

Implementación de una técnica de aprendizaje máquina sin el uso de un framework.

Imanol Muñiz Ramirez A01701713

Contents

Introducción	2
Contexto	2
Problema	2
Objetivo	2
Solución.....	2
Extracción de los datos	3
Transformación	3
Carga.....	4
Algoritmo	5
Resultados y conclusiones	5
Cambio de dataset.....	10
Extracción.....	11
Transformación	11
Carga.....	11
Resultados.....	11
Conclusiones.....	15
Modelo con framework.....	15
Marco teórico.....	15
Dataset.....	16
Resultados.....	16
Conclusiones.....	19

Introducción

Contexto

Teamfight Tactics es un juego online de 8 jugadores que consiste en construir el equipo más fuerte para derrotar al de los demás. Cada ronda se simula el enfrentamiento y obtienes monedas de acuerdo con los resultados. Entre las rondas puedes comprar personajes llamados campeones y formar sinergias entre ellas. Existen campeones con costos que van desde una moneda hasta cinco, siendo generalmente los más costosos los que tienen un mayor potencial. Cada campeón además de contar con una habilidad única también cuenta con estadísticas base que van acordes a su costo y rol. Por ejemplo, un personaje diseñado para resistir el daño enemigo puede contar con más puntos de vida que el que está diseñado para hacer daño, pero también entre campeones con mismo rol pero diferente costo, suele tener mayores atributos el que cuesta más.

Problema

Riot Games frecuentemente está diseñando las próximas versiones de su juego Teamfight Tactics (TFT) cambiando campeones y mecánicas. Dado que cada campeón tiene atributos únicos, suele ser complicado ajustarlos para el costo en el que tiene que encajar o viceversa. A veces se lanza el nuevo set o versión con ciertos campeones que resultan injustamente más poderosos provocando desequilibrios que afectan la variedad de opciones y consecuentemente la jugabilidad.

Objetivo

Desarrollar una herramienta de machine learning que nos permita asignar el costo de un campeón con base en sus estadísticas base.

Solución

El ETL se realizó con el código del archivo Data_Transformation.py y los datos limpios se almacenan en el documento TFT_Champion_Transformed.csv. Esto se hizo con la intención de no ejecutar la limpieza de los datos cada que ejecutamos el archivo de machine learning. A continuación se especifica más a detalle esta etapa.

Extracción de los datos

Los datos se descargaron de Kaggle a través de la siguiente URL, compartidos por un miembro de la comunidad.

https://www.kaggle.com/datasets/gyejr95/league-of-legends-tftteamfight-tacticschampion?select=TFT_Champion_CurrentVersion.csv

Los datos son tabulares. Las columnas son principalmente numéricas, aunque también hay textos. Son 52 instancias, una por cada campeón.

```
name,cost,health,defense,attack,attack_range,speed_of_attack,dps,skill_name,skill_cost,origin,class
gangplank,5,1000,30,60,1,1.0,60,gangplank_orbitalstrike,100/175,Space Pirate,['Mercenary', 'Demolitionist']
graves,1,650,35,55,1,0.55,30,graves_smokegrenade,50/80,Space Pirate,['Blaster']
neeko,3,800,35,50,2,0.65,33,neeko_popblossom,75/150,Star Guardian,['Protector']
```

Transformación

Primero fue necesario eliminar las columnas que no eran relevantes o el tipo de dato no era manejable para este algoritmo. Eliminamos name, class, origin y skill_name.

En el caso de la columna skill_cost observamos que el tipo de dato es un string que contiene dos enteros separados por una diagonal. Esta columna representa con cuánta energía o maná inicia el combate y cuánta necesita el personaje para ejecutar su habilidad. Un campeón que requiera una cantidad muy alta de maná para lanzar una habilidad o que inicie el combate con muy poca cantidad, podría repercutir en qué tan poderosa es esa unidad. Por lo tanto, esta columna nos es relevante para nuestro objetivo. A partir de esta creamos dos columnas llamadas inicial_maná y skill_cost.

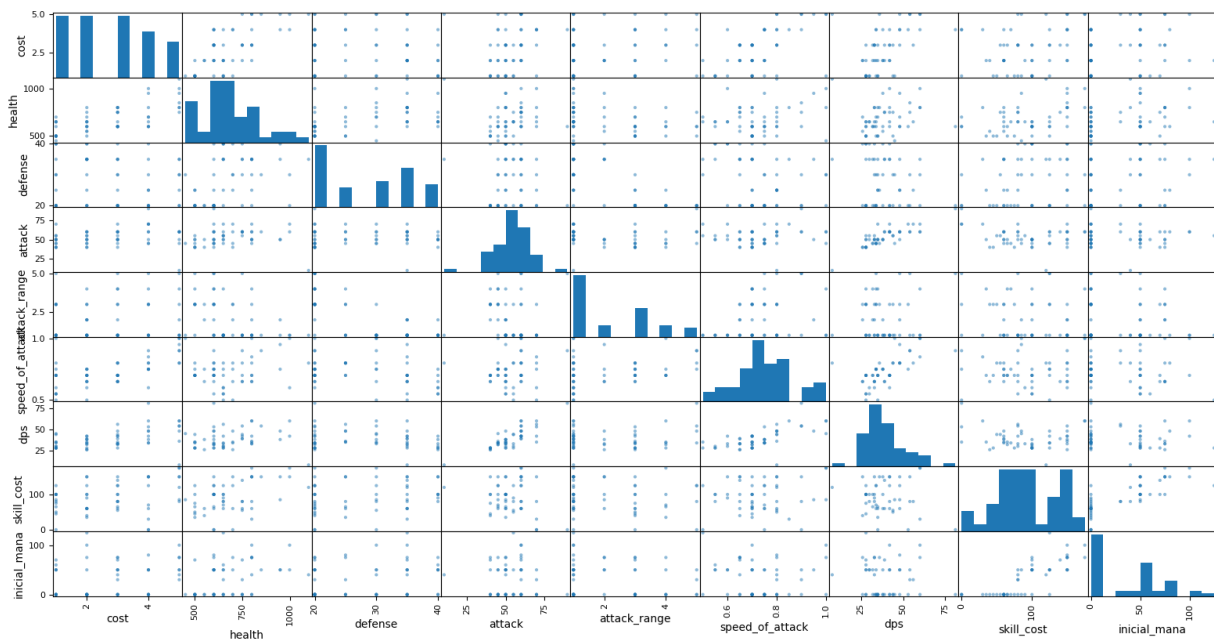
Revisando las instancias observamos que hay ciertos campeones que no utilizan maná para lanzar su habilidad. Por lo tanto, para estas unidades este atributo lo convertiremos a cero para que de esta forma solo repercuta el resto de sus estadísticas.

```
jhin,4,600,20,90,5,0.9,81,jhin_whisper,-,Dark Star,['Sniper']
jinx,4,600,20,70,3,0.75,53,jinx_getexcited,-,Rebel,['Blaster']
```

Por otra parte, nuestra columna objetivo (cost) la ponemos al final de la tabla para manejarla más fácilmente. Nos queda algo así:

	health	defense	attack	attack_range	speed_of_attack	dps	skill_cost	inicial_maná	cost
0	1000	30	60	1	1.00	60	175.0	100.0	5
1	650	35	55	1	0.55	30	80.0	50.0	1
2	800	35	50	2	0.65	33	150.0	75.0	3

Posteriormente realicé una matriz de dispersión para ver si alguna variable no mostraba una relación lineal y hacer el ajuste.



Del gráfico podemos identificar algunos patrones. Mientras en algunos hay una tendencia lineal clara (Ej. health vs cost) que hace notar que a mayor sea el costo del campeón, mayor serán sus puntos de vida, en otros pareciera una dispersión de los datos uniforme (Ej. defense vs cost), sin embargo, en el gráfico de defense vs health vemos que también a mayor defensa mayor vida lo que podría indicar que a mayor defensa también mayor el costo del personaje por lo que no es una columna que no esté aportando información.

