

---

# Implementación de una feed forward network para predecir si una persona es propensa a padecer de problemas del corazón

---

José Alejandro López Quel

21001127

Universidad Galileo

Statistical Learning II

## Abstract

1        En este paper, se presenta un caso de predicción de la posibilidad de llegar a tener  
2        problemas del corazón. Al determinar si una persona es propensa a padecer de  
3        problemas del corazón hay factores importantes como la presión arterial, uso de  
4        sustancias nocivas como alcohol y tabaco, obesidad, entre otras, los cuales pueden  
5        ayudar a predecir la existencia de futuras enfermedades. Por lo que en este trabajo,  
6        se analiza utilizar los datos antes mencionados para predecir estos comportamientos  
7        empleando una red neuronal del tipo feed forward.

## 8    1    Introducción

9        El aprendizaje profundo (deep learning) se ha utilizado para resolver problemas del mundo real en  
10       muchos ámbitos. La medicina no es una excepción. Aunque son controvertidos, se han propuesto y  
11       utilizado múltiples modelos con cierto éxito. Cada año mueren en Estados Unidos unas 610.000  
12       personas por enfermedades del corazón, lo que supone 1 de cada 4 muertes. Las enfermedades del  
13       corazón son la principal causa de muerte tanto en hombres como en mujeres.

14       Se considera que los factores de riesgo tradicionales de los problemas de corazón son el  
15       colesterol LDL alto, la hipertensión arterial, la diabetes mellitus, el tabaquismo, los antecedentes  
16       familiares de problemas cardíacos, la edad, la obesidad y un estilo de vida poco saludable. Sin  
17       embargo, estos problemas pueden controlarse eficazmente con un cambio de estilo de vida y la  
18       adopción de hábitos saludables, y por lo tanto esto equivale a ahorrar el alto costo de tratamientos  
19       médicos y hospitalizaciones si se detecta a tiempo. Con la detección precoz de estos problemas,  
20       los pacientes pueden recibir una serie de tratamientos aconsejados por los médicos para reducir  
21       el riesgo de futuros problemas cardíacos y aliviar o controlar los síntomas. En este contexto los  
22       historiales clínicos pueden considerarse un recurso de información útil para ayudar a los médicos en  
23       la detección o la predicción de problemas del corazón.

24       Los avances en el aprendizaje automático y la inteligencia artificial han motivado a muchos  
25       científicos de datos a utilizar estas tecnologías en la detección precoz de enfermedades de alto riesgo  
26       del corazón como las cardiopatías. El aprendizaje automático aplicado a los historiales clínicos  
27       puede ser una herramienta útil para predecir el evento de cardiopatía coronaria con síntomas de  
28       enfermedades cardíacas, así como para explorar las características clínicas más significativas y los  
29       factores de riesgo que pueden causar eventos como un infarto o incluso la muerte. Los médicos  
30       pueden aprovechar el aprendizaje automático para clasificar las características clínicas y y desvelar  
31       correlaciones y relaciones ocultas y no evidentes entre los datos de los pacientes.  
32         
33

34 Por lo que, en este artículo se decide estudiar un dataset de 462 historiales médicos, los cuales  
35 pertenecen a pacientes de Sudáfrica. Los objetivos de este estudio son:

- 36 • Investigar y experimentar modelos de redes neuronales del tipo feed forward para predecir  
37 si las personas son propensas o no a padecer de problemas del corazón.
- 38 • Identificar el modelo de aprendizaje automático más eficaz que logra el mejor rendimiento  
39 de predicción en el conjunto de datos dado.

## 40 2 Factores para determinar problemas del corazón

41 Existen múltiples factores de riesgo que pueden contribuir al desarrollo de la cardiopatía isquémica.  
42 En la Tabla 1 se muestra un resumen de los factores más importantes. Aunque la edad, el sexo y  
43 los antecedentes familiares son factores que no pueden cambiarse ni controlarse, reconocerlos como  
44 factores de riesgo puede capacitar a la persona para tomar la iniciativa para vigilar los  
45 factores controlables. El riesgo de problemas del corazón aumenta con la edad y los antecedentes  
46 familiares de enfermedades cardíacas. Los médicos controlan de forma rutinaria muchos factores  
47 de riesgo. Las lecturas de la presión arterial La presión arterial se mide en casi todos los encuentros  
48 clínicos. La presión arterial alta es uno de los efectos secundarios de la restricción del flujo  
49 sanguíneo. La obstrucción arterial que restringe el flujo sanguíneo en los vasos sanguíneos de la  
50 sangre en los vasos sanguíneos provoca hipertensión debido a la mayor resistencia del flujo sanguíneo  
51 cuando el corazón bombea. Las mediciones de la presión arterial superiores a 140/90 se asocian a un  
52 mayor riesgo de enfermedad cardíaca.

