

Topic Modeling

Project Overview

- Goal: Discover latent topics in unlabeled documents
- Unsupervised NLP with probabilistic modeling
- Outputs: Topics (top words) + document-topic distribution

Steps

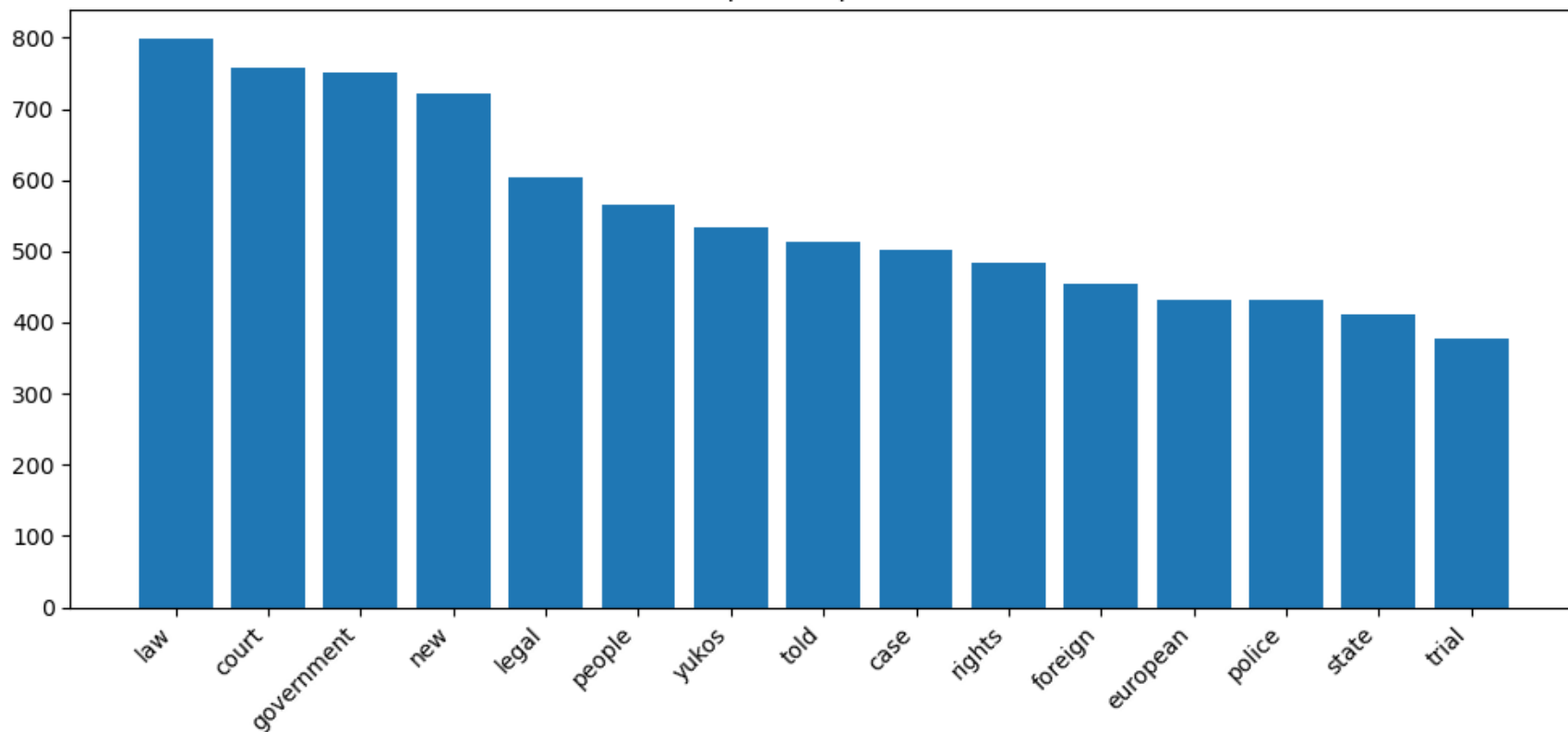
- Clean text (lowercase, remove punctuation/stopwords)
- Vectorize (Count, TF-IDF)
- Fit LDA to discover topics
- Evaluate with coherence & perplexity (if shown)
- Visualize top words + topic distribution

