

Predictor automatizado para la selección del entretenimiento electrónico

El cine, los libros e incluso los videojuegos pueden ser considerados como la octava maravilla del mundo. Queremos dar, por tanto, una recomendación fiel a los gustos del usuario para su disfrute superlativo.

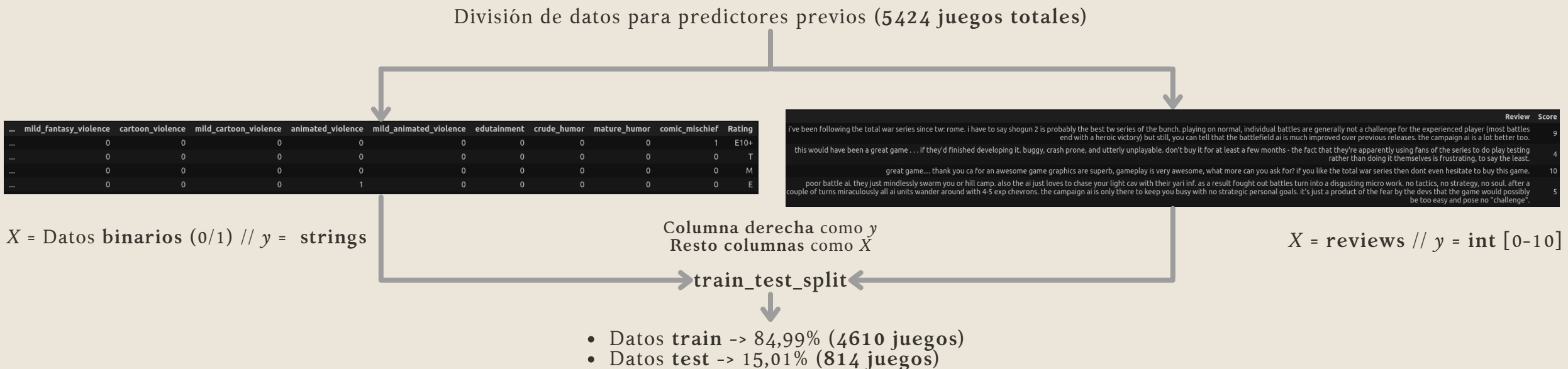
Conjunto de datos

	Nombre	Fecha Salida	Genero	Plataforma	Portable	Desarrolladora	Rating ESRB	metaScore	userScore	Multijugador	Online
0	The Legend of Zelda: Skyward Sword	2011	general	[Wii]	0	Nintendo	E10+	93.00	8.1	0	0
1	Total War: Shogun 2	2011	general	[PC]	0	Creative Assembly	T	90.00	8.4	1	1
2	Crysis	2007	action	[PC, PlayStation 3, Xbox 360]	0	Crytek	M	84.33	7.6	0	1
3	Allegiance	2000	sci-fi	[PC]	0	Microsoft Game Studios	E	86.00	7.5	1	0

He aquí una muestra de los datos empleados. Cada uno de estos consta de quince atributos. Trataremos de predecir el nombre del juego en cuestión. No obstante, contamos con muchos datos a tratar.

Por lo tanto, necesitamos realizar unas aclaraciones previas:

- La mayoría de variables, **genero**, **portable**, **online**, ..., son binarias, fácilmente tratables.
- Existen una serie de variables las cuales provienen como resultado de otras distintas, estas son el **Rating ESRB** y ambos **Scores**. Nos centraremos en clasificadores especializados para estas dos variables particulares.



Predictor ESRB

Usamos **GridSearchCV** para encontrar los mejores hiper-parámetros y predictores. Se usó **StratifiedKFold** para la división de *train* y *validación* dentro (2 splits).

Utilizamos todo tipo de clases de predictores:

- Supervisados:**
 - Clasificación:** Regresión Logística, Naive Bayes, Red Neuronal, SVMs, Ensembles.
 - Regresión:** Regresión lineal.
- No Supervisados:**
 - Kmeans.

