

ENTRE ALHÓNDIGAS, BESOS Y MOMIAS

Dr. A. Pastor López (CIMAT), Dr. Rafael Guerrero (DCEA-Universidad de Guanajuato)

Entregar (V1): Martes 4 de Noviembre de 2025 antes de las 19:00hrs

Contexto

Realiza los siguientes puntos en un notebook de Python *lo mejor organizado y claro posible*. Ponga su nombre al notebook (e.g., fulanito1_fulanita2.ipynb) y también en la primera celda del notebook junto ponga los nombres de TODO el equipo. Sube al classroom de moodle el notebook como un archivo, que deberá haber sido ejecutado en tú máquina y mostrar el resultado en las celdas.

Para este proyecto consideraremos el conjunto de datos recolectado por el equipo del Dr. Rafael Guerrero, profesor en la División de Ciencias Económico Administrativas de la Universidad de Guanajuato. Este conjunto de datos contiene un aproximado de diez mil opiniones de turistas en trip advisor en 10 sitios turísticos de la ciudad de Guanajuato. **El objetivo es realizar las siguientes actividades y contestar las preguntas.** Para esta tarea se puede usar cualquier librería o herramienta de Python (e.g., sklearn, keras, nltk, códigos de github de otras personas (citando), etc.). También puede reusar su código de tareas previas, o puede simplemente usar TfidfVectorizer, CountVectorizer, etc. de sklearn. Puede usar también SelectKBest como en el Lecture de DOR lo hizo el profesor, o usar su propio código de Chi2.¹

Para estas actividades usted determine el número de features (palabras) de alguna forma según su intuición. Usualmente el top 10k con base en frecuencia podría ser buena elección si su hardware es suficiente para llevar a cabo las actividades. Si su reducción es a 5k o menos términos, algo con base en Chi2, Ganancia de Información o valores TFIDFs podría venir mejor para no perder tanta información y llegar a buenas conclusiones. Este Proyecto/Examen puede ser individual o en equipos de hasta 4 personas máximo. NO SE VALE COPIAR; EN CASO DE ENCONTRAR COPIAS ES CERO PARA TODOS LOS MIEMBROS DE LOS EQUIPOS INVOLUCRADOS.

Actividades (50pts)

1. (+2.5pts extra) Construya estadísticas básicas respecto a la opinión de cada lugar turístico. **Preprocese y limpie el texto según sus intuiciones y argumente brevemente sobre ello.** Considere scores de 4 a 5 como **positivos**, calificaciones de 3 como **neutros** y las de 2 a 1 como **negativos**. Es interesante ver:

¹Recomiendo ampliamente usar lo más posible las funciones de Sklearn, para aprender a usarlas además de que son muy eficientes al llevar todo en matrices sparse. Esto hará que puedas manipular vocabularios enormes y más rápido.

- (a) Promedios de calificación por lugar, y desviaciones estándar en los scores
 - (b) Basado en palabras: longitud promedio de opiniones y desviaciones estándar
 - (c) Histogramas de edades de opiniones por lugar
 - (d) Histograma de tipo de visitantes (nacional o internacional) por lugar
 - (e) **Sugiere dos más interesantes para ti.**
2. (5pts) Utilizando una estrategia de feature selection (se sugiere χ^2 o ganancia de información) visualice con *word_cloud* (https://amueller.github.io/word_cloud/) nubes de palabras el top k (se sugiere 50) de palabras más relevantes para cada uno de los 10 lugares. Note que serán 10 nubes, una por lugar.
3. (15pts) Para cada uno de los 10 sitios turísticos, haga un descubrimiento automático de los 3 tópicos con Latent Semantic Analysis (LSA) (investiga, estudia y aprende por su cuenta LSA) más relevantes y 10 palabras contenidas en cada tópico de cada uno de los siguientes subgrupos:
- (a) Hombres
 - (b) Mujeres
 - (c) Turistas Nacionales
 - (d) Turistas Internacionales
 - (e) Jóvenes (elige un rango de edad interesante con base en sus estadísticas)
 - (f) Mayores (elige un rango de edad interesante con base en sus estadísticas)

Antes de aplicar LSA, asegúrese de hacerlo sobre una matriz lo más grande posible (para su hardware) de TFIDF Normalizada a L2. Note que para cada sitio turístico deberá saber cuales son los 3 temas de interés y sus palabras, para cada uno de estos subgrupos. Recomiendo Gensim; para rápido y fácil. Otra sugerencia puede ser usar la función TruncatedSVD de sklearn para obtener la descomposición de matrices como se sugiere en el siguiente video para implementar LSA: <https://www.youtube.com/watch?v=hB51kkus-Rc>. También podría llevar a cabo svd con numpy.

4. (10pts) Para cada uno de los 10 sitios turísticos construya tres Bolsas de Palabras tfidf de la siguiente manera: i) 1000 términos con mayor peso (χ^2 con respecto al género), ii) 2000 bigramas con mayor peso (χ^2 con respecto al género), y iii) 1000 trigramas con mayor peso (χ^2 con respecto al género). Luego concatene las tres representaciones que fueron calculadas de forma independiente, con sus propios tfidfs según su espacio y su propio L2. Finalmente sobre todo ese espacio concatenado de 4000 características aplique ganancia de información o χ^2 y obtenga los 1000 features más relevantes. Muestre una nube de palabras con el top 50 features relevantes para cada lugar turístico (10 nubes en total, quiero ver, para cada sitio, en la misma nube los mejores uni-bi-trigramas cuando se aplica feature selection en el espacio global (con respecto al género)).

5. (10pts) Diseñe un análisis temporal (formato libre) que muestre opiniones positivas, negativas y neutras a través de los meses y años para todos los sitios turísticos. En pocas palabras mostrar la evolución de las opiniones a través del tiempo.
6. (10pts) Haga una partición 70% (train) 30% (test) para evaluar clasificación en el dataset con respecto a las 3 clases (positiva, negativa y neutro). Evalúe un SVM Lineal con grid-search con las siguientes representaciones y HAGA UNA TABLA COMPARATIVA:
 - Bolsa de Palabras Binaria con L2.
 - Bolsa de Palabras Frecuencia con L2.
 - Bolsa de Palabras TFIDF con L2.
 - Proponga una representación basada en word vectors DOR con L2
 - Proponga una representación basada en word vectors TCOR con L2
 - Proponga una representación basada en word vectors Word2Vec con L2
 - (Opcional; 10 puntos extra en el examen) Proponga un método de clasificación de su elección para derrotar al mejor de todos los anteriores. Puede usar códigos de otros autores, puede usar LLMs con prompting, puede usar modelos de Hugging-Face/Paperswithcode.

1 (50pts) Preguntas: Conteste lo más detallado posible lo siguiente, dando argumentos y conclusiones claras según su análisis previo. Cada respuesta entre 150 (mínimo) y 300 (máximo) palabras.

1. (10pts) ¿De los sitios turísticos, cuál diría usted que es el más polémico y **la razón de ello?**
2. (10pts) En cuanto al sitio más polémico, ¿Cómo es la diferencia de opinión y temas entre turistas nacionales e internacionales?
3. (10pts) ¿Cuál diría que es el sitio que le gusta más a las mujeres y por qué?
4. (10pts) ¿Cuál diría que es el sitio que le gusta más a las personas jóvenes y por qué?
5. (10pts) ¿Qué otras observaciones valiosas puede obtener de su análisis? (e.g., ¿identificó de qué se queja la gente? ¿qué tipo de cosas le gustó a la gente?, etc.)