Method

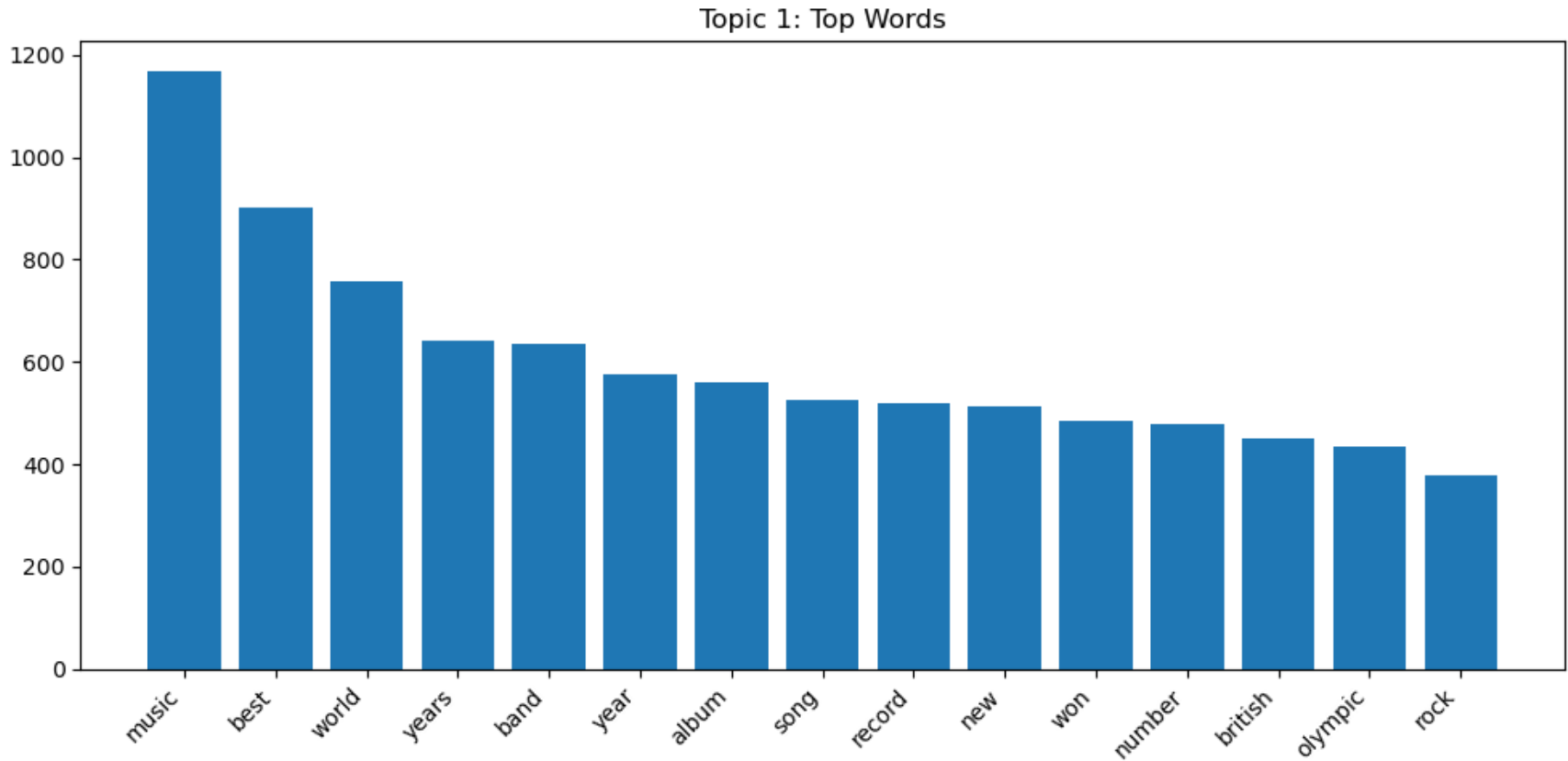
- Model: Latent Dirichlet Allocation (LDA)
- Features: CountVectorizer and TF-IDF
- Preprocessing: lowercasing,
digits/punctuation removal, stopwords