Si observamos el conjunto de datos podemos observar un último problema. Existen atributos que pueden ser 1000 y otros que van entre 0 y 1. Al trabajar con valores grandes y potencias ocasiona desbordamientos de variables y que sea más complicado encontrar los hiper parámetros adecuados resultando en que el modelo no converja. Para solucionar este problema aplicamos un escalamiento a todas las variables para convertirlas en números entre 0 y 1. Al final obtuvimos una tabla con esta forma.

```
health,defense,attack,attack_range,speed_of_attack,dps,skill_cost,inicial_mana,cost
0.8461538461538463,0.5,0.625,0.0,1.0,0.7123287671232876,1.0,0.8,5
0.3076923076923077,0.75,0.5625,0.0,0.10000000000000009,0.3013698630136986,0.45714285714285713,0.4,1
0.5384615384615385,0.75,0.5,0.25,0.30000000000000004,0.3424657534246575,0.8571428571428571,0.6,3
```

Carga

Los datos originales utilizados se encuentran disponibles en el repositorio en el archivo TFT_Champion_CurrentVersion.csv: <https://github.com/ViejoAgrio/Machine-Learning>

Algoritmo

El archivo No_framework.py en la carpeta de scripts contiene el código que obtiene los parámetros “ b_i ” de una función para calcular el costo de cada campeón según sus estadísticas base. Esta función la obtenemos por el método de regresión lineal, utilizando como modelo la función $f(x) = b_n x_n + b_{n-1} x_{n-1} + \dots + b_1 x_1 + b_0$ donde $f(x)$ representa el costo del campeón y las x_i cada una de las estadísticas base de este. Para calcular el error primero utilizamos una función que convierte la evaluación de la función $f(x)$ a su entero más cercano (simulando una regresión logística múltiple) y posteriormente hacemos el promedio del error al cuadrado (MSE) entre los datos reales y la predicción del modelo. Por otra parte, para conocer la dirección en que deben ajustarse los parámetros tras cada época usamos el algoritmo de gradiente descendiente.

El entrenamiento y las pruebas se hacen por medio de validación cruzada, después de revolver las instancias, separamos subconjuntos de tamaño “ x ” y entrenamos con el resto, luego evaluamos con este subconjunto de “ x ” instancias. Repetimos hasta haber evaluado cada fila. Hacemos esto para minimizar el sesgo que podría ocasionar entrenar y probar con el mismo conjunto de datos. De esta forma evitamos caer en que el modelo solo haya “memorizado” los datos.

El funcionamiento del código está más especificado en los comentarios de este.

Resultados y conclusiones

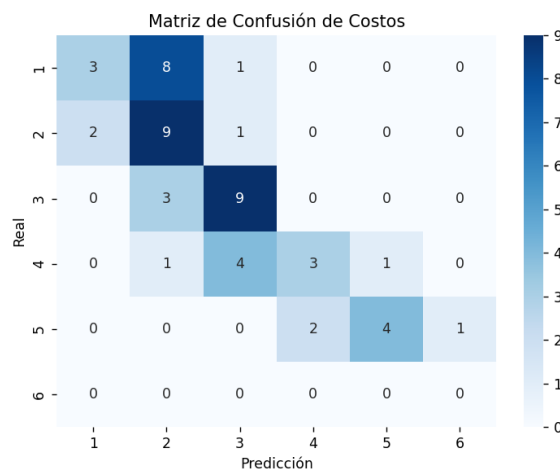
Las predicciones del modelo se ven de esta forma.

```

Predictions: [3, 4, 5, 4, 2]
Actual: [3, 5, 5, 4, 1]
Block 0: Error = 0.4000
Predictions: [3, 4, 1, 3, 2]
Actual: [1, 5, 1, 3, 1]
Block 1: Error = 1.2000
Predictions: [2, 3, 4, 3, 6]
Actual: [2, 3, 3, 3, 5]
Block 2: Error = 0.4000
Predictions: [3, 1, 3, 1, 2]
Actual: [3, 2, 2, 2, 1]
Block 3: Error = 0.8000
Predictions: [4, 2, 2, 3, 2]
Actual: [4, 1, 1, 3, 4]
Block 4: Error = 1.2000
Predictions: [1, 3, 2, 3, 3]
Actual: [1, 3, 1, 4, 4]
Block 5: Error = 0.6000
Predictions: [2, 2, 3, 4, 2]
Actual: [2, 2, 3, 5, 1]
Block 6: Error = 0.4000
Predictions: [3, 3, 3, 4, 3]
Actual: [2, 3, 4, 4, 4]
Block 7: Error = 0.6000
Predictions: [2, 2, 4, 2, 2]
Actual: [1, 3, 4, 1, 2]
Block 8: Error = 0.6000
Predictions: [5, 2, 2, 5, 2]
Actual: [5, 3, 2, 5, 2]
Block 9: Error = 0.2000
Predictions: [2, 2]
Actual: [2, 2]
Block 10: Error = 0.0000

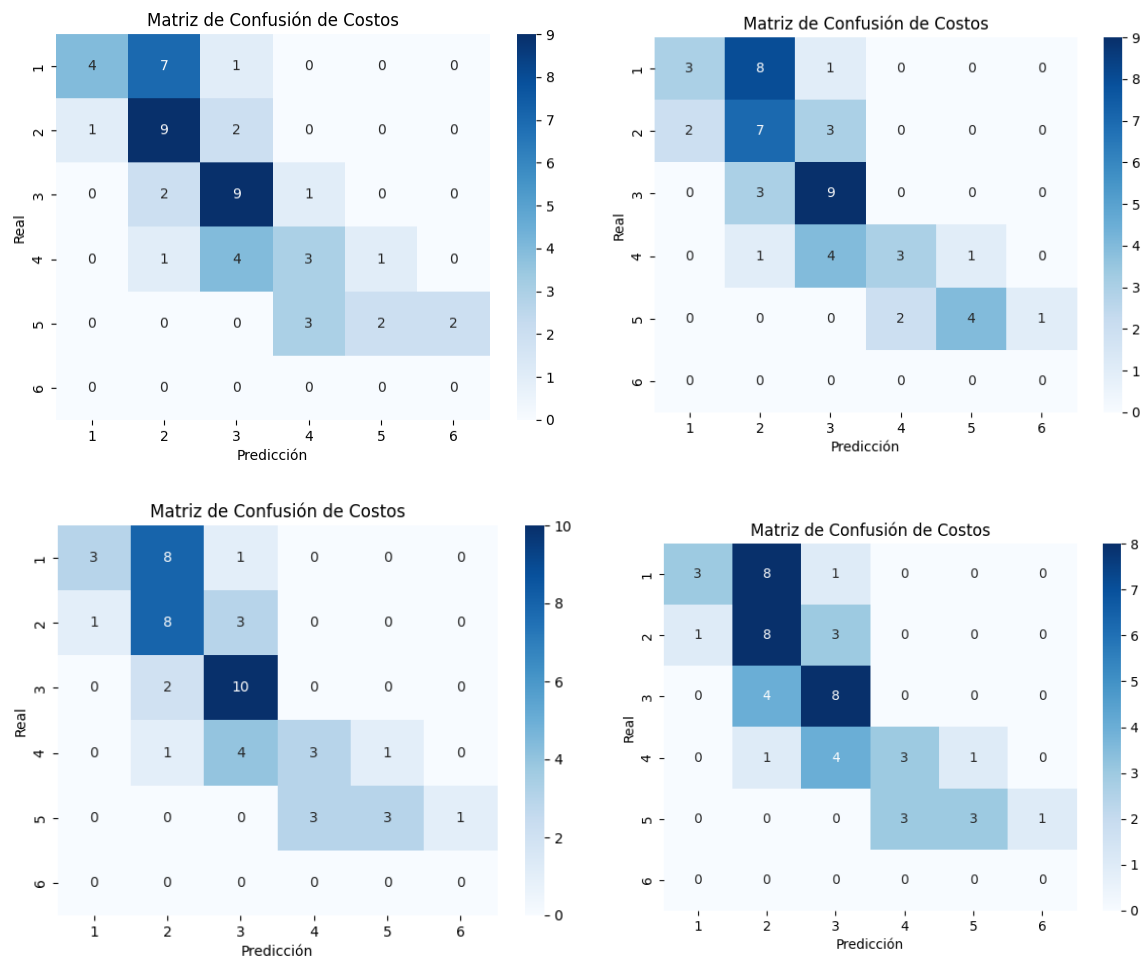
```

El promedio de error para esta simulación fue 0.5818, lo que nos indica que generalmente una predicción está errada por la raíz de esa cantidad (0.761). En consecuencia, podemos afirmar que el modelo predice correctamente algunos costos y los erróneos varían en su mayoría por una categoría hacia arriba o abajo. Esto lo podemos ver más fácilmente en la matriz de confusión.

