

Practica 2: Tipología y ciclo de vida de los datos UOC

Álvaro Monforte Marín

January 13, 2023

Contents

1	Pregunta 1	3
2	Pregunta 2	4
3	Pregunta 3	5
4	Pregunta 4	6
5	Pregunta 5	8
6	Pregunta 6	9
7	Pregunta 7	10
8	Pregunta 8	11

1 Pregunta 1

Descripción del dataset. ¿Por qué es importante y qué pregunta/problema pretende responder?

El dataset usado en esta practica se encuentra alojado en esta [pagina web](#) proporcionado en el documento de la practica. A su vez su origen y pertenencia se debe a la Universidad de California Irvine es a un repositorio que se puede encontrar aqui [pagina web del repositorio original](#).

Este dataset nos presenta 14 atributos, entre los cuales tenemos el catorceavo (output) como label objetivo para estudiar como los otros 13 a priori pueden relacionarse o no entre si con tal de obtener una estimacion acerca de si la persona tiene mas riesgo o no de padecer un ataque cardiaco. Ademas posee 303 pacientes implicados (las instancias).

Atributo	Tipo de dato	Informacion ofrecida
Age	int64	Edad del paciente
Sex	int64	Sexo del paciente
cp	int64	Dolor de pecho $\in [0, 4]$
trestbps	int64	Presion de sangre en reposo
chol	int64	Serum de colesterol
fbs	int64	Azucar en sangre (ayunas)
restecg	int64	Electrocardiograma en reposo
thalach	int64	Maximo ratio de pulsaciones
exang	int64	Angina inducida por ejercicio
oldpeak	float64	Depresión del ST ind. ej. rel.
slp	int64	Pendiente del pico de ejercicio de ST
caa	int64	Numero de tejidos coloreados
thall*	int64	"Thallium Stress Test"
output**	int64	Estrechamiento del diámetro

- '*' : Nos informa sobre la condicion del flujo de la sangre en ese momento. Es un parametro medico bastante complejo.
- '**': Nos habla de la probabilidad de padecer un ataque cardiaco. Es el label.

Las muertes por causas isquemicas en el corazon superan los 8 millones de personas anuales. Luego es de alta importancia estudiar cualquier dataset que permita un entendimiento temprano de este tipo de padecimientos. [Fuente](#)

;

2 Pregunta 2

Integración y selección de los datos de interés a analizar. Puede ser el resultado de adicionar diferentes datasets o una subselección útil de los datos originales, en base al objetivo que se quiera conseguir.

Respeto a la integracion y seleccion, creo que todos los datos son sumamente importantes. Quizas ya que tengo en cuenta la presion en sangre, ademas del dano en los tejidos, voy a obviar el azucar en sangre al menos en este estudio. Vease [cierta informacion](#) que ayuda a entender cuan correlados pueden estar estos eventos.

Luego prosigo de inicio con las 303 instancias y ahora 13 atributos a investigar.

3 Pregunta 3

Limpieza de los datos.

- 3.1. ¿Los datos contienen ceros o elementos vacíos?
Gestiona cada uno de estos casos.

El código se ha realizado en lenguaje Python, ayudándome de las librerías Pandas en general aunque se ha hecho también uso de la librería Scipy para algunos cometidos de forma principal, además de sklearn y pingouin para algunos cálculos estadísticos.

Los datos actualmente no contienen elementos vacíos pero sí ceros. Esto se ha logrado mediante la ejecución de la suma de todos los posibles Nan dentro de un dataframe en el que se encuentra la selección de todos los datos. Respecto a los ceros, hay que tener en cuenta que algunos datos en su forma original eran categóricos sin orden (taggeados) finitos 0 o 1; Sin embargo, otros eran datos de tipo likert. Hay alguna excepción de categóricos sin orden de tipo finito por clase, que ayudaran a establecer diferencias respecto a un grupo control con valor 0. Luego los elementos nulos o cero no han sido tratados con dinámicas especiales. Se podría haber también estipulado si existía ciertas correlaciones a la hora de eliminar NaN pero se ha procedido de tal manera que se ha eliminado la instancia completa perdiendo ciertos datos que podrían ser de interés.

- 3.2. Identifica y gestiona los valores extremos

Para la identificación de valores extremos se podría haber procedido por un método gráfico Q-Q que además de ayudar a la eliminación de ciertos outliers también nos podría haber ayudado a entender la normalidad de los datos. Pero se ha hecho uso de un análisis numérico eliminando aquellos en cuantiles superiores al 75% o al 25%, que también sería visible en un diagrama de Box-plot. Con ello, aunque perdiendo información (no siempre los outliers son inútiles), pero como estudio temprano está bien ya que podemos centralizar los datos.

Tras la realización de la limpieza se ha obtenido un total de 229 instancias.

4 Pregunta 4

Análisis de los datos.

