

Tarea 3

INSTRUCCIONES

Estas pruebas están basadas en la filosofía de aprender haciendo y la intención es que tratando de realizar los problemas, ejercicios o análisis, los estudiantes ejerciten e interioricen los conocimientos teóricos que han adquirido mediante el estudio de los textos, los tutoriales y los materiales complementarios. Al tratarse de una prueba evaluable se trata de un trabajo individual, es de este modo, pensando los problemas por su cuenta e investigando, como lograrán una buena comprensión de la materia.

Entrega: Las soluciones a los problemas han de ser redactados en un documento realizado mediante una aplicación informática que permita incluir las fórmulas y los desarrollos de los problemas. Existen múltiples opciones como el editor de ecuaciones de Microsoft Word Office, o software libre como LibreOffice. También es posible redactarlas empleando LaTeX. Es recomendable este modo ya que ofrece la mejor presentación para textos científicos y su aprendizaje será de enorme utilidad para aquellos que desarrollan una carrera en Ciencias e Ingeniería. Si deciden tomar esta opción en el curso virtual hay un enlace con un excelente tutorial de LaTeX y una plantilla que pueden usar si lo desean.

La empresa Makichan ha puesto recientemente la posibilidad de acceder de forma gratuita al software Scientific Word. Se trata de un editor de LaTeX que facilita enormemente el manejo de ecuaciones:

<https://www.mackichan.com/techtalk/v60/FreeSW.htm>

En ningún caso se permitirá la entrega de documentos redactados a mano y escaneados.

A la plataforma debe subirse un archivo .zip .rar .tar o equivalente en el que estén contenidos la memoria en pdf con el formato adecuado y también la fuente: el archivo .doc o .tex o equivalente del que se ha obtenido el pdf. También debe subir el archivo de LiveScript de Matlab.mlx con el que ha realizado los cálculos y exportado a pdf.

Trabajo individual: Se trata de una prueba calificable y por tanto individual. No se permitó el intercambio de información entre los estudiantes, que además iría en detrimento de su aprendizaje. El equipo docente, como es su obligación, velará porque esto sea así en la corrección y utilizando las herramientas antiplagio de las que dispone la institución. En caso de ser detectado plagio se verá obligado a actuar al respecto tomando las medidas pertinentes sobre todos los agentes implicados.

1**ESTUDIO DE CONJUNTOS DE DATOS**

Esta tercera tarea está enfocada en el tratamiento en detalle de un conjunto de datos. De este conjunto debe obtener la información que considere pertinente o interesante implementando las diferentes técnicas que se han estudiado hasta el momento: métodos estadísticos, regresión lineal y logística, clusterización, reglas de asociación, arboles de decisión, etc... Recuerde que no todas las técnicas son de utilidad en todos los problemas y conjuntos de datos, e incluso puede carecer de sentido aplicarlas.

El conjunto de datos que se va a estudiar es clásico como ejercicio en Ciencia de Datos y Aprendizaje Automático, es posible que encuentre en la red mucha información al respecto y es libre de estudiarla y aprender de ella, pero no es válido reproducirla literalmente. Debe aportar sus propias soluciones y deben justificar las técnicas que usen en función de los resultados que quieran obtener o la información que crean que resulta interesante. Aunque hay ciertos conjuntos de datos -como el que nos ocupa en este ejercicio- que llevan mucho tiempo en circulación y han sido muy estudiados es habitual que de ellos se extraiga nueva información, muchas veces aplicando nuevas técnicas o realizando de forma inteligente transformaciones de los datos, seleccionando subconjuntos interesantes, etc. La información que se puede obtener de un conjunto de datos no está necesariamente relacionada de forma directa con el tamaño del mismo. La información obtenible depende de la calidad de los datos, lo representativos que sean del problema a tratar y la habilidad del científico de datos para aplicar las técnicas adecuadas y llegar a conclusiones con los resultados obtenidos.

Tenga en cuenta que puede resultar interesante hacer transformaciones de los datos, por ejemplo, obtener una nueva columna que sea el ratio entre dos cantidades, su producto o su suma.

En las partes en las que realice ejercicios de clasificación debe tener en cuenta que en la primera tarea programó medidas de aciertos y fallos en la regresión logística, como se indicó, son de suma importancia ya que su uso se extiende a todo problema de clasificación. En estos ejercicios debe implementar, donde sea necesario, medidas equivalentes para medir la eficiencia de los modelos que ha propuesto e implementado.

Todos los cálculos y transformaciones deben ser realizados mediante el software Matlab.

