

Curso de aprendizaje automatizado

PCIC, UNAM

Tarea 3: Regresión y clasificación lineal

Fecha límite: 15 de abril.

Formato: Archivo ZIP con código fuente y reporte en PDF.

Forma de entrega: Enviar correo electrónico a gibranfp@unam.mx, bere.mctic@gmail.com y richardt.pcic@gmail.com con asunto *Tarea 3: Regresión y clasificación lineal* y archivo ZIP adjunto.

Descripción

Realiza los siguientes ejercicios de regresión y clasificación lineal:

Predicción de precios de casas

A partir de la base de datos de precios de casas de Boston¹ (*Boston Housing Data*), realiza la regresión de los precios de las casas con las siguientes variantes:

- Mínimos cuadrados con expansión polinomial de diferentes grados.
- Mínimos cuadrados con expansión polinomial de grado 20 y penalización por norma ℓ_1 y ℓ_2 con diferentes valores de λ .
- Mínimos cuadrados con expansión polinomial de grado 2 y selección de atributos.

Grafica el error cuadrático medio en entrenamiento y validación con respecto al grado del polinomio, valor de λ y número de atributos. Todos los modelos deberán ser evaluados con 10 repeticiones de validación cruzada de 5 particiones.

Predicción de juegos

Un club del juego de Go recopiló los resultados de varias partidas entre diferentes jugadores, almacenados en el archivo `juegos_entrenamiento.txt`, con el objetivo de predecir el resultado de partidas futuras, ejemplos de las cuales se encuentran en el archivo `juegos_validacion.txt`. Los

¹La base de datos se encuentra disponible en <http://archive.ics.uci.edu/ml/machine-learning-databases/housing/housing.data> y la descripción en <http://archive.ics.uci.edu/ml/machine-learning-databases/housing/housing.names>

archivos `juegos_entrenamiento.txt` y `juegos_validacion.txt`² contienen 3 columnas: la primera corresponde al identificador del jugador A, la segunda al identificador del jugador B y la tercera es el resultado de la partida (1 si ganó el jugador A o 0 si ganó el jugador B). En el club hay un total de D jugadores, por lo que cada identificador es un número entero entre 1 y D . La predicción del resultado de un juego se puede plantear como un problema de clasificación: dados 2 jugadores (A y B) se requiere predecir si A ganó ($y = 1$) o si fue B ($y = 0$). Realice los siguientes ejercicios:

- Entrena y evalúa un clasificador bayesiano ingenuo. Al ser un modelo generativo (modela la probabilidad conjunta $P(\mathbf{x}, y)$), es posible generar partidas artificiales con los parámetros calculados. Genera nuevas partidas que sigan la distribución modelada.
- Entrena y evalúa un clasificador de regresión logística³. Para esto es necesario reparametrizar las entradas. Explica el procedimiento y la lógica de la reparametrización que realizaste. La Selecciona y visualiza los valores de los parámetros. Grafica las curvas ROC y de precisión-exhaustividad y reporta sus áreas bajo la curva.
- Compara el clasificador bayesiano ingenuo y regresión logística en este problema. ¿Qué ventajas y desventajas tienen los modelos entrenados? ¿Qué pasaría si se entrena el clasificador bayesiano ingenuo con los vectores reparametrizados o si se entrena un modelo de regresión logística usando los vectores de entrada originales? ¿Consideras que las presuposiciones de cada clasificador son apropiadas para los datos del problema? ¿Para este tipo de problemas cuál de los dos recomendarías y por qué?
- Deriva la regla de actualización para el algoritmo del descenso por gradiente de un clasificador donde $\hat{y} = \text{sigm}(\boldsymbol{\theta}^\top \mathbf{x})$ y la función de pérdida sea

$$E(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n \left(\hat{y}^{(i)} - y^{(i)} \right)^2.$$

Discute las diferencias entre este clasificador y el de regresión logística y compara sus rendimientos en la tarea de predicción de juegos. (**Ejercicio opcional, 2 puntos extras**)

Regresión logística vs clasificador bayesiano ingenuo

Compara los métodos de regresión logística⁴ y el clasificador bayesiano ingenuo en las siguientes tareas:

- *Clasificación de spam*⁵
- *Clasificación de tumores de seno*⁶

Discute qué modelo seleccionarías y por qué. Todos los modelos deberán ser evaluados con 10 repeticiones de validación cruzada estratificada de 5 particiones.

²Estos archivos se encuentran disponibles en http://turing.iimas.unam.mx/~gibranfp/cursos/aprendizaje_automatizado/data/regl_data.zip

³Se espera que el clasificador por regresión logística se programe usando únicamente las bibliotecas NumPy y SciPy de Python.

⁴Se espera que el clasificador de regresión logística se programe usando únicamente las bibliotecas NumPy y SciPy de Python.

⁵Con el conjunto de datos disponible en http://turing.iimas.unam.mx/~gibranfp/cursos/aprendizaje_automatizado/data/nb_data.zip

⁶Con el conjunto de datos disponible en <http://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.data> y la descripción en <http://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.names>