

PROYECTO 1: MINERÍA DE TEXTO PARA TURISMO

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

Entregar: Viernes 12 de Mayo de 2023 antes de las 23:59:59

Contexto

Realiza los siguientes puntos en un notebook de Python *lo mejor organizado y claro posible*. Ponga su nombre al notebook (e.g., `adrian_pastor_lopez_monroy.ipynb`) y también en la primera celda del notebook junto con el número de Tarea. Sube al classroom el notebook como un archivo, que deberá haber sido ejecutado en tu 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 es INDIVIDUAL.

Actividades (50pts)

1. (2.5pts) Construya estadísticas básicas respecto a la opinión de cada lugar turístico. **Pre-procese 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:
 - (a) Promedios de calificación por lugar, y desviaciones estándar en los scores
 - (b) Basado en palabras: longitud promedio de opiniones y desviaciones estándar

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.

- (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. (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 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. Como sugerencia puede 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. Otra sugerencia para LSA es investigar a hacerlo con la librería gensim de python.

4. (5pts) Para cada uno de los 10 sitios turísticos, haga una nube de palabras que muestre las palabras más asociadas a sus opiniones negativas utilizando χ^2 . Puede usar funciones de sklearn.
5. (15pts) Para cada uno de los 10 sitios turísticos construya tres Bolsas de Palabras de la siguiente manera: i) 1000 términos con mayor peso tfidf, ii) 2000 bigramas con mayor tfidf, y iii) 1000 trigramas con mayor tfidf. 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).

6. (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.

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, cual diría usted que es el más polémico y **la razón de ello?**
2. (10pts) En cuanto al sitio más polémico, ¿Como es la diferencia de opinión y temas entre turistas nacionales e internacionales?
3. (10pts) ¿Cual diría que es el sitio que le gusta más a las mujeres y por qué?
4. (10pts) ¿Cual 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 que se queja la gente? ¿qué tipo de cosas le gustó a la gente?, etc.)