Table 1: Factores de riesgo asociados a las enfermedades cardíacas

Factor de riesgo	Controlable
Presión arterial sistólica	Si
Tabaco acumulativo	Si
Colesterol LDL	Si
Antecedentes familiares	No
Obesidad	Si
Alcohol en sangre	Si
Edad	No

## 53 3 Descripción del dataset

54 El conjunto de datos para este estudio se ha obtenido de Sudáfrica, el cual es un subconjunto de  
55 un conjunto de datos más amplio. Tiene un total de 462 observaciones médicas (instancias) y 10  
56 características, 9 como características clínicas independientes, y 1 es la variable objetivo, una clase  
57 binaria etiquetada como 0 o 1, es decir, se ha detectado un evento de problema relacionado al  
58 corazón, es decir, se ha detectado para las observaciones médicas como positivo o negativo. Los  
59 datos corresponden a un grupo de hombres de una zona de alto riesgo de alto riesgo de enfermedades  
60 cardíacas en Sudáfrica. Cada paciente de alto riesgo fue supervisado en el conjunto de datos y las  
61 características recuperadas fueron las siguientes: presión arterial sistólica (Sbp), tabaco acumulado  
62 en kg (Tobacco), colesterol malo también conocido como colesterol de lipoproteínas de baja densidad  
63 (Ldl), adiposidad, antecedentes familiares de enfermedades cardíacas (Famhist), comportamiento  
64 tipo A (TypeA), obesidad, consumo actual de alcohol (Alcohol), edad de inicio (Age) y enfermedad  
65 coronaria (Chd) (sí=1 o no=0).

### 66 3.1 Preprocesamiento

67 El conjunto de datos original se encuentra en formato .dat, por lo que se ha convertido a .csv, y se  
68 ha editado el nombre de las columnas para que sea más expresivas. Se codifica los valores de texto  
69 categóricos existentes en el conjunto de datos original en valores numéricos para que puedan ajustarse  
70 a los modelos de aprendizaje automático.

	sbp	tobacco	ldl	adiposity	famhist	typea	obesity	alcohol	age	chd
0	160	12.00	5.73	23.11	Present	49	25.30	97.20	52	1
1	144	0.01	4.41	28.61	Absent	55	28.87	2.06	63	1
2	118	0.08	3.48	32.28	Present	52	29.14	3.81	46	0
3	170	7.50	6.41	38.03	Present	51	31.99	24.26	58	1
4	134	13.60	3.50	27.78	Present	60	25.99	57.34	49	1

Figure 1: Dataset antes de procesamiento

	sbp	tobacco	ldl	adiposity	famhist	typea	obesity	alcohol	age	chd
0	160	12.00	5.73	23.11	1	49	25.30	97.20	52	1
1	144	0.01	4.41	28.61	0	55	28.87	2.06	63	1
2	118	0.08	3.48	32.28	1	52	29.14	3.81	46	0
3	170	7.50	6.41	38.03	1	51	31.99	24.26	58	1
4	134	13.60	3.50	27.78	1	60	25.99	57.34	49	1

Figure 2: Dataset después de procesamiento

### 3.2 División de datos en entrenamiento y validación

Se divide el dataset en 85% de los datos para el entrenamiento y el resto (15%) para validación.

```
X_train, X_validate, y_train_cat, y_validate_cat = train_test_split(x, y, test_size=0.15)

X_train.shape
(272, 9)

y_train_cat.shape
(272,)

y_validate_cat.shape
(48,)
```

Figure 3: División del dataset

## 4 Metodología y experimentación

La utilidad de un buen modelo de predicción de problemas del corazón depende en gran medida de su precisión y la estabilidad. Para lograrlo, se divide la investigación en tres experimentos, en cada uno se emplea distintas capas para obtener las predicciones más adecuadas y se evalúa tanto para el dataset de entrenamiento como para el dataset de validación.