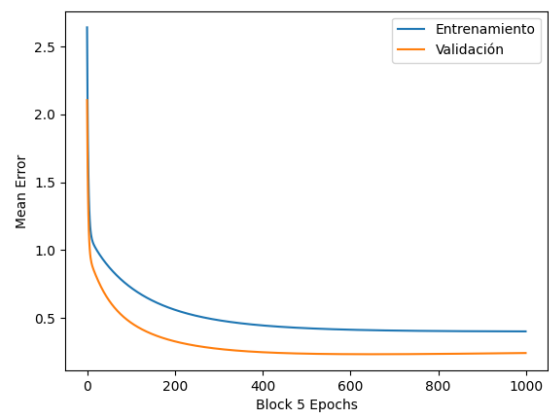
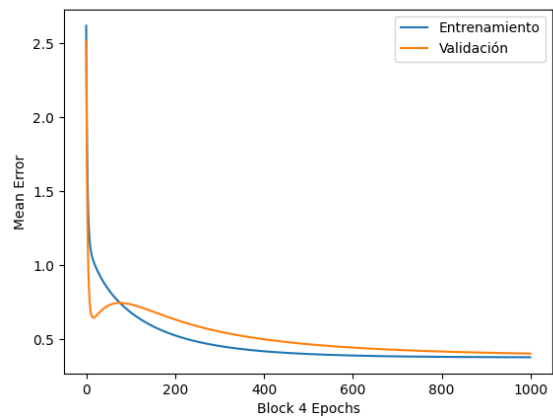
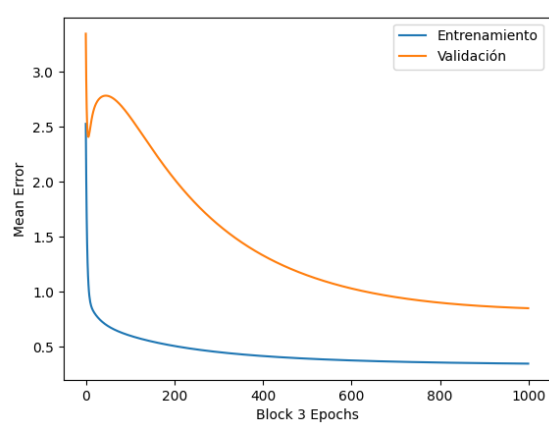
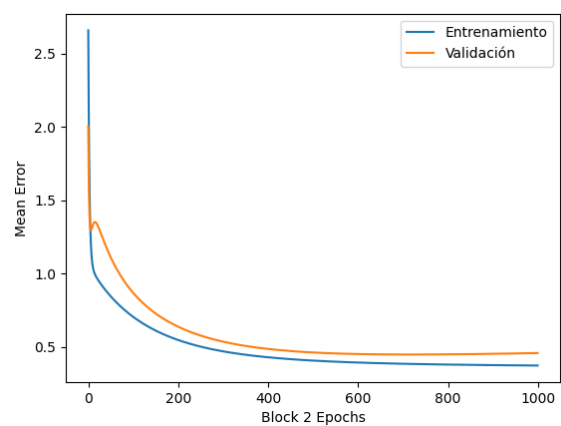
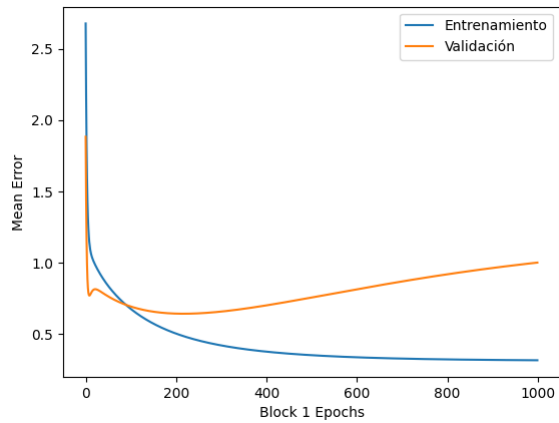
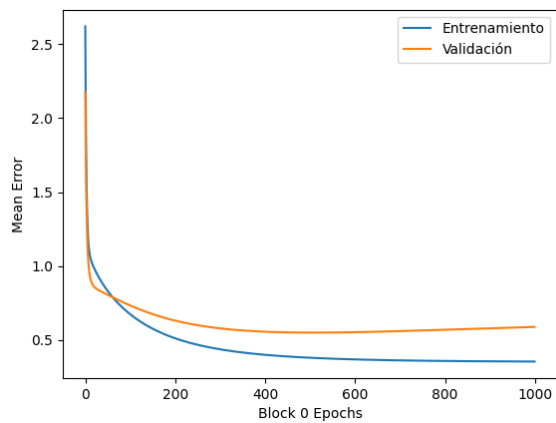


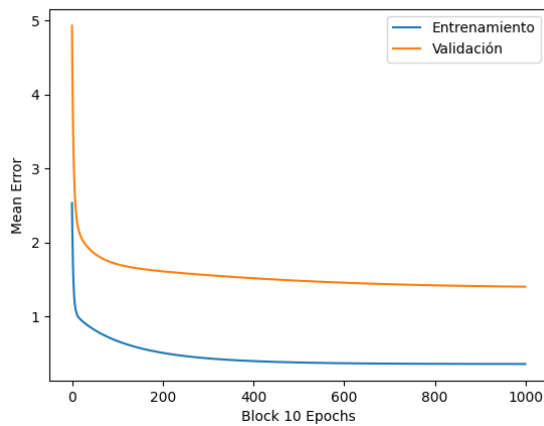
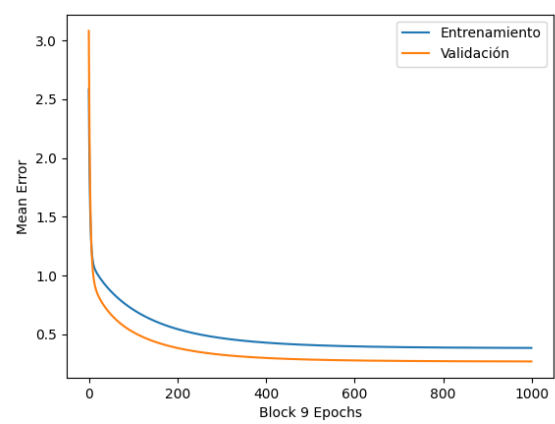
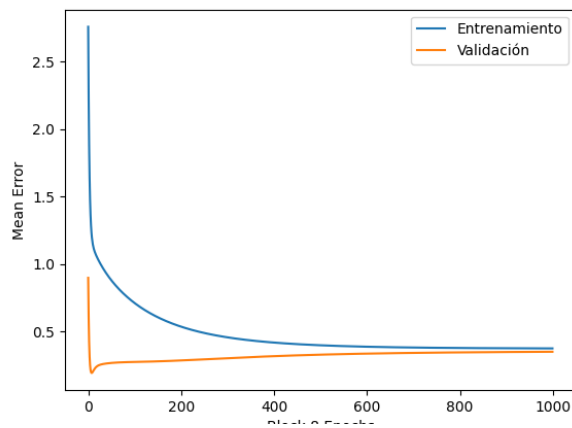
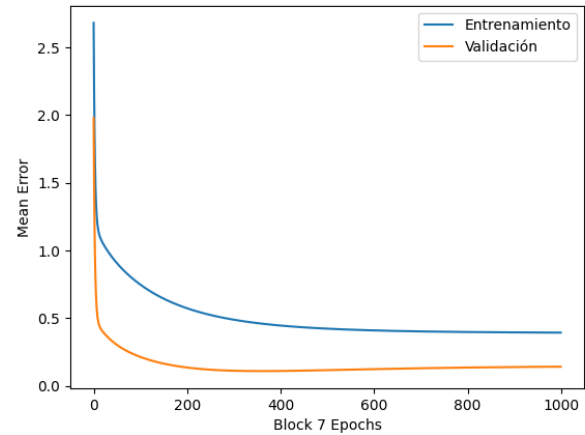
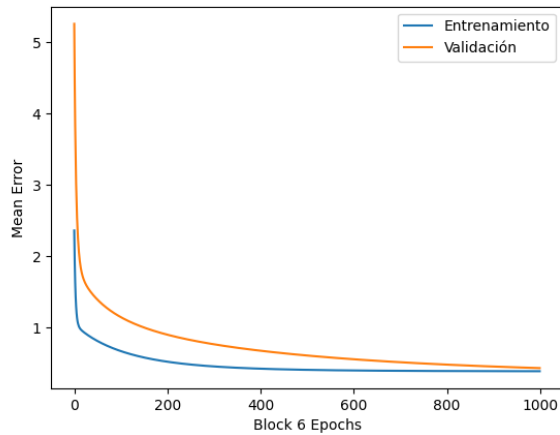
Tras múltiples ejecuciones del código encontramos que el patrón es el mismo. Suele fallar en un 75% de veces prediciendo los campeones de costo 1 asignándoles 2 y suele acertar en un 75% de las veces en el caso de predecir los de coste 3. Para los costes 4 y 5 erra en

poco más del 50% de los casos y en el caso de coste 2 acierta en el 66% de veces. Generalmente acierta 27 de 52 instancias (Precisión: 51%).