Visualization

Topic 0: Top Words

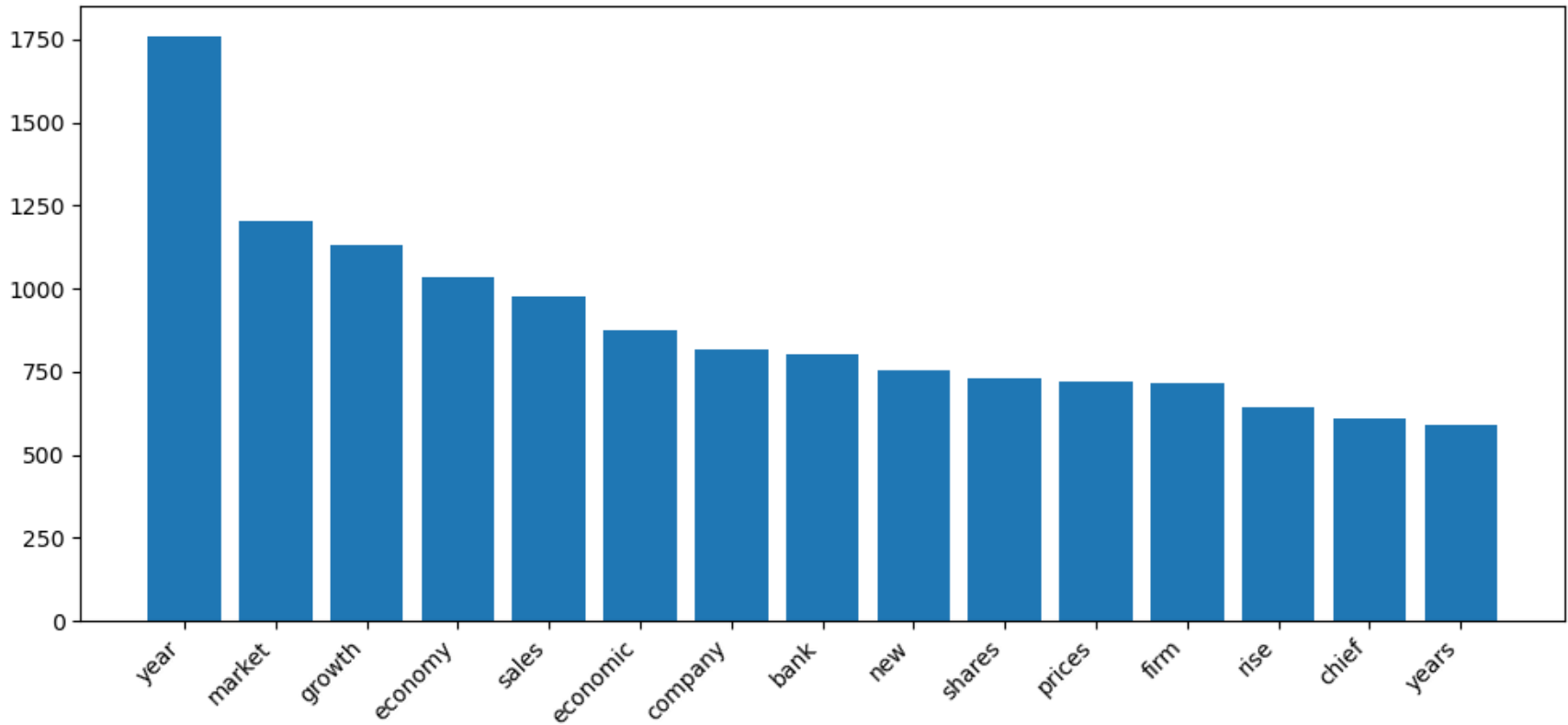


Visualization 2