### 4.1 Experimento 1

Se emplea un modelo sencillo de 2 capas, cada una de ellas empleando como activación ReLU. Se emplea *binary crossentropy* como función de pérdida y *ADAM* como optimizador. Se emplea la métrica de *accuracy* para validar el modelo y se emplean *checkpoint* y *earling stopping* ambos

82 parámetros monitoreando el valor de *accuracy* para los datos de validación con el parámetro de  
 83 paciencia en 1000 iteraciones. Para el entrenamiento se utilizan 500 *epochs* y un *batch size* de 10.

Model: "Experimento_1"		
Layer (type)	Output Shape	Param #
dense_96 (Dense)	(None, 12)	120
dense_97 (Dense)	(None, 9)	117
dense_98 (Dense)	(None, 1)	10
Total params: 247		
Trainable params: 247		
Non-trainable params: 0		

Figure 4: Sumario del modelo para el Experimento 1

#### 84 4.1.1 Resultados

85 Se obtiene que el valor de *accuracy* para los datos de entrenamiento es de 0.6985, mientras que para  
 86 el set de validación es de 0.7708.

```

28/28 [=====] - 0s 1ms/step - loss: 0.5844 - accuracy: 0.6985
5/5 [=====] - 0s 2ms/step - loss: 0.5467 - accuracy: 0.7708
=====

Modelo: Experimento_1.h5

Train accuracy: 69.9%

Test accuracy: 77.1%

=====

```

Figure 5: Resultados para el Experimento 1

## 87 4.2 Experimento 2

88 Se emplea un modelo de 2 capas, cada una de ellas empleando como activación ReLU, la primera  
 89 empleando 64 unidades, seguida de una capa de regularización de batch, posteriormente la segunda  
 90 capa empleando 32 unidades seguido de un dropout del 25 por ciento, por ultimo se obtienen los  
 91 datos empleando sigmoid. Se emplea *binary crossentropy* como función de perdida y *ADAM* como  
 92 optimizador. Se emplea la métrica de *accuracy* para validar el modelo y se emplean *checkpoint* y  
 93 *earling stopping* ambos parámetros monitoreando el valor de *accuracy* para los datos de validación  
 94 con el parámetro de paciencia en 1000 iteraciones. Para el entrenamiento se utilizan 2000 *epochs* y  
 95 un *batch size* de 10.

#### 96 4.2.1 Resultados

97 Se obtiene que el valor de *accuracy* para los datos de entrenamiento es de 0.9301, mientras que para  
 98 el set de validación es de 0.9167.

Model: "Experimento_2"		
Layer (type)	Output Shape	Param #
dense_99 (Dense)	(None, 64)	640
batch_normalization_32 (Batch Normalization)	(None, 64)	256
dense_100 (Dense)	(None, 32)	2080
dropout_30 (Dropout)	(None, 32)	0
dense_101 (Dense)	(None, 1)	33
Total params: 3,009		
Trainable params: 2,881		
Non-trainable params: 128		

Figure 6: Sumario del modelo para el Experimento 2

```

9/9 [=====] - 1s 2ms/step - loss: 0.2297 - accuracy: 0.9301
2/2 [=====] - 0s 5ms/step - loss: 0.3666 - accuracy: 0.9167
=====

Modelo: Experimento_2.h5

Train accuracy: 93.0%

Test accuracy: 91.7%

=====

```

Figure 7: Resultados para el Experimento 2

### 99 4.3 Experimento 3

100 Se emplea un modelo de 4 capas, cada una de ellas empleando como activación ReLU, la primera  
101 empleando 16 unidades, seguida de una capa de regularización de batch y un dropout del 45 por  
102 ciento, posteriormente la segunda capa empleando 8 unidades seguido de un dropout del 45 por ciento,  
103 después se utiliza una capa de 2 unidades seguida de una capa de regularización de batch y por ultimo  
104 se obtienen los datos empleando sigmoid. Se emplea *binary crossentropy* como función de pérdida y  
105 *ADAM* como optimizador. Se emplea la métrica de *accuracy* para validar el modelo y se emplean  
106 *checkpoint* y *earling stopping* ambos parámetros monitoreando el valor de *accuracy* para los datos de  
107 validación con el parámetro de paciencia en 1000 iteraciones. Para el entrenamiento se utilizan 2000  
108 *epochs* y un *batch size* de 20.

