



## DATA SCIENCE: ONLINE NEWS POPULARITY

**Abstract:** This dataset summarizes a heterogeneous set of features about articles published by Mashable in a period of two years (Mashable's primary focus is on technology, lifestyle and entertainment news. It is a global, multi-platform media and entertainment company).

**Internship Task:** Please carry out an Exploratory Data Analysis and create compelling story based on the given dataset; also predict the number of shares in social networks (popularity).

### Attribute Information:

1. url: URL of the article (non-predictive)
2. n\_tokens\_title: Number of words in the title
3. n\_tokens\_content: Number of words in the content
4. n\_unique\_tokens: Rate of unique words in the content
5. n\_non\_stop\_words: Rate of non-stop words in the content
6. n\_non\_stop\_unique\_tokens: Rate of unique non-stop words in the content
7. num\_hrefs: Number of links
8. num\_self\_hrefs: Number of links to other articles published by Mashable
9. num\_imgs: Number of images
10. num\_videos: Number of videos
11. average\_token\_length: Average length of the words in the content
12. num\_keywords: Number of keywords in the metadata
13. data\_channel\_is\_lifestyle: Is data channel 'Lifestyle'?
14. data\_channel\_is\_entertainment: Is data channel 'Entertainment'?
15. data\_channel\_is\_bus: Is data channel 'Business'?
16. data\_channel\_is\_socmed: Is data channel 'Social Media'?
17. data\_channel\_is\_tech: Is data channel 'Tech'?
18. data\_channel\_is\_world: Is data channel 'World'?
19. weekday\_is\_monday: Was the article published on a Monday?
20. weekday\_is\_tuesday: Was the article published on a Tuesday?
21. weekday\_is\_wednesday: Was the article published on a Wednesday?
22. weekday\_is\_thursday: Was the article published on a Thursday?

23. weekday\_is\_friday: Was the article published on a Friday?
24. weekday\_is\_saturday: Was the article published on a Saturday?
25. weekday\_is\_sunday: Was the article published on a Sunday?
26. is\_weekend: Was the article published on the weekend?
27. LDA\_00: Closeness to LDA topic 0
28. LDA\_01: Closeness to LDA topic 1
29. LDA\_02: Closeness to LDA topic 2
30. LDA\_03: Closeness to LDA topic 3
31. LDA\_04: Closeness to LDA topic 4
32. global\_sentiment\_polarity: Text sentiment polarity
33. global\_rate\_positive\_words: Rate of positive words in the content
34. global\_rate\_negative\_words: Rate of negative words in the content
35. avg\_positive\_polarity: Avg. polarity of positive words
36. avg\_negative\_polarity: Avg. polarity of negative words
37. title\_sentiment\_polarity: Title polarity
38. shares: Number of shares (TARGET)

Citation: K. Fernandes, P. Vinagre and P. Cortez. A Proactive Intelligent Decision Support System for Predicting the Popularity of Online News. Proceedings of the 17th EPIA 2015 - Portuguese Conference on Artificial Intelligence, September, Coimbra, Portugal. <http://archive.ics.uci.edu/ml/datasets/Online+News+Popularity>