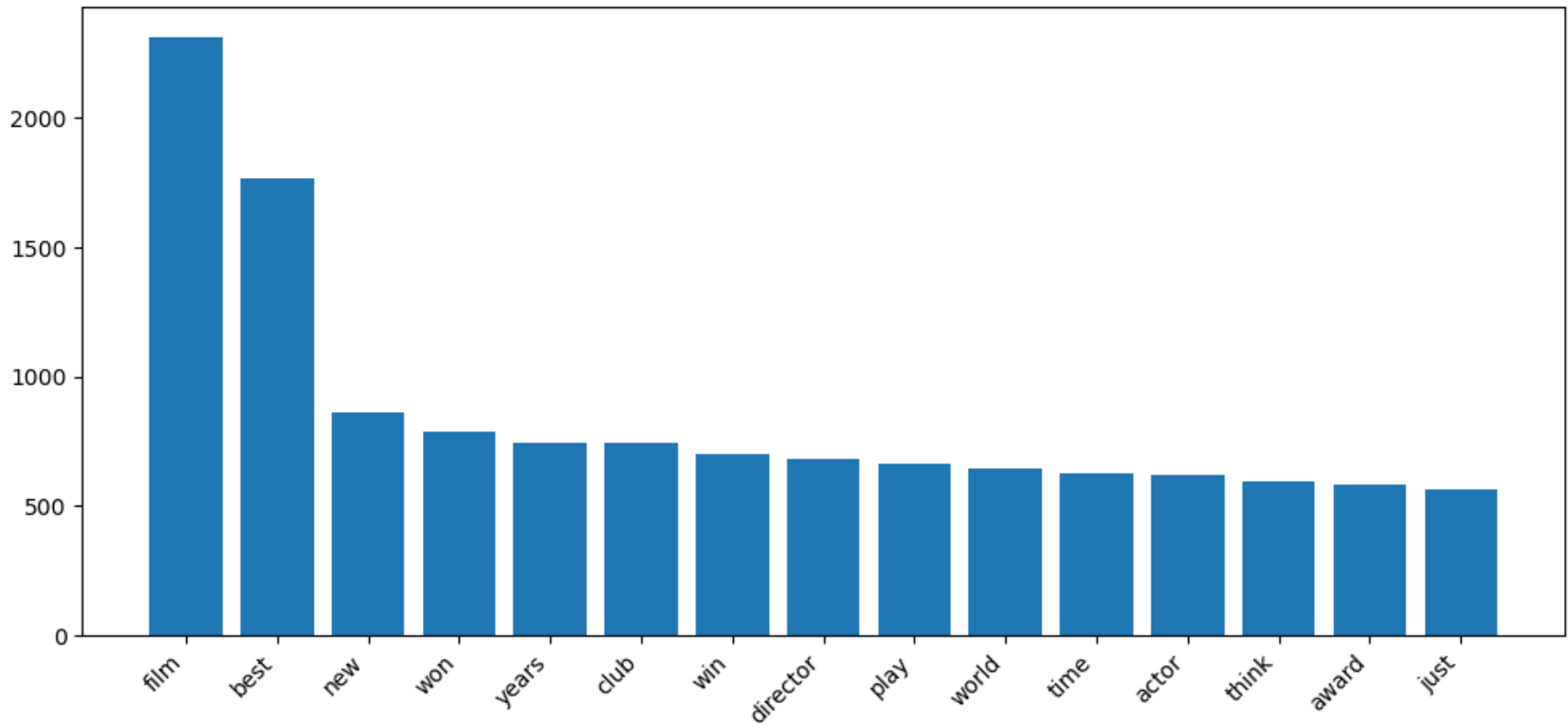
Visualization 3

Topic 2: Top Words



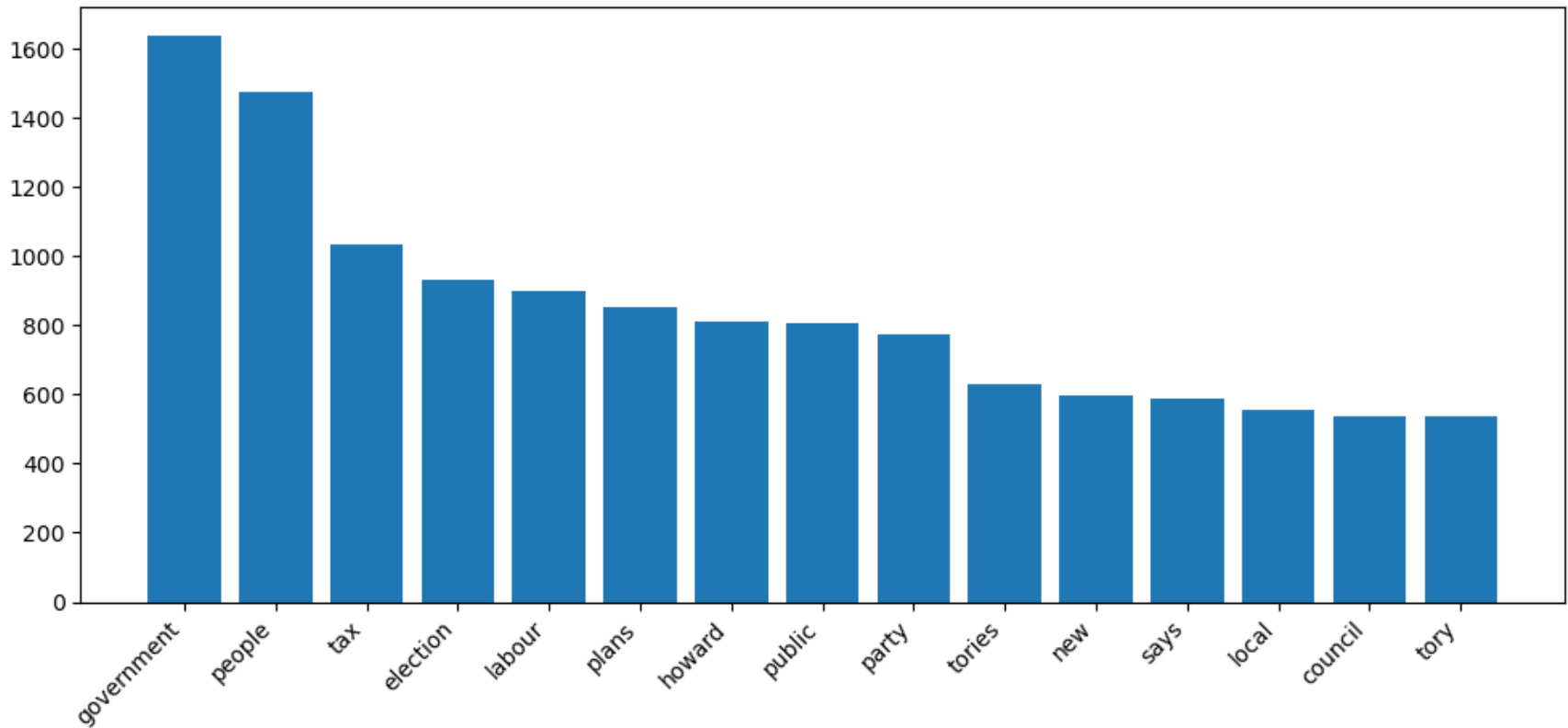
Visualization 4

