

- **N-gramas:** esta propuesta se hizo tanto para determinar el género como la variedad lingüística. Lo que hicimos fue modificar el

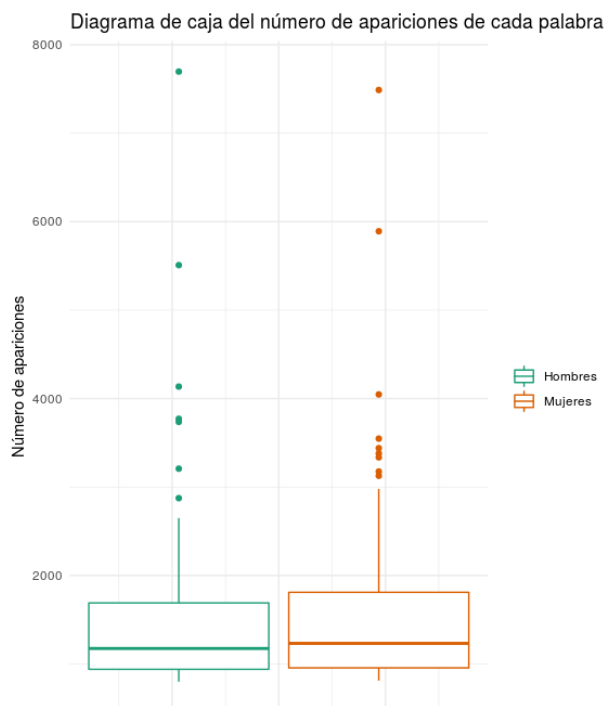


Figura 2: Diagrama de cajas por genero para las cien palabras más frecuentes.

script dado en R para generar una bolsa con los n-gramas más frecuentes. Debido a la cantidad de información en R no conseguimos llegar a un resultado ya que no terminaba de ejecutarse. La hipótesis fue que habíamos visto en una asignatura del máster como el aplicar 2-gramas y 3-gramas mejoraba el resultado usando solamente bolsa de palabras.

- **Terminaciones de cada palabra:** esta propuesta se hizo para el problema de determinación del genero. Lo que hicimos fue implementar un script en R para quedarse con los sufijos más discriminantes, como las terminaciones -o -a, pero al igual que con los n-gramas no conseguimos llegar a un resultado en R. La hipótesis es que como en el idioma español existen sufijos distinto para reflejar el genero masculino o femenino, bajo el supuesto de que un autor al referirse a sí mismo usará las palabras de acuerdo a su sexo, entonces el usar la frecuencia de aparición de los sufijos seleccionados podría discriminar el sexo del autor del twitt.
- **Longitud del Twitt:** esta propuesta se hizo para el problema de determinación del genero. La idea es determinar el genero únicamen-

te usando la variable que determina el número de letras, símbolos de puntuación, emoticonos y espacios en blanco. La hipótesis está basada en estudios psicológicos que muestran que en términos generales las mujeres hablan más que los hombres, trasladando esta hipótesis a nuestro caso, supusimos que las mujeres escribirían twitts más largos.

- **Bolsa de palabras con las n palabras más frecuentes:** esta propuesta se hizo tanto para determinar el género como la variedad lingüística. Usamos la baseline usando el algoritmo de clasificación *RandomForest* y variando el número n de palabras. La hipótesis fue que dado que la bolsa de palabras ya aportaba buenos resultados usando el algoritmo *Support Vector Machines* (SVM) y con la experiencia de los casos vistos en el máster, de que el algoritmo *RandomForest* proporciona mejores resultados, decidimos cambiar el algoritmo.
- **Bolsa de palabras determinadas previamente en base a temáticas:** esta propuesta se hizo para el problema de determinación del genero. Se buscaron palabras de temas como deportes, creyendo que esta bolsa discriminaría a los twitts escritos por hombres, y otros como moda, bajo la hipótesis de que discriminaría mucho mejor a los twitts escritos por mujeres. Al ser una elección subjetiva de temas para hombre y mujeres, al igual que es complejo elegir unas palabras que representen estos temas, la aproximación fue muy superficial. La hipótesis fue que los hombres y mujeres hablan de temas diferentes concretos.
- **Term frequency – Inverse document frequency (Tf-idf):** esta propuesta se hizo para el problema de determinación del genero así como para la variedad lingüística. El método es similar al de bolsa de palabras, salvo por la diferencia de que en este se tiene en cuenta también la frecuencia de ocurrencia del término en la colección de documentos, así que lo que tenemos es una medida numérica que expresa cuán relevante es una palabra para un documento en una colección. Este método lo implementamos y aplicamos en Python, ya que nos resultó más fácil al haber trabajado antes con este método en esa plataforma. La

Figura 3: Siguiendo el orden de izquierda a derecha y de arriba a abajo las variedades lingüísticas son: Argentina, Chile, Colombia, México, Perú, España, Venezuela.

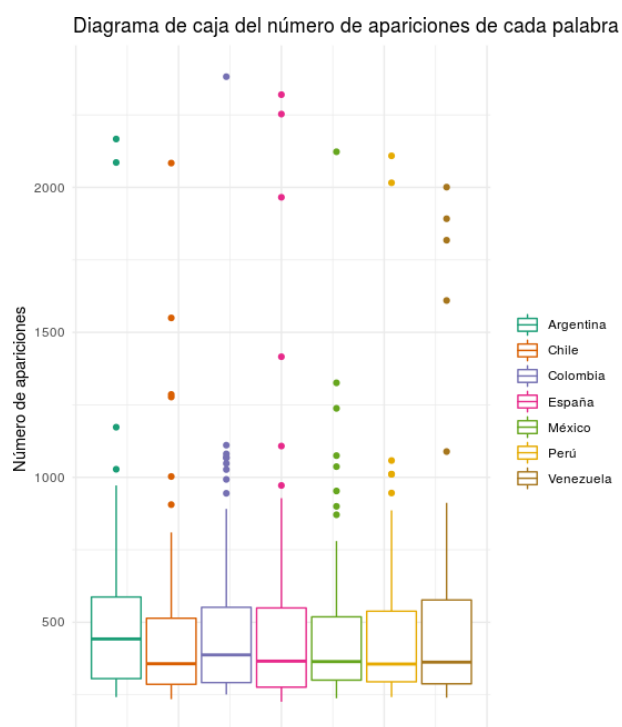


Figura 4: Diagrama de cajas por variedad lingüística para las cien palabras más frecuentes.