La tasa de aprendizaje utilizada fue de 0.1 y cada bloque de entrenamiento fue sometido a 1000 épocas. Estos hiper parámetros fueron seleccionados porque llegaban al límite de aprendizaje rápidamente. Al experimentar con números de épocas muy grandes el error no bajaba de 0.38 por lo que no tenía caso aumentarlas. Los gráficos de error vs época del entrenamiento y validación se ven de esta forma

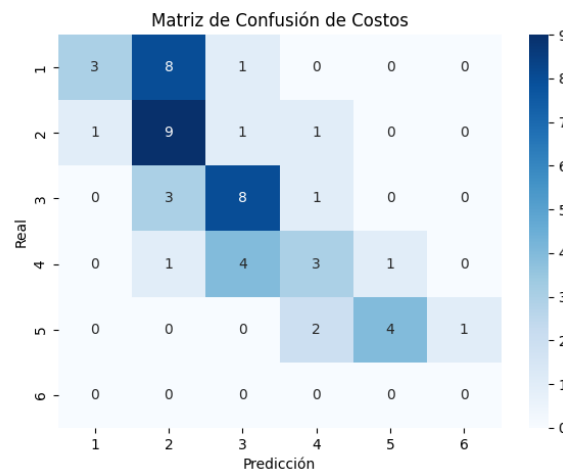




- Podemos concluir que el modelo sí aprende pues la curva de entrenamiento baja y así se mantiene.
- El comportamiento del modelo varía mucho dependiendo del bloque de datos. En algunos el error de la validación es incluso menor al de entrenamiento lo que podría indicar que el modelo generaliza bien o que ese bloque tenía ejemplos muy sencillos, mientras en otros cómo en el bloque 1 podríamos interpretar un problema de

overfitting pues el error de validación aumenta ante un caso que desconoce. Esto puede suceder debido a que ciertos bloques contuvieron campeones que no tenían un claro patrón entre sus estadísticas y sus costos, pues existen más atributos de los que depende el costo del campeón como su escalado en el tiempo y el potencial de su habilidad.

- Los casos dónde el error de validación baja y luego sube son pocos y ligeros, lo que indica que no hay un grave problema de overfitting.
- De la matriz de confusión podemos intuir que el modelo es muy simple para los datos. Mientras que campeones de costos bajos predice alto, para costos altos predice abajo y en medio es dónde tiene mayor afinidad un comportamiento no lineal, por lo que concluimos que este modelo no tiene la capacidad de ajustarse mucho más a los datos.
- Se intentó mejorar el modelo añadiendo features no lineales, por ejemplo: health, speed_of_attack, defense, attack al cubo o usando tangente hiperbólica y se ajustó un poco mejor, pero también aumentó la dispersión. La precisión fue la misma.



- Por otra parte, la columna objetivo son valores discretos, por lo que se adaptaría mejor un modelo lógico.

Cambio de dataset

Con el objetivo de poder comparar el comportamiento del modelo en entrenamiento, validación y prueba, era necesario incrementar el número de instancias del dataset. Anteriormente contábamos con 52 que eran pocas para poder tener un conjunto de entrenamiento y de prueba. Para solucionar esto unimos los datos de las versiones 14 y 15 del juego obteniendo un dataset de 125 instancias.

Extracción

Los datos fueron extraídos de las páginas web https://wiki.leagueoflegends.com/en-us/TFT:List_of_champions/Base_statistics/Set_14 y <https://www.datafft.com/database#unit>. Los datos originales se encuentran en los archivos TFT_set_14_raw.txt y TFT_set_15.csv.

Transformación

Para los datos de la versión 14 utilizamos un script de Python (Data_transformation_set_14.py) que nos ayudó a pasar los datos copiados de la página web a un csv con el formato que tenían los de la versión 15. Los datos de la versión 15 fueron extraídos manualmente por lo que en este proceso se hizo la discriminación de features y el acomodo de las columnas.

Posteriormente, juntamos los datos de ambas versiones en el archivo TFT_set_14_y_15.csv. Para poderlo utilizar en nuestro modelo era necesario someterlo a un escalado previamente. Para esto utilizamos el script MinMaxScaler.py. Nos quedó de esta forma:

```

total_health,health2,initial_mmu_size,initial_atk2,atk3,atk6,atk7,atk8,atk9,atk10,atk11,atk12,atk13,atk14,atk15,atk16,atk17,atk18,atk19,atk20,atk21,atk22,atk23,atk24,atk25,atk26,atk27,atk28,atk29,atk30,atk31,atk32,atk33,atk34,atk35,atk36,atk37,atk38,atk39,atk40,atk41,atk42,atk43,atk44,atk45,atk46,atk47,atk48,atk49,atk50,atk51,atk52,atk53,atk54,atk55,atk56,atk57,atk58,atk59,atk60,atk61,atk62,atk63,atk64,atk65,atk66,atk67,atk68,atk69,atk70,atk71,atk72,atk73,atk74,atk75,atk76,atk77,atk78,atk79,atk80,atk81,atk82,atk83,atk84,atk85,atk86,atk87,atk88,atk89,atk90,atk91,atk92,atk93,atk94,atk95,atk96,atk97,atk98,atk99,atk100,atk101,atk102,atk103,atk104,atk105,atk106,atk107,atk108,atk109,atk110,atk111,atk112,atk113,atk114,atk115,atk116,atk117,atk118,atk119,atk120,atk121,atk122,atk123,atk124,atk125,atk126,atk127,atk128,atk129,atk130,atk131,atk132,atk133,atk134,atk135,atk136,atk137,atk138,atk139,atk140,atk141,atk142,atk143,atk144,atk145,atk146,atk147,atk148,atk149,atk150,atk151,atk152,atk153,atk154,atk155,atk156,atk157,atk158,atk159,atk160,atk161,atk162,atk163,atk164,atk165,atk166,atk167,atk168,atk169,atk170,atk171,atk172,atk173,atk174,atk175,atk176,atk177,atk178,atk179,atk180,atk181,atk182,atk183,atk184,atk185,atk186,atk187,atk188,atk189,atk190,atk191,atk192,atk193,atk194,atk195,atk196,atk197,atk198,atk199,atk200,atk201,atk202,atk203,atk204,atk205,atk206,atk207,atk208,atk209,atk210,atk211,atk212,atk213,atk214,atk215,atk216,atk217,atk218,atk219,atk220,atk221,atk222,atk223,atk224,atk225,atk226,atk227,atk228,atk229,atk230,atk231,atk232,atk233,atk234,atk235,atk236,atk237,atk238,atk239,atk240,atk241,atk242,atk243,atk244,atk245,atk246,atk247,atk248,atk249,atk250,atk251,atk252,atk253,atk254,atk255,atk256,atk257,atk258,atk259,atk260,atk261,atk262,atk263,atk264,atk265,atk266,atk267,atk268,atk269,atk270,atk271,atk272,atk273,atk274,atk275,atk276,atk277,atk278,atk279,atk280,atk281,atk282,atk283,atk284,atk285,atk286,atk287,atk288,atk289,atk290,atk291,atk292,atk293,atk294,atk295,atk296,atk297,atk298,atk299,atk300,atk301,atk302,atk303,atk304,atk305,atk306,atk307,atk308,atk309,atk310,atk311,atk312,atk313,atk314,atk315,atk316,atk317,atk318,atk319,atk320,atk321,atk322,atk323,atk324,atk325,atk326,atk327,atk328,atk329,atk330,atk331,atk332,atk333,atk334,atk335,atk336,atk337,atk338,atk339,atk340,atk341,atk342,atk343,atk344,atk345,atk346,atk347,atk348,atk349,atk350,atk351,atk352,atk353,atk354,atk355,atk356,atk357,atk358,atk359,atk360,atk361,atk362,atk363,atk364,atk365,atk366,atk367,atk368,atk369,atk370,atk371,atk372,atk373,atk374,atk375,atk376,atk377,atk378,atk379,atk380,atk381,atk382,atk383,atk384,atk385,atk386,atk387,atk388,atk389,atk390,atk391,atk392,atk393,atk394,atk395,atk396,atk397,atk398,atk399,atk400,atk401,atk402,atk403,atk404,atk405,atk406,atk407,atk408,atk409,atk410,atk411,atk412,atk413,atk414,atk415,atk416,atk417,atk418,atk419,atk420,atk421,atk422,atk423,atk424,atk425,atk426,atk427,atk428,atk429,atk430,atk431,atk432,atk433,atk434,atk435,atk436,atk437,atk438,atk439,atk440,atk441,atk442,atk443,atk444,atk445,atk446,atk447,atk448,atk449,atk450,atk451,atk452,atk453,atk454,atk455,atk456,atk457,atk458,atk459,atk460,atk461,atk462,atk463,atk464,atk465,atk466,atk467,atk468,atk469,atk470,atk471,atk472,atk473,atk474,atk475,atk476,atk477,atk478,atk479,atk480,atk481,atk482,atk483,atk484,atk485,atk486,atk487,atk488,atk489,atk490,atk491,atk492,atk493,atk494,atk495,atk496,atk497,atk498,atk499,atk500,atk501,atk502,atk503,atk504,atk505,atk506,atk507,atk508,atk509,atk510,atk511,atk512,atk513,atk514,atk515,atk516,atk517,atk518,atk519,atk520,atk521,atk522,atk523,atk524,atk525,atk526,atk527,atk528,atk529,atk530,atk531,atk532,atk533,atk534,atk535,atk536,atk537,atk538,atk539,atk540,atk541,atk542,atk543,atk544,atk545,atk546,atk547,atk548,atk549,atk550,atk551,atk552,atk553,atk554,atk555,atk556,atk557,atk558,atk559,atk560,atk561,atk562,atk563,atk564,atk565,atk566,atk567,atk568,atk569,atk570,atk571,atk572,atk573,atk574,atk575,atk576,atk577,atk578,atk579,atk580,atk581,atk582,atk583,atk584,atk585,atk586,atk587,atk588,atk589,atk590,atk591,atk592,atk593,atk594,atk595,atk596,atk597,atk598,atk599,atk600,atk601,atk602,atk603,atk604,atk605,atk606,atk607,atk608,atk609,atk610,atk611,atk612,atk613,atk614,atk615,atk616,atk617,atk618,atk619,atk620,atk621,atk622,atk623,atk624,atk625,atk626,atk627,atk628,atk629,atk630,atk631,atk632,atk633,atk634,atk635,atk636,atk637,atk638,atk639,atk640,atk641,atk642,atk643,atk644,atk645,atk646,atk647,atk648,atk649,atk650,atk651,atk652,atk653,atk654,atk655,atk656,atk657,atk658,atk659,atk660,atk661,atk662,atk663,atk664,atk665,atk666,atk667,atk668,atk669,atk670,atk671,atk672,atk673,atk674,atk675,atk676,atk677,atk678,atk679,atk680,atk681,atk682,atk683,atk684,atk685,atk686,atk687,atk688,atk689,atk690,atk691,atk692,atk693,atk694,atk695,atk696,atk697,atk698,atk699,atk700,atk701,atk702,atk703,atk704,atk705,atk706,atk707,atk708,atk709,atk710,atk711,atk712,atk713,atk714,atk715,atk716,atk717,atk718,atk719,atk720,atk721,atk722,atk723,atk724,atk725,atk726,atk727,atk728,atk729,atk730,atk731,atk732,atk733,atk734,atk735,atk736,atk737,atk738,atk739,atk740,atk741,atk742,atk743,atk744,atk745,atk746,atk747,atk748,atk749,atk750,atk751,atk752,atk753,atk754,atk755,atk756,atk757,atk758,atk759,atk760,atk761,atk762,atk763,atk764,atk765,atk766,atk767,atk768,atk769,atk770,atk771,atk772,atk773,atk774,atk775,atk776,atk777,atk778,atk779,atk780,atk781,atk782,atk783,atk784,atk785,atk786,atk787,atk788,atk789,atk790,atk791,atk792,atk793,atk794,atk795,atk796,atk797,atk798,atk799,atk800,atk801,atk802,atk803,atk804,atk805,atk806,atk807,atk808,atk809,atk810,atk811,atk812,atk813,atk814,atk815,atk816,atk817,atk818,atk819,atk820,atk821,atk822,atk823,atk824,atk825,atk826,atk827,atk828,atk829,atk830,atk831,atk832,atk833,atk834,atk835,atk836,atk837,atk838,atk839,atk8
```