#### 109 4.3.1 Resultados

110 Se obtiene que el valor de *accuracy* para los datos de entrenamiento es de 0.6801, mientras que para  
111 el set de validación es de 0.7917.

Model: "Experimento_3"		
Layer (type)	Output Shape	Param #
=====		
dense_102 (Dense)	(None, 16)	160
batch_normalization_33 (Batch Normalization)	(None, 16)	64
dropout_31 (Dropout)	(None, 16)	0
dense_103 (Dense)	(None, 8)	136
dropout_32 (Dropout)	(None, 8)	0
dense_104 (Dense)	(None, 2)	18
batch_normalization_34 (Batch Normalization)	(None, 2)	8
dense_105 (Dense)	(None, 1)	3
=====		
Total params: 389		
Trainable params: 353		
Non-trainable params: 36		

Figure 8: Sumario del modelo para el Experimento 3

```

9/9 [=====] - 0s 2ms/step - loss: 0.6118 - accuracy: 0.6801
2/2 [=====] - 0s 5ms/step - loss: 0.5827 - accuracy: 0.7917
=====

Modelo: Experimento_3.h5

Train accuracy: 68.0%

Test accuracy: 79.2%

=====

```

Figure 9: Resultados para el Experimento 3

#### 112 4.4 Discusión de resultados

113 Se obtiene que el modelo que mejor se desempeña es el del Experimento 2, el cual es el que obtiene  
 114 un mayor valor de *accuracy* tanto para los datos de entrenamiento como con los datos de validación.  
 115 Se puede observar en las gráficas siguientes el comportamiento del aprendizaje para cada uno de los  
 116 modelos. Se puede observar que para el Experimento 1 el *accuracy* se va incrementando y luego  
 117 empieza a disminuir a partir de la iteración 400. Se tiene que para el Experimento 2 el *accuracy*  
 118 incrementa con cada iteración y posteriormente se detiene para evitar *overfitting*. Por último, el  
 119 Experimento 3 se puede observar que se incrementa de forma lenta y luego se mantiene el mismo  
 120 nivel de *accuracy*.