Resultados obtenidos

Predictor	Accuracy
0	svm 0.706388
1	ovo 0.706388
2	randomForest 0.705160
3	bagging 0.703931
4	adaboost 0.701474
5	decisionTree 0.696560
6	ova 0.681818
7	regressionLogistica 0.675676
8	redNeuronal 0.662162
9	nb 0.624079
10	kmeans 0.276413
11	regressionLineal 0.199017

Hiper-parametros SVM
• C: 1
• gamma: 5
• kernel: rbf
• decision_function_shape: ovo

- Los hiper-parámetros obtenidos en **SVM** indican que se intenta **obtener una varianza baja**.
- El **kernel** utilizado, **rbf**, se utiliza para poder ajustar al máximo los parámetros.

Hiper-parametros RandomForest
• n_estimators: 20
• criterion: gini
• min_samples_leaf: 2
• min_samples_split: 10
• max_features: None
• class_weight: None
• bootstrap: True

- Los hiper-parámetros obtenidos en el mejor **Ensemble** indican que se necesitan **varios datos iniciales** para crear unos árboles decentes.
- Un claro uso del **bootstrap** para poder crear un predictor que **no sobre-aprenda**.

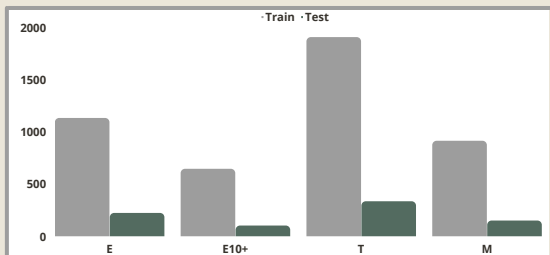
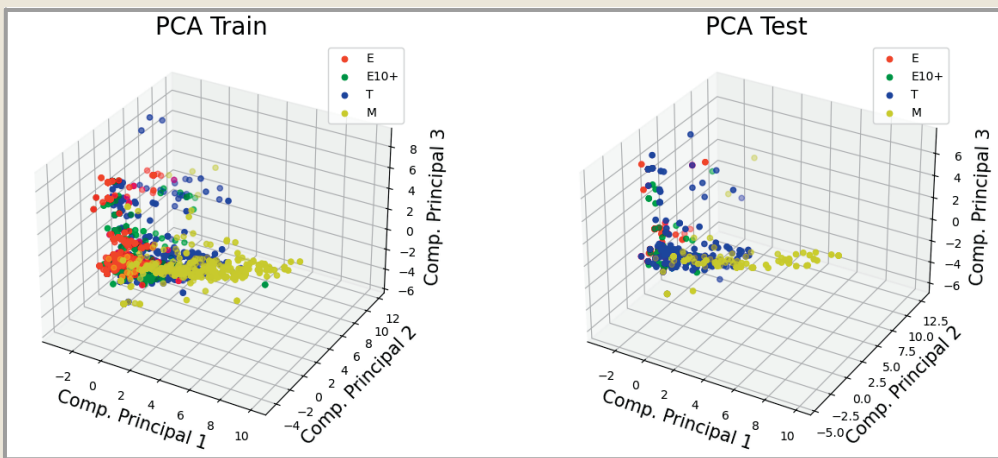
Los videojuegos son usados como hobby o profesión por todo tipo de personas. Es muy importante que cada uno de ellos tenga la mejor experiencia posible y se ajuste a las preferencias personales.

¿Seremos capaces de predecir correctamente?

Datos a utilizar

Existen un total de 48 características. Haremos un tratamiento previo:

- Escalamiento** para tener datos en rangos parecidos.
- Reducción de dimensionalidad **PCA** hasta una dimensión representable.



