid, scoring="r2")
gridsearchCV.fit(X_train,y_train)
```

```
[23]
...
> GridSearchCV
> estimator: SVR
  > SVR
```

Se muestran los resultados

Los mejores parámetros para obtener el mejor score

```
gridsearchCV.best_params_
```

```
[22]
... {'C': 100, 'gamma': 'auto', 'kernel': 'rbf'}
```

Resultado R2

```
gridsearchCV.best_score_
```

```
[23]
... 0.2324619900598311
```

Ilustración 3. Máquina de vectores de Soporte con GridSearch

KNN

```
param_grid=[{'n_neighbors':[1,2,3,4,5,6,7,8,9,10,20,30,40,50,60,70,80,90,100],'algorithm':['auto','ball_tree','kd_tree','brute'],'weights':['uniform','distance']}
gridsearchCV=GridSearchCV(estimator=KNeighborsRegressor(),param_grid=param_grid, scoring="r2")
gridsearchCV.fit(X_train,y_train)
```

```
[1]
> GridSearchCV
> estimator: KNeighborsRegressor
  > KNeighborsRegressor
```

Se muestran los resultados

Los mejores parámetros para obtener el mejor score

```
gridsearchCV.best_params_
```

```
[1]
{'algorithm': 'brute', 'n_neighbors': 20, 'weights': 'distance'}
```

Resultado R2

```
gridsearchCV.best_score_
```

```
[3]
0.3758430078582195
```

Ilustración 4. K-Vecinos más cercanos con GridSearchCV

Decision Tree

```
param_grid=[{'criterion': ["squared_error","friedman_mse","poisson"],'max_depth': range(1,10),'min_samples_split': range(3,10),'min_samp:
gridsearchCV=GridSearchCV(estimator=DecisionTreeRegressor(),param_grid=param_grid, scoring="r2")
gridsearchCV.fit(X_train,y_train)
```

[27] Python

GridSearchCV

- estimator: DecisionTreeRegressor
 - DecisionTreeRegressor

Se muestran los resultados

Los mejores parametros para obtener el mejor score

```
gridsearchCV.best_params_
```

[28] Python

```
{'criterion': 'squared_error',
 'max_depth': 7,
 'min_samples_leaf': 3,
 'min_samples_split': 5}
```

Resultado R2

+ Código + Markdown

```
gridsearchCV.best_score_
```

[29] Python

```
0.27058275134148896
```

Ilustración 5. Decision Tree con GridSearchCV

Los datos proporcionados en el conjunto de datos están bastante descorrelacionados entre sí, para obtener un mejor resultado el conjunto de datos puede mejorar si desea analizarlo mediante de un modelo de aprendizaje automático.

Enlace para video con explicación más detallada:

<https://www.youtube.com/watch?v=9gelV7wpy8Q>