Finalmente debe redactar una memoria en la que debe seguir el esquema de presentar una introducción, luego un desarrollo y finalmente unas conclusiones.

Un esquema posible para la memoria puede ser:

Introducción: Donde se presenta el problema, los datos disponibles y lo que se pretende averiguar.

Tratamiento de datos/Resultados: Donde se explican las técnicas que se usan, y los resultados que se obtienen. De forma crítica por que se escogen unas y no otras, como los resultados de alguna de ellas motivan la elección de otras, etc.

Conclusiones: Un resumen de las conclusiones más importantes que se han obtenido mediante el tratamiento de los datos. Las conclusiones no deben ser muchas, no más de cuatro o cinco y cada una de ellas debe, en unas pocas líneas, explicar un resultado importante.

En muchos trabajos académicos es bastante habitual que lo último que se escriba sea la introducción: durante la redacción del tratamiento de datos y resultados se va perfilando que es lo más importante del trabajo realizado y que conclusiones se pueden sacar, posteriormente se pueden escribir estas y finalmente la introducción que de una idea del trabajo con el que se va a encontrar el lector.

La memoria del estudio a tratar **no debe ocupar más de 12 páginas en total**, incluyendo las figuras: pueden ser necesarias menos. No se trata de hacer un trabajo largo, sino de aprender a seleccionar que técnicas usar y que resultados es importante poner en relevancia en el trabajo. Seguramente haya muchos cálculos de prueba que no merezca la pena incluir en la memoria o resultados menores que sólo merezcan ser comentados en unas líneas, pero que no requieran la atención de aparecer en figuras, tablas, etc. No se trata de hacer el mayor número posible de cálculos y aplicar una gran cantidad de técnicas, sino de ver cuáles son interesantes y se valorará la capacidad de síntesis al exponer los resultados y la justificación de utilizar un método u otro. Se debe tener una actitud crítica.

Consulte el documento de orientaciones para la redacción de las tareas donde encontrará indicaciones sobre diversos detalles, formato, etc.

Le resultará también útil la Guía para la Redacción de Trabajos Académicos que la Biblioteca de la UNED ha confeccionado:

https://uned.libguides.com/trabajos_academicos

Los conjuntos de datos están disponibles en el curso virtual ya que se han incluido en algunos ejemplos de los tutoriales de Matlab del material complementario y el conjunto Iris dataset viene incluido en Matlab.

1. 1 Conjunto de datos: Bank Dataset.

El conjunto de datos de Bank Dataset es un conjunto de datos mediano con 2000 entradas. Contiene información de los clientes de un banco, con los que se ha contactado o se pretende contactar, para una campaña de marketing en la que se intenta que se suscriban a un determinado producto.

age: edad

job: tipo de trabajo

marital: estado marital

education: nivel de educación

default: si han incurrido en impago

balance: balance de cuenta

housing: casa en propiedad

loan: si tienen un préstamo

contact: tipo de teléfono de contacto.

day: día del contacto

month: mes del contacto

duration: duración en segundos de la llamada de contacto

campaign: número de contactos realizados durante esta campaña de marketing

pdays: número de días que han pasado desde que el cliente fue contactado por última vez

previous: número de contactos realizados previos esta campaña de marketing

poutcome: resultado de la anterior campaña de marketing con el cliente: si fue exitosa, falló o no hubo respuesta.

subscribed: si el cliente se ha suscrito al producto.

Es un conjunto de datos que se ha usado frecuentemente como modelo en problemas de predicción, pero también puede tener interés desde el punto de estudiar si entre las diferentes columnas de datos, puede establecerse algún tipo de relación, es posible clasificar a los clientes en función de sus ingresos

Hay total libertad para experimentar con los datos y obtener información que considere relevante del conjunto o determinar si unas técnicas funcionan mejor que otras para un determinado propósito, etc.

Junto al conjunto de datos de modelización y entrenamiento, con 2000 entradas, hay otro conjunto

extraído de la misma fuente que es el conjunto de test con 100 entradas. En este caso es muy importante medir el éxito de los posibles análisis que se hayan podido realizar, con un conjunto de test que sea diferente de aquel con el que se han entrenado los modelos.

En el conjunto de test hay categorías de datos que no están muy bien representadas ya que es pequeño. Puede ser interesante fraccionar el conjunto de entrenamiento con la idea de extraer de él un conjunto de test diferente.

Como el resto de los conjuntos de datos utilizados en el texto base se puede descargar en la página web del libro

<https://www.wiley.com/en-us/Data+Science+and+Big+Data+Analytics%3A+Discovering%2C+Analyzing%2C+Visualizing+and+Presenting+Data-p-9781118876138>