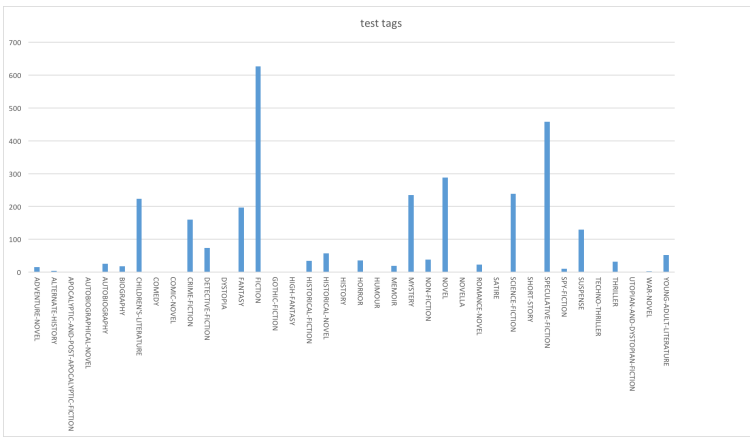


預測 test 出來的 tags 分佈情況



觀察上述兩張圖片可以發現其分佈非常像 (tag 在 train 中出現較多次則在 test 中也會出現較多次)。

4. (1%) 本次作業中使用何種方式得到 word embedding? 請簡單描述做法。

答：

透過至 Stanford 大學網頁下載 Stanford 大學 NLP 組提供的 GloVe (Global Vectors for Word Representation)

GloVe 是透過整合基於統計 word vector 模型的方法以及基於預測 word vector 模型的非監督式學習方法，其 GloVe 的損失函數為：

$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

其中 f 為函數用來降低高頻詞對模型的干擾。

上述方法為 Stanford 如何訓練出 glove.6B.XX.txt (XX 代表維度, ex: 50, 100, 200, 300) 的簡易介紹，要使用此 txt 很簡單。因為此 txt 使用空白分隔，故透過 python split 即可將一行切成 array, array[0] 代表 word, array[1 - XX] 則代表 word vector，即可使用。

5. (1%) 試比較 bag of word 和 RNN 何者在本次作業中效果較好。

答：

在實作 bag of word 是紀錄 input 的每個 word 出現次數，並轉成一個維度很大的 vector，作為 input vector，再使用 fully connected neural network 的 model 下去訓練，最終透過此方法訓練出來的 model 丟到 kaggle 結果為 0.46121。而透過 RNN (第 1 題所使用的架構方法) 訓練出來的 model 丟到 kaggle 結果為 0.50343。故我在本次作業中使用 RNN 得到的效果會比較。