Cabe destacar que hemos agregado nuevas columnas con el objetivo de diferenciar mejor los costos de los campeones. Ahora añadimos el escalado de vida y ataque.

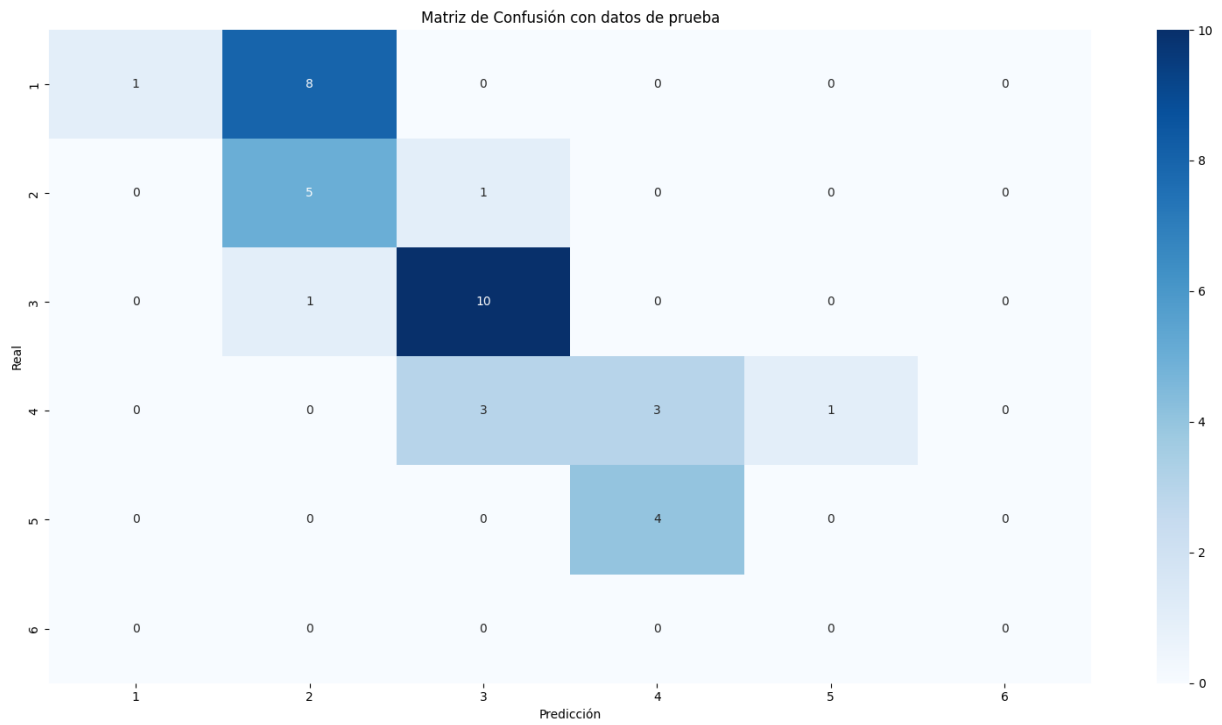
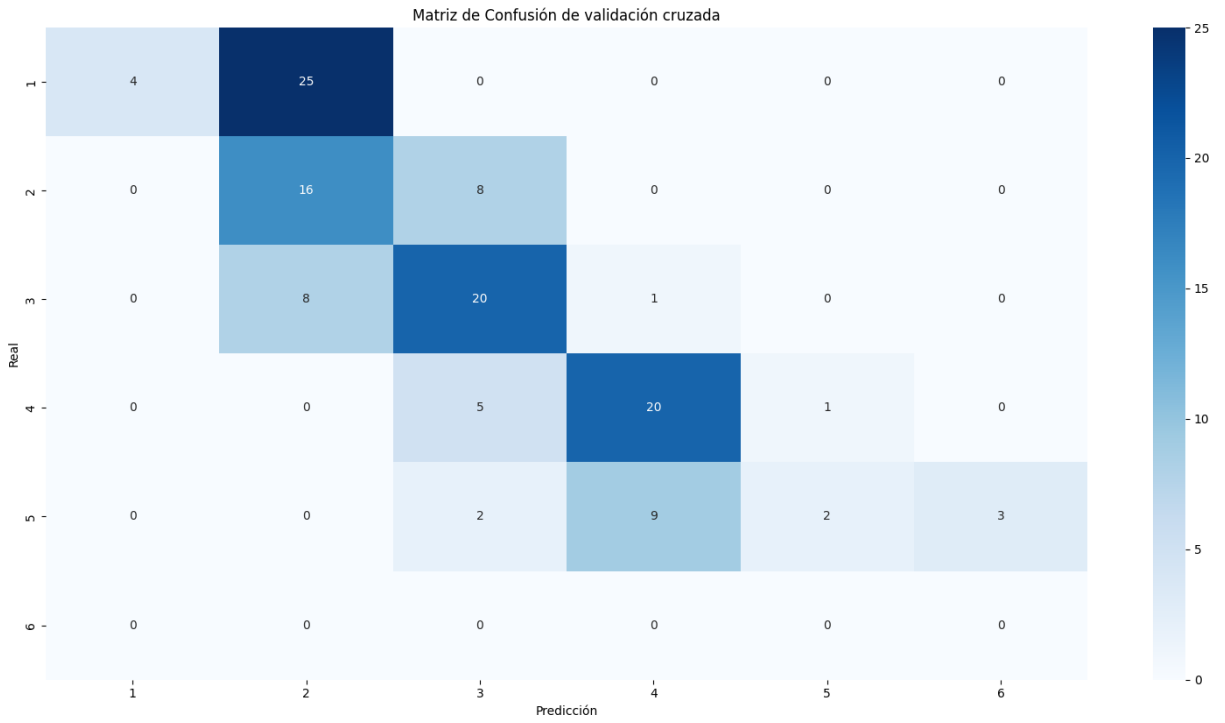
Finalmente resolvimos las instancias y tomamos 37 para pruebas y 88 para entrenamiento. Cada conjunto lo separamos en los archivos TFT_set_14_y_15_test.csv y TFT set 14 y 15 train.csv respectivamente.