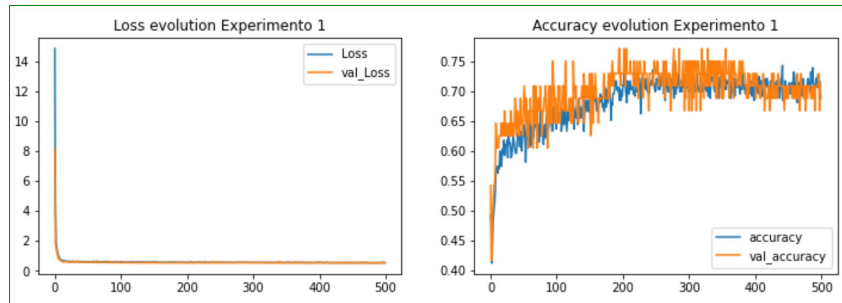


Figure 10: Gráficas de loss y accuracy para el Experimento 1

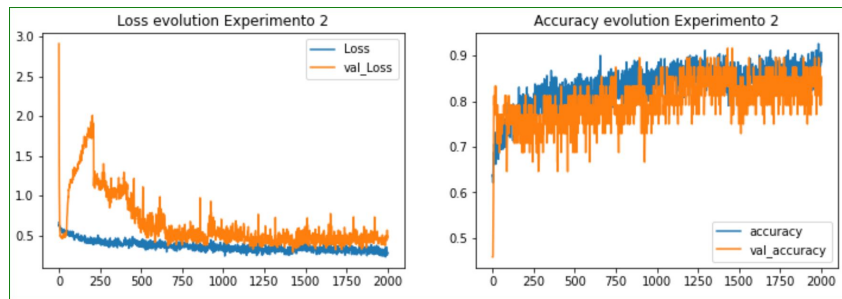


Figure 11: Gráficas de loss y accuracy para el Experimento 2

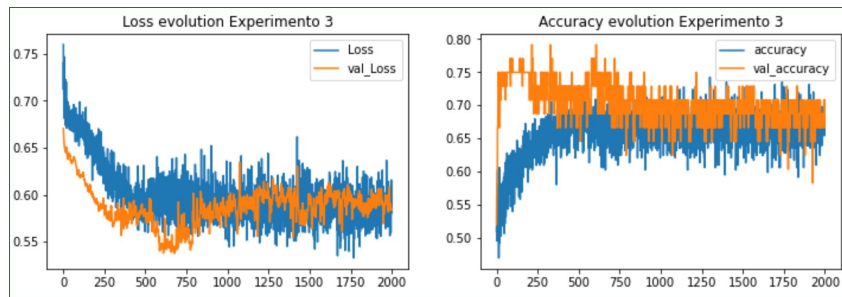


Figure 12: Gráficas de loss y accuracy para el Experimento 3

121 Empleando Keras se obtiene los siguientes resultados del reporte de clasificación, se puede observar  
 122 que el valor de precisión para los valores 0 (Persona no propensa a padecer problemas del corazón)  
 123 es de 0.87, mientras que para los valores de 1 (Persona propensa a padecer problemas del corazón) es  
 124 de 0.68. Verificando el valor de recall y de f1-score se encuentra que todos son arriba de 0.71.

125 Los resultados demuestran que, con datos suficientes y características clínicas seleccionadas, las  
 126 técnicas de aprendizaje automático son capaces de predecir la ocurrencia de eventos de problemas del  
 127 corazón con un porcentaje de accuracy elevado. La selección de las características detalladas en este  
 128 trabajo han demostrado y confirmado que las características clínicas y los factores de riesgo como el  
 129 tabaco, el colesterol LDL, la presión arterial sistólica, la adiposidad y los antecedentes familiares  
 130 se encuentran entre las características más importantes que ayudan a la detección temprana y a la  
 131 predicción de la presencia de eventos de cardiopatía coronaria a partir de los registros médicos. Los  
 132 médicos pueden aprovechar el análisis exploratorio de datos realizado en el conjunto de datos para  
 133 mostrar las correlaciones y relaciones entre los datos de los pacientes.

134  
 135 El éxito del aprendizaje automático depende en gran medida de la riqueza de los datos que  
 136 representan el fenómeno considerado. Aunque el conjunto de datos seleccionado tiene las  
 137 características y los factores de riesgo más conocidos para predecir enfermedades relacionadas con el  
 138 corazón, con un conjunto bastante rico de características, más datos y más variables pueden ayudar a

	precision	recall	f1-score	support
0.0	0.87	0.71	0.78	28
1.0	0.68	0.85	0.76	20
accuracy			0.77	48
macro avg	0.77	0.78	0.77	48
weighted avg	0.79	0.77	0.77	48

Figure 13: Resultados reporte de clasificación experimento con mejores resultados.

139 mejorar los resultados de la predicción. Si se dispusiera de otros conjuntos de datos externos con las  
140 mismas características y procedentes de distintas regiones, no solo limitado a un país en específico,  
141 se podrían llegar a utilizado para validar los resultados obtenidos en este trabajo.

## 142 5 Trabajo futuro

143 Como trabajo futuro, se plantea aplicar el enfoque desarrollado en este trabajo de aprendizaje  
144 automático en otros conjuntos de datos de enfermedades cardiovasculares, cáncer y enfermedades  
145 infecciosas, ya que podrían llegar a obtenerse resultados que puedan llegar a ser favorables para  
146 el área de detección de enfermedades tempranas. También se plantea implementar los modelos  
147 obtenidos como servicio web e integrarlos en una aplicación para que los médicos puedan evaluar su  
148 utilidad en el mundo real.

## 149 Referencias

- 150 [1] Lakhani, P., Sundaram, B. Deep learning at chest radiography: Automated classification of pulmonary  
151 tuberculosis by using convolutional neural networks. Radiology 284, 574–582 (2017).  
152 [2] Yasaka, K., Akai, H. Deep learning with convolutional neural network for differentiation of liver masses at  
153 dynamic contrast-enhanced CT: A preliminary study. Radiology 286, 887–896 (2018).  
154 [3] Indrakumari, R., Poongodi, T., Jena, S. R. (2020). Heart Disease Prediction using Exploratory Data Analysis.  
155 Procedia Computer Science, 173, 130–139.