Se observan varios datos:

- Clases E y M claramente diferenciadas
- Clase T en punto medio
- Clase E10+ prácticamente indiferenciable



Predictor Reviews

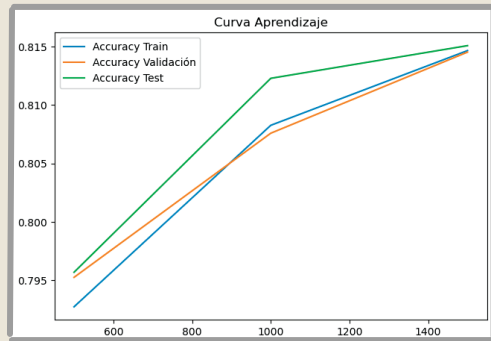
Bag of Words

Usamos un predictor basado en Naive Bayes Multinomial, siguiendo el algoritmo *Bag of Words*

Necesitamos extraer las palabras más representativas, los **tokens** a analizar. Usaremos la función **CountVectorizer** y **TfidfVectorizer**.

Hiper-parametros vectorizer
• strip_accents: ascii
• analyzer: word
• max_features: [500, 1000, 1500]
• binary: True

Realizamos 3 iteraciones distintas para saber cual es el mejor número de palabras para extraer.



Se observan los siguientes datos:

- Con el aumento de atributos, aumenta también el acierto
- No es un porcentaje de cambio tan drástico, pero al ser datos grandes, la diferencia entre la primera aproximación y la última en cuanto al conjunto *Test*, es de un aumento de acierto de 551 review

La normalización se ha realizado por lo siguiente:

- Son valoraciones subjetivas por lo que no existe una distinción clara entre clases.
- Al ser 10 notas (clases) distintas, la diferencia entre valorar un juego como 8 o 9 es muy pequeña, por lo que consideraremos como correcta una predicción que se desvíe como mucho dos puntos de la nota real.

De los resultados podemos decir lo siguiente:

- Existe un mayor porcentaje de acierto en los extremos.
- Por otro lado, aquellos que se encuentran en "puntos muertos", como valoraciones entre el 3 y el 7 devuelven una peor precisión. Podemos achacar esto a la complejidad subjetiva del problema.

Conclusión

KMeans

Tras haber entrenado los predictores previos creamos el definitivo, recomendador de videojuegos.

Hiper-parametros SVM
• n_clusters: len(listaJuegos)/2
• init: k-means++
• algorithm: lloyd

Se utiliza un clasificador no supervisado ya que no nos interesa obtener un "porcentaje de acierto", sino tener juegos clasificados.

- Utilizamos la mitad de clusters de la cantidad total de juegos.
- Usamos un algoritmo como **KMeans** para poder obtener recomendaciones. Obtenemos el cluster al que más cerca se encuentre y damos los juegos que se encuentren en ese cluster como recomendaciones.

Características X:

añoLanzamiento / genero / portable / característicasESRB / review / multijugador / online

Ejemplo a recomendar: [2002, 'fantasy', 1, 'T', 10, 1, 0]

- Ambos valores han sido dados con los predictores anteriores.

Recomendado:

Nombre	Plataforma	Desarrolladora	Grupo
101	Play Station	Predator 2 Games, Inc. Game Boy Advance, Nintendo Wii	Predator 2 Games, Inc.

Líneas Futuras

- Viendo el porcentaje de acierto obtenido en el **predictor ESRB**, necesitaríamos de alguna forma mejorar este número. Se podría realizar mediante la eliminación de algunas características o su unión en grupos. Por otro lado, hacer otro tipo de estandarización/normalización y/o utilizar otro tipo de reducción de dimensión.
- La clase **E10+** ha demostrado ser especialmente conflictiva, afectando al rendimiento de todo el sistema. Concluimos con el ya mencionado problema del mundo real, es una clase que de forma instintiva se solapa con las demás, produciendo fallos.
- Por último, en el **Bag of Words**, podríamos hacer un preprocesamiento distinto u otro tipo de predictor no basado en porcentajes para mejorar la precisión. Para ello se debería usar una menor cantidad de ejemplos para no realizar un sobre-aprendizaje.