

MACHINE LEARNING PARA ECONOMISTAS - TRABAJO PRÁCTICO N°4

Clasificación y regularización de desocupación usando la EPH

Fecha de entrega: 17 de diciembre de 2024

Profesor: Walter Sosa Escudero

Asistente: Tomás Pacheco

Contenido: Análisis de hogares para los determinantes de la desocupación, problema de clasificación de desempleo entre cohortes usando métodos de regularización y elección de hiperparámetros por cross-validation.

Reglas de formato y presentación

Al finalizar el trabajo práctico deben hacer un último *commit* en su repositorio de GitHub con el mensaje Entrega final del TP.

- Asegúrense de haber creado una carpeta llamada TP4. Deben entregar un reporte (pdf) y el código (puede ser Jupiter Notebook o .py). Ambos deben estar dentro de esa carpeta.
- La última versión en el repositorio es la que será evaluada. Por lo que es importante que:
 - No envíen el mensaje hasta no haber terminado y estar seguros de que han hecho el *commit* y *push* a la versión final que quieren entregar.
 - No hagan nuevos *push* después de haber entregado su trabajo. Esto generaría confusión acerca de qué versión es la que quieren que se les corrija.
- El informe debe ser entregado en formato PDF, con los gráficos e imágenes en este mismo archivo. Puede tener una extensión máxima de hasta 8 paginas (no se permite Apéndice). Se espera una buena redacción en la resolución.
- Entregar el código con los comandos utilizados, identificando claramente a que inciso corresponde cada comando.

Parte I: Análisis de la base de hogares y tipo de ocupación

Ahora que ya están familiarizados con la Encuesta Permanente de Hogares (EPH) y la desocupación, vamos a complejizar un poco la construcción de las tasas del desempleo. Relacionaremos la información a nivel hogar.

1. Explore el diseño de registro de la base de hogar: a priori, ¿qué variables creen pueden ser predictivas de la desocupación y sería útil incluir para perfeccionar el ejercicio del TP3? Mencionen estas variables y justifiquen su elección.
2. Descarguen la base de microdatos de la EPH correspondiente al primer trimestre de 2004 y 2024 en formato `.dta` y `.xls`, respectivamente. La base de hogares se llama `Hogar_t104.dta` y `usu_hogar_T124.xls`, respectivamente. Eliminen todas las observaciones que no corresponden al aglomerado con el que están trabajando¹ y unan ambos trimestres en una sola base. Unan, a la base de la encuesta individual de cada año, la base de la encuesta de hogar. Asegúrese de estar usando las variables `CODUSU` y `NRO_Hogar` para el merge.
3. Limpie la base de datos tomando criterios que hagan sentido. Explique cualquier decisión como el tratamiento de valores faltantes (missing values), extremos (outliers), o variables categóricas. Justifique sus decisiones.
4. Construyan variables (mínimo 3) que no estén en la base pero que sean relevantes para predecir individuos desocupados (por ejemplo, la proporción de personas que trabajan en el hogar).
5. Presenten estadísticas descriptivas de cinco variables de la encuesta de hogar que ustedes creen que pueden ser relevantes para predecir desocupación. Comenten las estadísticas obtenidas.
6. En el TP3 calcularon la tasa de desocupación según INDEC y economía laboral, para el 1er trimestre de 2024. Utilice una sola observación por hogar y sumen el ponderador `PONDERA` que permite expandir la muestra de la EPH al total de la población que representa ¿Cuál es la tasa de hogares con desocupación para su aglomerado? ¿se asemeja dicha tasa a la reportada [en el INDEC en sus informes](#)?

Parte II: Clasificación y regularización

El objetivo de esta parte del trabajo es nuevamente intentar predecir si una persona está desocupada o no. Esta vez utilizando distintas variables de características individuales y preguntas de la encuesta de hogar. A su vez incluiremos ejercicios de regularización y de validación cruzada.

1. Para cada año, partan la base respondieron en una base de prueba y una de entrenamiento (`X_train`, `y_train`, `X_test`, `y_test`) utilizando el comando `train_test_split`. La base de entrenamiento debe comprender el 70 % de los datos, y la semilla a utilizar (*random state instance*) debe ser 101. Establezca a desocupado como su variable dependiente en la base de entrenamiento (vector `y`). El resto de las variables serán las variables independientes (matriz `X`). Recuerden agregar la columna de unos (1).
2. Expliquen como elegirían λ por validación cruzada. Detallen por qué no usarían el

¹ver Tabla (1)

conjunto de prueba (test) para su elección.

3. En validación cruzada, ¿cuáles son las implicancias de usar un k muy pequeño o uno muy grande? Cuando $k = n$ (con n el número de muestras), ¿cuántas veces se estima el modelo?
4. Implementen la penalidad, L1 como la de LASSO y L2 como la de Ridge, para regresión logística usando la opción `penalty` y reporten la matriz de confusión, la curva ROC, los valores de AUC y de Accuracy para cada año.² ¿Cómo cambiaron los resultados con respecto al TP3? ¿Mejor o empeoró la performance de regresión logística con regularización?
5. Realicen un barrido en $\lambda = 10^n$ con $n \in \{-5, -4, -3, \dots, +4, +5\}$ y utilicen 10-fold CV para elegir el λ óptimo en regresión logística con Ridge y con LASSO. ¿Qué λ seleccionó en cada caso? Usando la librería de `seaborn`, generen `box plot` mostrando la distribución del error de predicción para cada λ . Cada box debe corresponder a un valor de λ y contener como observaciones el error medio de validación para cada partición. Además, para la regularización LASSO, generen un line plot, pero ahora graficando el promedio de la proporción de variables ignoradas por el modelo en función de λ , es decir la proporción de variables para las cuales el coeficiente asociado es cero³.
6. En el caso del valor óptimo de λ para LASSO encontrado en el inciso anterior, ¿qué variables fueron descartadas? ¿Son las que hubieran esperado? ¿Tiene relación con lo que respondieron en el inciso 1 de la Parte I?
7. Elijan alguno de los modelos de regresión logística donde hayan probado distintos parámetros de regularización y comenten: comparen los resultados de 2004 versus 2024, ¿qué método de regularización funcionó mejor: Ridge o LASSO? Comenten mencionando el error cuadrático medio (ECM).

Cuadro 1: Grupos y aglomerados

Grupo	Nombre aglomerado	Código
1	Posadas	07
2	Mar del Plata - Batán	34
3	Gran Córdoba	13
4	Gran La Plata	02
5	Ciudad de Bs As	32
6	Bahía Blanca - Cerri	04
7	Gran San Juan	27
8	Gran Tucumán - T. Viejo	29
9	Partidos del GBA	33

²En la tutorial 6, vimos el método de regularización en regresión lineal donde la variable dependiente es numérica. En este caso, nuestra variable dependiente es binaria (ocupado, desocupado), por lo que usamos la regresión logística y aprovechamos la opción de penalidad para aplicar los métodos de regularización visto en clase.

³*Hint:* a mayor penalidad, esperamos que más coeficientes sean 0, por lo tanto, esta figura debe tener una forma de "S".