Topic 3: Top Words



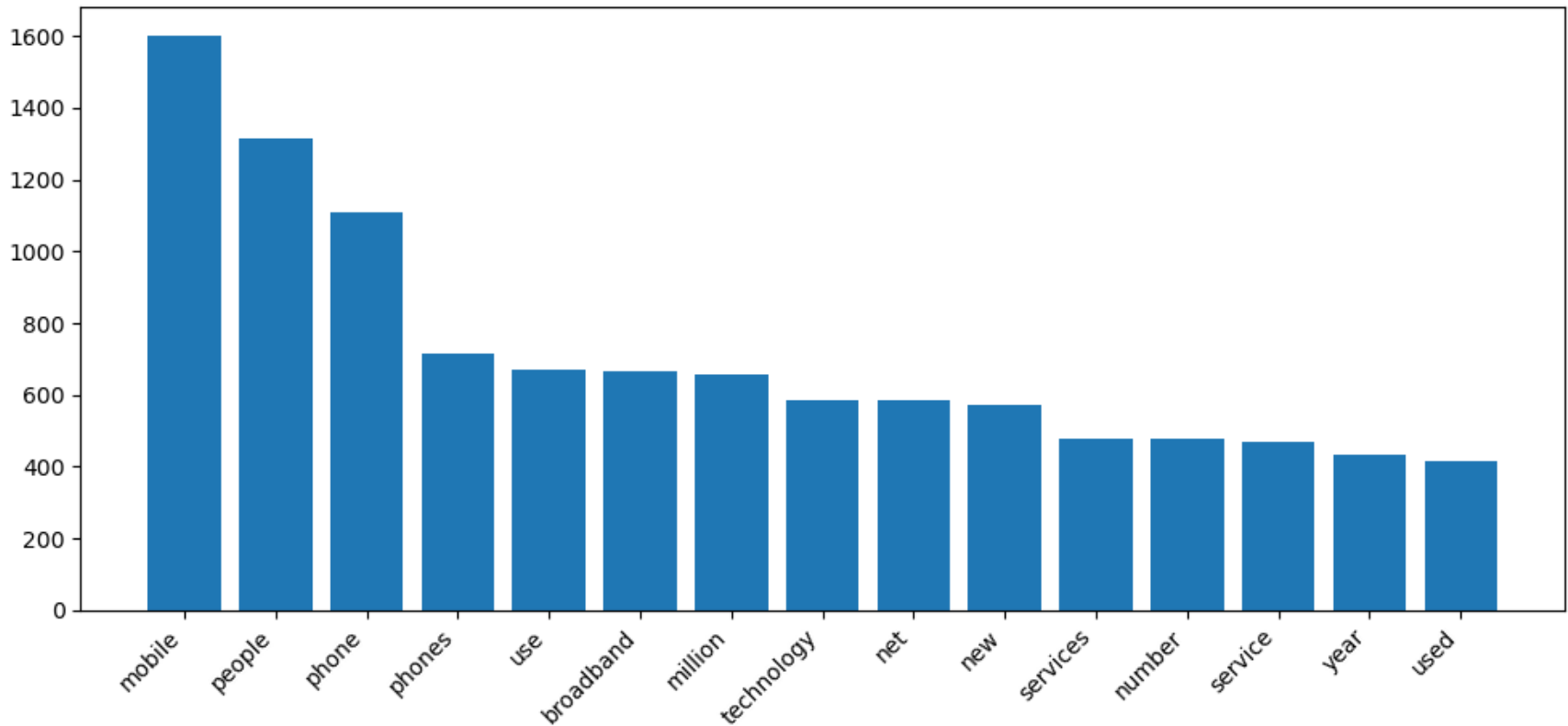
Visualization 5

Topic 4: Top Words

