

escoger que parámetro se utilizará para compararlos y, en este estudio, debido a la predominancia de la secuencia genética BRCA frente al resto, apareciendo en más del doble de ocasiones que la segunda secuencia genética más común, se ha tomado la decisión de utilizar el valor Kappa para evitar cualquier desbalance en los resultados.

Para entrar más en detalle los paradigmas se separarán en dos grupos: los que tienen un porcentaje de Kappa inferior a 0.95 y los que están por encima de ese valor.

Dentro del primer grupo se encuentran los árboles de decisión, las redes neuronales de más de una capa y la SVM radial. Este subgrupo también sería susceptible a ser dividido en dos ya que los resultados obtenidos de los árboles de decisión y de la red de dos capas son bastante buenos, entre el 0.80 y el 0.85, y distan mucho de los pésimos valores ofrecidos por la red neuronal de 3 capas, 0.45, y de la SVM, 0. Estos dos últimos los desecharemos desde el primer momento. En el caso del paradigma C4.5 no se observa una gran variabilidad en los datos si cambiamos los parámetros, por lo general son bastante buenos, lo cual no ocurre con RPart y RPart2 que tienden a sobreajustar mucho los resultados y cuando se utilizan parámetros más reales cae bastante su exactitud. En el caso de la red neuronal de 2 capas vemos como a medida que aumenta el número de neuronas van mejorando sus resultados, pero siguen estando lejos tanto en exactitud como en tiempo si los comparamos con la red de 1 capa.

Si pasamos al segundo subgrupo nos encontramos con resultados muy interesantes, rodando un Kappa de 0.973 en el caso de la red neuronal de 1 capa, y las SVM lineales y polinomial, un poco menos exacto es el Knn con un Kappa de 0.959. Con estos valores de Knn podemos asegurar que no hay variables inservibles, ya que sino Knn saldría perjudicado. Estudiando los resultados de la red neuronal podemos ver como tiende a mejorar a mayor cantidad de neuronas por capa, pero también aumenta el tiempo empleado en calcularlo. Por último, las maquinas de vector soporte han dado muy buenos resultados tanto en el conjunto de datos de entrenamiento como en el de prueba; también podemos asegurar que las regiones de los datos están bien alejadas unas de otras ya que el rendimiento ofrecido por el margen duro y el margen blando no ha influido en los resultados.

4. Conclusión

Por ultimo voy a tratar de explicar porque he escogido el paradigma Knn como el mejor para este problema y como he llegado a esta conclusión.

Para empezar es necesario decir que todos los paradigmas que entraban en el subgrupo de un Kappa inferior a 0.95 también cumple que, si comparamos sus resultados mediante binom-test frente a los resultados del Knn, nos sale siempre un valor inferior al 0.05, por lo que se puede asegurar que Knn es mejor. También se han comparado su desempeño en el entrenamiento mediante t-test frente al Knn y se puede asegurar que Knn es mejor que todos ellos.

Sin embargo estas decisiones no pueden ser aplicadas al grupo que obtuvo unos resultados de Kappa superiores a 0.95, ya que de entre ellos ninguno puede

asegurarse que es mejor que los demás mediante estos test; de hecho los resultados mostrados por Knn son ligeramente peores. Tampoco se puede argumentar que los valores varían mucho entre entrenamiento y pruebas, ya que están bastante parejos. Por estas razones cualquiera de los paradigmas de este subgrupo serian totalmente validos.

Como ultimo pilar para apoyar la elección de Knn, se ha utilizado el tiempo empleado en entrenar el paradigma. Si se observa la **Figura ??** se puede observar que solo aparecen 4 paradigmas y esto es porque la red neuronal ha sido descartada desde un principio ya que desvirtuó el resto de resultados, por lo que no se puede comparar en cuanto al tiempo empleado con los demás. De los 4 restantes, el polinomial es el peor, seguido de las dos versiones del lineal y el mejor es el Knn con un tiempo estimado de 0 segundos. Esto quiere decir que, con la muestra aplicada al estudio, el tiempo empleado por el Knn es despreciable. Todos los tiempos aumentarán a medida que crezca el número de variables, he incluso Knn llegará a necesitar cierto tiempo de computación, pero este valor será muy inferior a los demás si tenemos en cuenta el conjunto de datos inicial, de 20.000 variables, donde el Knn sería el más optimo en cuanto a resultados frente a tiempo.

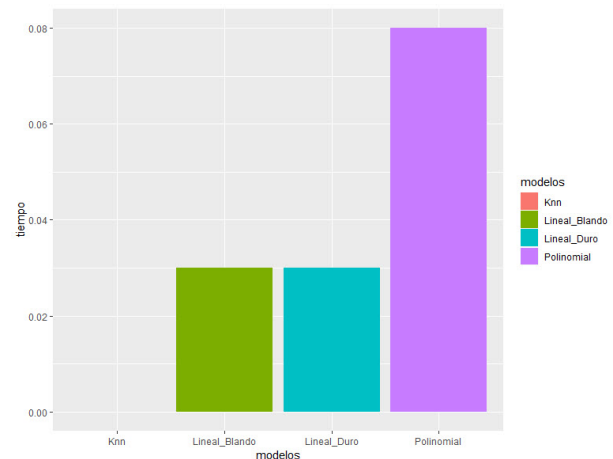


Figura 2: Tiempos empleados en el entrenamiento de los mejores modelos.