- 4.1. Selección de los grupos de datos que se quieren analizar/comparar (p. ej., si se van a comparar grupos de datos, ¿cuáles son estos grupos y qué tipo de análisis se van a aplicar?)
- 4.2. Comprobación de la normalidad y homogeneidad de la varianza.
- 4.3. Aplicación de pruebas estadísticas para comparar los grupos de datos. En función de los datos y el objetivo del estudio, aplicar pruebas de contraste de hipótesis, correlaciones, regresiones, etc. Aplicar al menos tres métodos de análisis diferentes.

La selección de los grupos de interés se puede realizar entorno a las variables categóricas finitas sin orden, para así formalizar subgrupos. En este caso le doy cierto peso a la variable *exang* que determina si hay angina inducida por ejercicio, el por qué de su importancia, tenemos varios papers (estudios) que presentan a esta sintomatología o padecimiento como algo relevante. Lease el artículo [1]. Por ello puedo comparar *exang* contra *trtb*s, *chol* y *thalachh*. Donde tanto como *trtb*s como *chol* son normales mientras que *thalachh* no se distribuye de forma normal. Luego esto sería un indicio para no estudiar de forma paramétrica a *exang* contra *thalachh* si hago una comparación de medias (donde ahora comentare la homocedasticidad) o al menos usar un test que sea robusto como el de Levene.

Una vez se vea si se prestan a una correlación entre ellas, estaría bien contrastar como se comportan respecto al target de la variable *output*, ya que nos queremos adentrar en el estudio de la varianza explicada por estas de la variable a predecir *output*.

Para la normalidad se hizo uso del test de D'Agostino, este tiene como parámetros cruciales la curtosis y la asimetría, hay que tener en cuenta que al descartar outliers centralizando los datos respecto a la mediana podemos provocar una ayuda en la mejor de la curtosis ya que este mide la 'densidad' de puntos respecto al primer momento de inercia, luego al centralizarse a priori debería haber una mayor densidad respecto al momento de giro.

Para la homocedasticidad se hace uso de Bartlett cuando haya normalidad y del test de Levene cuando no la haya, el test de Levene además tiene la propiedad de gran robustez en el momento de que no se presente normalidad.

Para la correlación de las variables, tenemos que tener en cuenta el grupo de supuestos que se cumple, en este caso cuando sean normales se hace uso de la correlación de Pearson, cuando no sea así se hará uso de la de Spearman que es no paramétrica.

Para el test de diferencia de medias t-test, debemos tener en cuenta la homocedasticidad de los datos, cuando esta no se cumple se debe añadir una corrección de Levene

a los grados de libertad de las variables para generar la nueva hipótesis de que sean homocedásticas y así ejecutarlo de tal manera que lo sean, en el algoritmo de python, gracias al módulo pingouin ya está implementada esta corrección pero se debe señalar en mi código con un flag.

Para el estudio de la regresión lineal, las variables deben ser independientes cosa que se puede contrastar con la correlación que existan entre ellas en un mapa de calor. También se debe cumplir que haya una distribución normal de los errores (cosa que normalmente en el campo médico se puede asumir). Para ayudarme a realizar la regresión lineal con el módulo de sklearn me he ayudado de la página web [cienciadedatos](#) que lo explica detalladamente con todo lujo de detalles.

5 Pregunta 5

Representación de los resultados a partir de tablas y gráficas. Este apartado se puede responder a lo largo de la práctica, sin necesidad de concentrar todas las representaciones en este punto de la práctica.

6 Pregunta 6

Resolución del problema. A partir de los resultados obtenidos, ¿cuáles son las conclusiones? ¿Los resultados permiten responder al problema?

7 Pregunta 7

Código. Hay que adjuntar el código, preferiblemente en R, con el que se ha realizado la limpieza, análisis y representación de los datos. Si lo preferís, también podéis trabajar en Python.

8 Pregunta 8

Vídeo. Realizar un breve vídeo explicativo de la práctica (máximo 10 minutos), donde ambos integrantes del equipo expliquen con sus propias palabras el desarrollo de la práctica, basándose en las preguntas del enunciado para justificar y explicar el código desarrollado. Este vídeo se deberá entregar a través de un enlace al Google Drive de la UOC (<https://drive.google.com/...>), junto con enlace al repositorio Git entregado.

References

- [1] Alon Eisen, Deepak L. Bhatt, P. Gabriel Steg, Kim A. Eagle, Shinya Goto, Jianping Guo, Sidney C. Smith, E. Magnus Ohman, Benjamin M. Scirica, and null null. Angina and future cardiovascular events in stable patients with coronary artery disease: Insights from the reduction of atherothrombosis for continued health (reach) registry. *Journal of the American Heart Association*, 5(10):e004080, 2016. doi: 10.1161/JAHA.116.004080. URL <https://www.ahajournals.org/doi/abs/10.1161/JAHA.116.004080>.