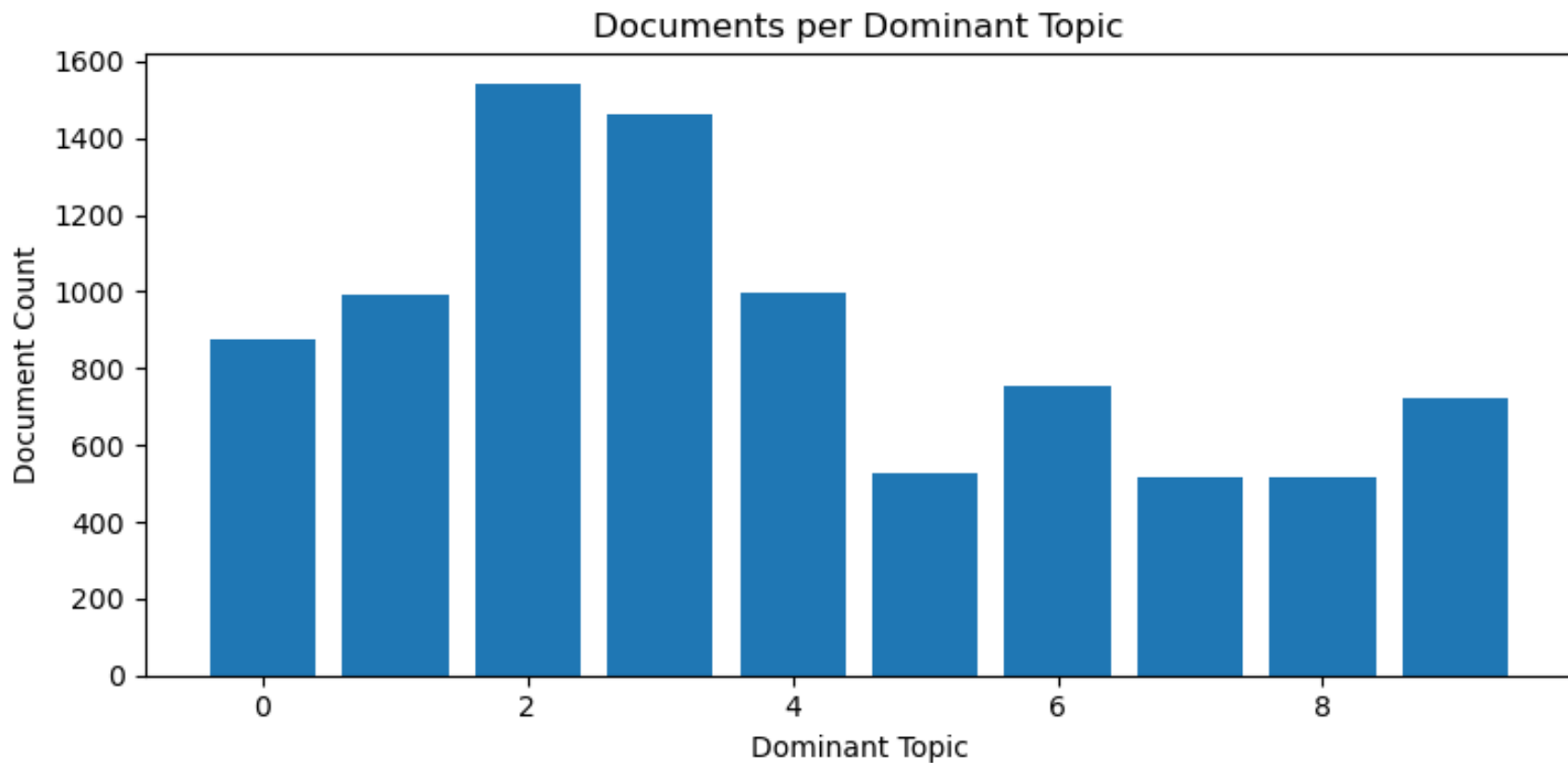


Visualization 6

Topic 5: Top Words



Visualization 7



Conclusion

- LDA reveals coherent topics via top words
- Figures reused directly from the notebook
- All numbers and visuals come from this HTML only