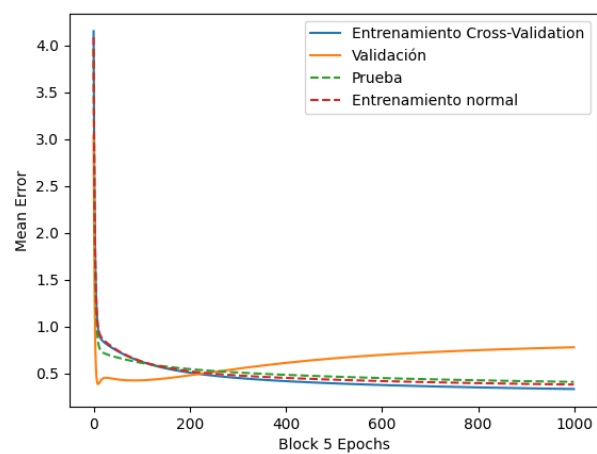
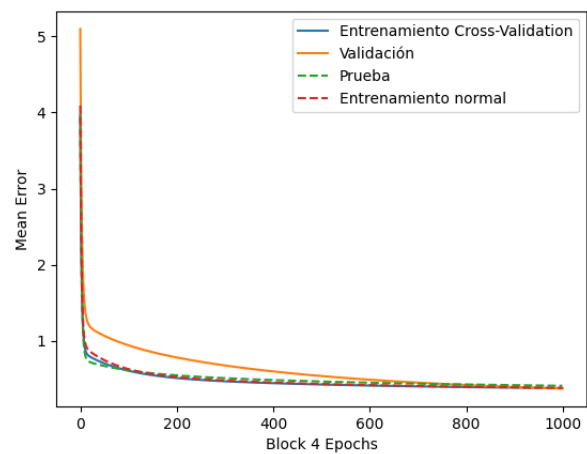
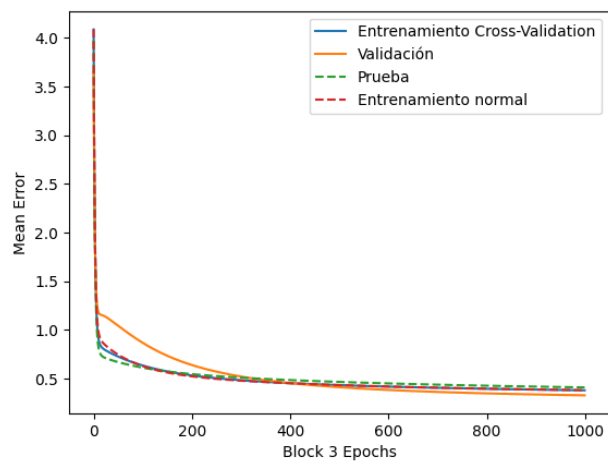
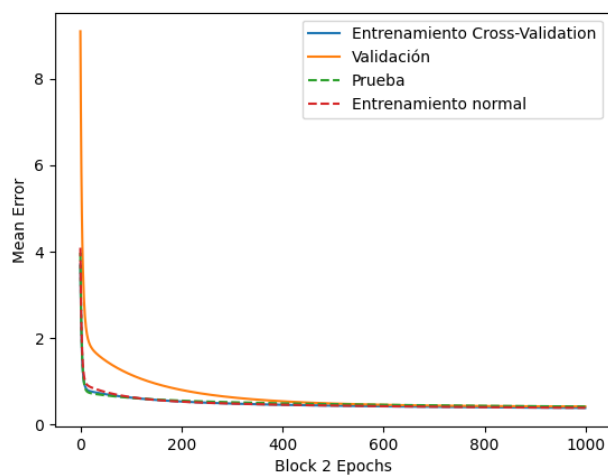
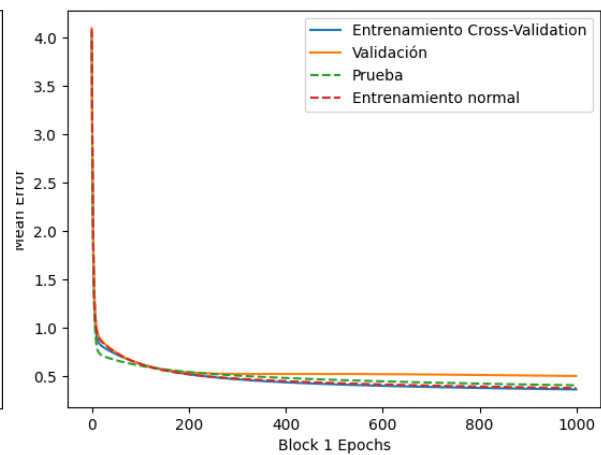
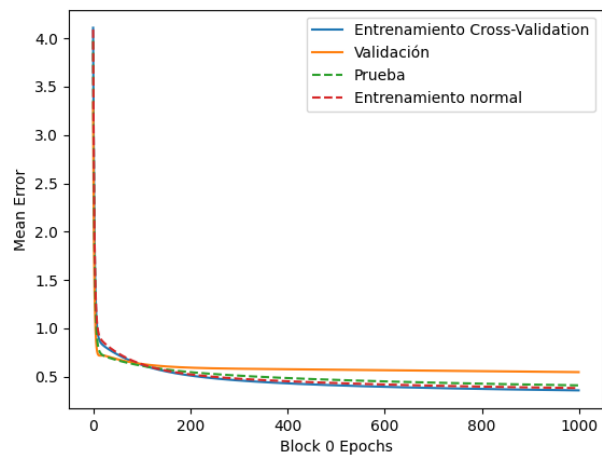
Carga

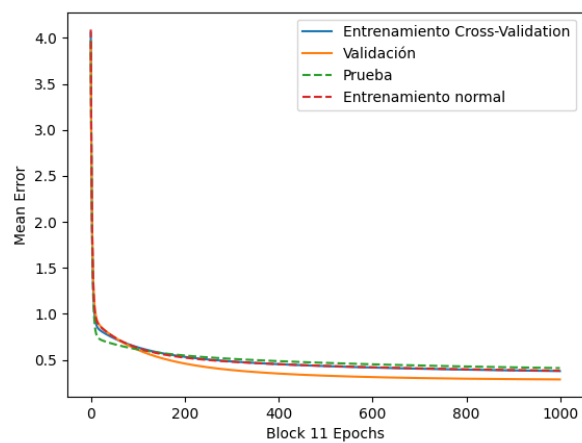
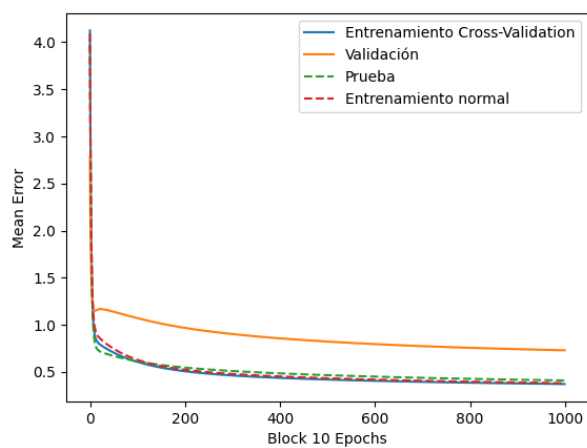
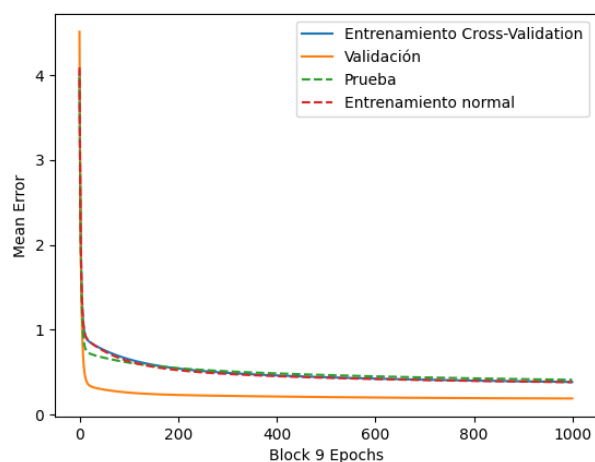
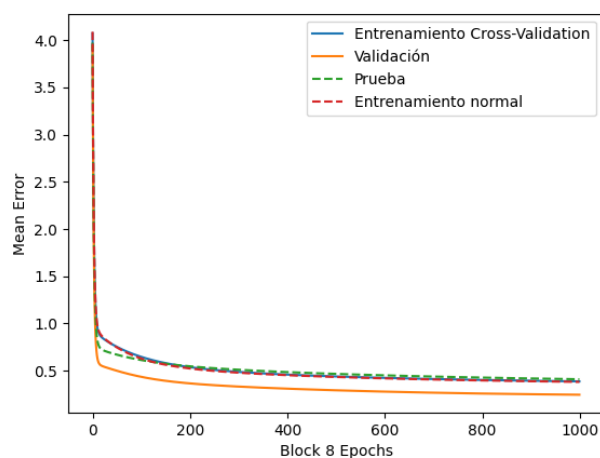
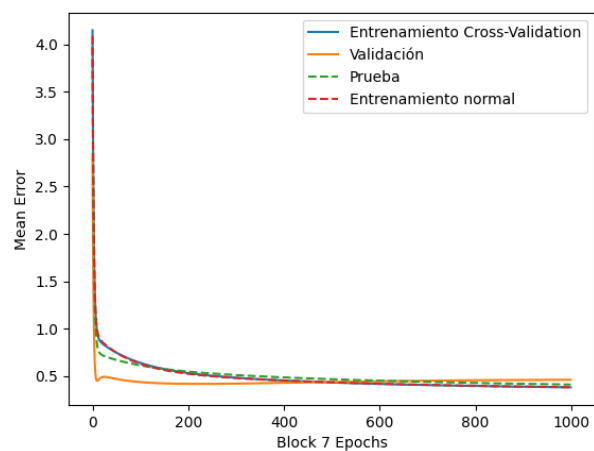
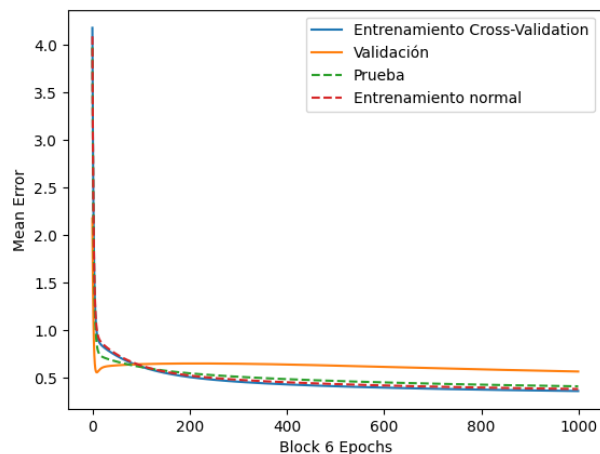
Cada conjunto de datos que se ha ido utilizando están en el repositorio <https://github.com/ViejoAgrio/Machine-Learning> en la carpeta de Datasets.

Resultados

Con estos cambios podemos comparar el comportamiento del modelo en el entrenamiento, validación y pruebas. Los resultados fueron los siguientes.







	Test	Validación
Error promedio	0.502	0.53
Precisión	51%	49%

Conclusiones

- Las curvas de error entre el entrenamiento con cross validation y el entrenamiento normal son muy similares, lo que indica que el modelo no depende de cómo se parta el dataset y consecuentemente a que no se sobre ajusta a los datos.
- A pesar del cambio de dataset a uno con más instancias y features, esto no influyó en un cambio significativo en el error promedio o precisión de las predicciones, en consecuencia, podemos intuir que es el modelo el que no se adecua para el aprendizaje de los patrones en estos datos.
- De igual forma, el patrón en las matrices de confusión se mantuvo.

Modelo con framework

Con el objetivo de encontrar un modelo que nos sea útil para diseñar y/o predecir las estadísticas base de los campeones, intentaremos utilizar otro modelo más adecuado para este tipo de datos.

Dado que nuestra variable objetivo son múltiples valores discretos y anteriormente probamos con una regresión logística múltiple, ahora intentaremos con un algoritmo basado en árboles de decisión: Random forest de clasificación con cross entropy para calcular el error.

El código de esta implementación se encuentra en el archivo `With_Framework.py` en la carpeta de scripts.

Marco teórico

El Random Forest es un algoritmo de aprendizaje supervisado basado en la técnica de *ensemble learning*, cuyo objetivo es mejorar la capacidad predictiva y la estabilidad de los modelos individuales. Combina múltiples árboles de decisión para generar un modelo más robusto. La idea central es que, en lugar de depender de un único árbol —que puede ser muy sensible a cambios en los datos, se construye un conjunto (o "bosque") de árboles y se toma una decisión conjunta, reduciendo la varianza y mejorando la generalización.

El proceso de construcción de un Random Forest comienza con la técnica de bagging (Bootstrap Aggregating). A partir del conjunto de datos original, se generan múltiples subconjuntos de entrenamiento mediante muestreo aleatorio con reemplazo. A cada subconjunto se le entrena un árbol de decisión independiente. Este procedimiento asegura que cada árbol vea una versión ligeramente distinta de los datos, lo cual introduce diversidad en el conjunto de modelos y evita que todos aprendan los mismos patrones.

Además del bagging, el Random Forest introduce otra fuente de aleatoriedad: en cada nodo del árbol, en lugar de evaluar todas las variables disponibles para decidir la mejor división, se selecciona al azar un subconjunto de features. Esto significa que diferentes árboles pueden tomar decisiones basadas en diferentes atributos, lo cual incrementa la diversidad entre los modelos y reduce la correlación entre ellos. Este paso es crucial para que el ensemble se beneficie de la “sabiduría de la multitud”.

En la fase de predicción, el funcionamiento depende del tipo de problema. En clasificación, cada árbol vota por una clase y el Random Forest escoge la clase mayoritaria. En regresión, se calcula el promedio de las predicciones individuales.

El Random Forest presenta varias ventajas: maneja bien datos de alta dimensión, es resistente al ruido, puede capturar relaciones no lineales y proporciona medidas de importancia de variables, lo cual ayuda a interpretar los resultados. Sin embargo, también tiene limitaciones, como el hecho de que puede volverse computacionalmente costoso con grandes cantidades de árboles o datos, y que su interpretabilidad es menor que la de un único árbol de decisión.

Dataset

El dataset utilizado fue el mismo que se expone en la sección de este documento: “Cambio de dataset”. La única diferencia es que para esta implementación se usó la versión sin escalamiento de los datos pues el framework lo hace en automático.

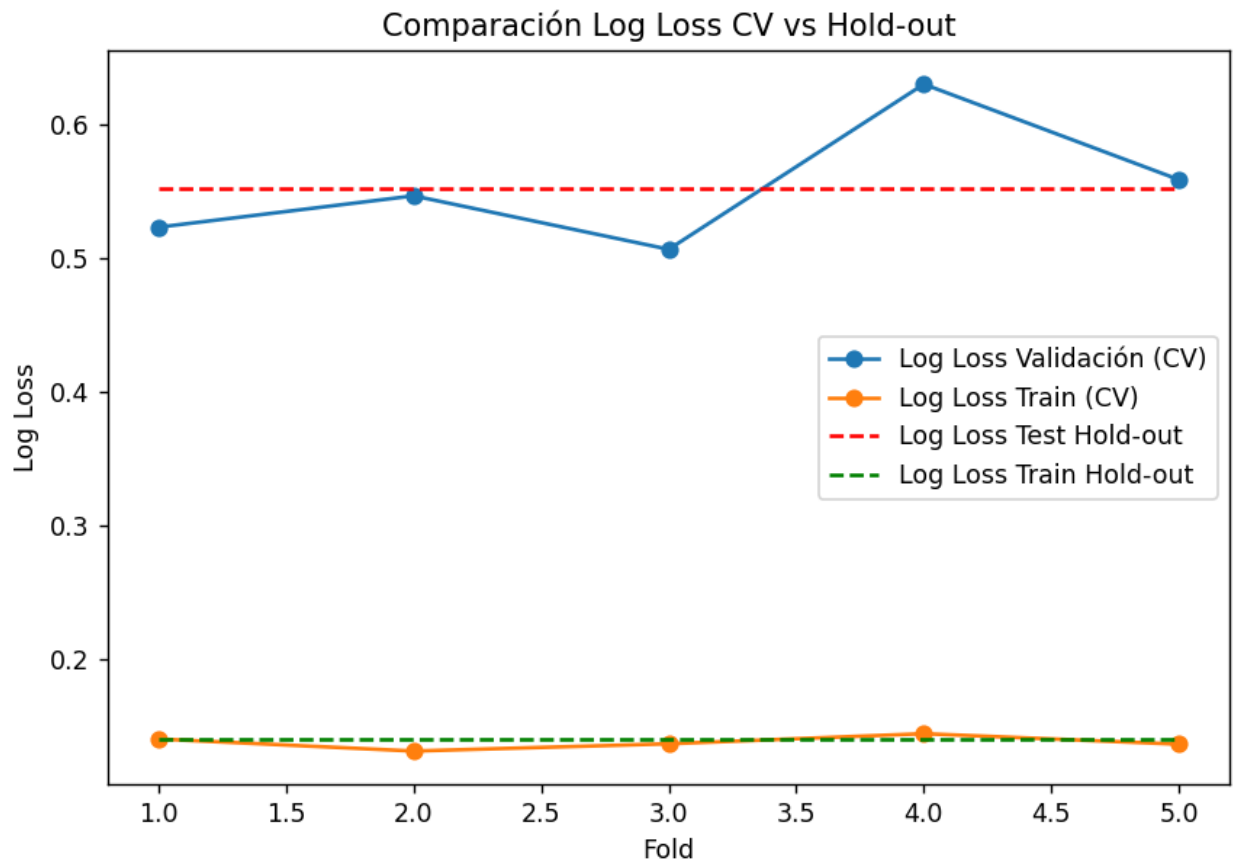
Resultados

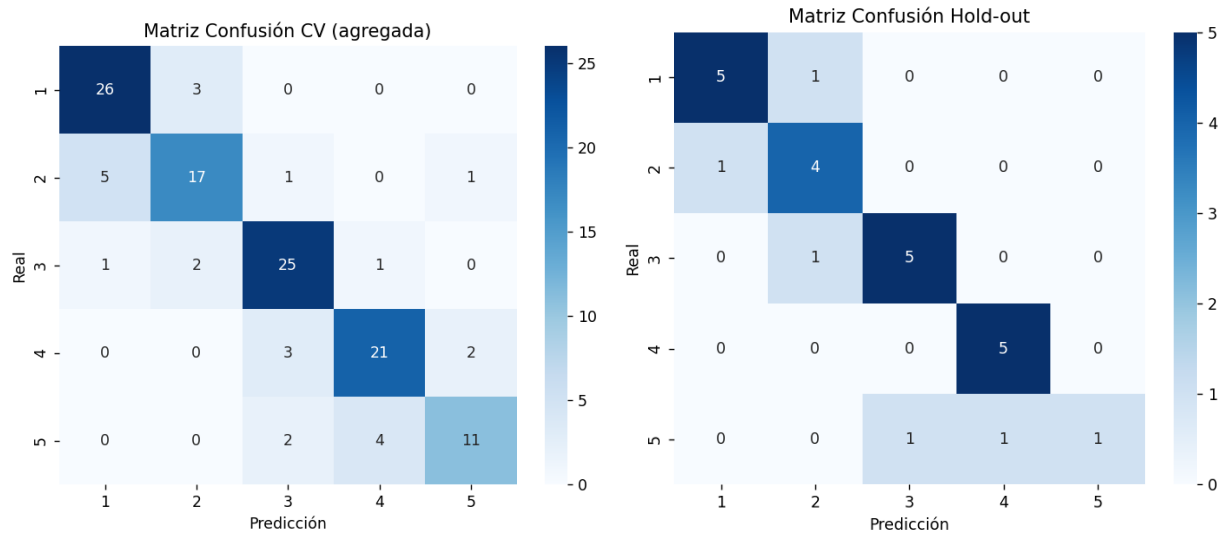
Tras el primer intento podemos ver rápidamente un gran incremento en la precisión del modelo, ahora obtenemos alrededor de un 80% tanto en validación como en test cuando con el modelo anterior era alrededor de 55%. Sin embargo, también es muy evidente un problema de overfitting dado que generalmente en el training con o sin validación cruzada obtenemos un 100% de precisión. Esto sugiere que se están memorizando los datos pues cuando lo exponemos a casos que no conoce falla significativamente más y la confianza con la que hace las aseveraciones de una clase en promedio descienden de 86% a 58% aproximadamente ($\text{Logloss} = 0.14$, $\exp(-0.14) = 0.86$. $\text{Logloss} = 0.54$, $\exp(-0.54) = 0.58$), es decir que para predecir el costo de un campeón generalmente esa clase tiene un 58% y entre las otras se reparten el 42% restante, por lo que normalmente hay un claro favorito para la predicción.


```
Fold 1: Train acc=1.0000 LogLoss=0.1406 | Val acc=0.6800 LogLoss=0.5235
Fold 2: Train acc=1.0000 LogLoss=0.1316 | Val acc=0.7600 LogLoss=0.5469
Fold 3: Train acc=1.0000 LogLoss=0.1370 | Val acc=0.8400 LogLoss=0.5069
Fold 4: Train acc=1.0000 LogLoss=0.1446 | Val acc=0.8800 LogLoss=0.6304
Fold 5: Train acc=1.0000 LogLoss=0.1367 | Val acc=0.8400 LogLoss=0.5590
```

Para solucionar el problema de overfitting se intentó limitar la profundidad de los árboles y aumentar su número (1000 árboles de 4 nodos de profundidad), sin embargo, aunque si acortó la diferencia entre la precisión en train y validación ligeramente (94.8% para train y 82% para validación), el logloss aumentó de forma significativa (la seguridad con la que predecía descendió de 58% a 42%), por lo que no se solucionó el problema y no tenemos un incremento importante de precisión.

```
Fold 1: Train acc=0.9500 LogLoss=0.4782 | Val acc=0.7200 LogLoss=0.7155
Fold 2: Train acc=0.9700 LogLoss=0.4296 | Val acc=0.7600 LogLoss=0.6944
Fold 3: Train acc=0.9400 LogLoss=0.4589 | Val acc=0.8000 LogLoss=0.6696
Fold 4: Train acc=0.9600 LogLoss=0.4548 | Val acc=0.8800 LogLoss=0.7675
Fold 5: Train acc=0.9300 LogLoss=0.4575 | Val acc=0.9200 LogLoss=0.7699
```

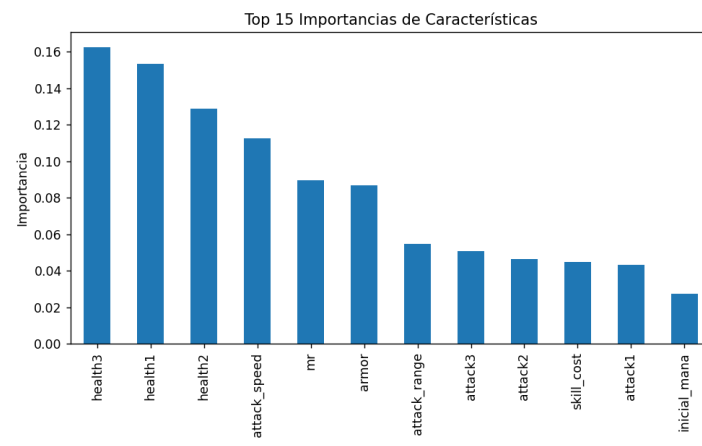




A diferencia del modelo de regresión, podemos apreciar que el patrón que se presentaba en el que fallaba con los valores extremos asignándoles un costo errado hacia el centro aquí no ocurre.

Algo que llama la atención es el caso del campeón de costo 2 que marcó como costo 5 en la matriz de confusión de validación. Este se trata de un personaje con una característica muy particular que está vinculada con su habilidad dentro del juego y en una de sus estadísticas lo hace parecer sumamente poderoso. Es un caso que vale la pena destacar pues es un gran ejemplo que ilustra el tratamiento de estas instancias ruidosas en un modelo de este tipo.

Finalmente, en la imagen de abajo podemos ver la importancia de cada uno de las features del dataset. Al parecer fue sumamente relevante añadir el escalado de vida que tenían los personajes pues cómo vemos las características más importantes para decidir el costo de un campeón fueron estas.



Conclusiones

- La hipótesis de que el modelo de random forest se ajustaría mejor que una regresión resultó ser verdadera.
- A pesar de que este modelo tiene un gran problema de overfitting, su precisión a la hora de las pruebas sigue siendo bastante superior, por lo que para resolver el problema planteado utilizaríamos este modelo.
- El uso del framework sklearn agilizó en gran medida la codificación del modelo, sin embargo, quedan muchos espacios desconocidos sobre cómo implementa las partes. ¿Cómo decide la profundidad de los árboles cuando no la indicamos? ¿Qué estrategia sigue para el armado de los árboles? ¿Las features que usa cada árbol son totalmente aleatorias? ¿Cuántas toma cada árbol? Aunque es probable que tenga opciones para ajustar cada una de estas partes, la comprensión y exploración de las diferentes técnicas son un tema que da para otro documento de investigación completo.