



Readability evaluation with cognitive features

Readability and language variations in a historical context

The study of readability investigates how easily a text can be read. Approaches from cognitivism and distributional semantics, investigating the relations between words, demonstrate how a good measure of cohesion in the text makes its comprehension easier.

- Our project takes into account the presence of **multiword expressions (MWEs)** in texts while most of the approaches deal only with single words.
- People consider as “individual expressions” groups of up to four words, **frequency** of MWEs affects the **processing performance**.

Since well known MWEs are processed faster, the presence of specific collocations or groups of words can be a relevant feature for measuring readability. This can lead to readability measures closer to how the human mind really works while reading.

Reading comprehension is not only related to the characteristics of a text/speech.

- Comprehension takes into account also **cultural background** (intended as other written works).
- A text can be easily understood in some years because it contains concepts close to the common way of thinking, while it could be difficult to read in other years because of the different mindset.

We are collaborating with I.S.I.G. (Italian-German Historical Institute), the history research group at FBK. The goal of the joint project is to analyze the political discourses of Alcide De Gasperi (first Prime Minister of Italy) with text processing tools.

We are working at developing a novel approach to exploit Google n-grams, in order to understand to what extent the language and the relations between concepts in De Gasperi’s corpus are aligned with the language and the concepts appeared in the years in which he wrote. This is useful also in order to track the evolution of a concept over the time.

A.L.C.I.D.E.: Analysis of Language and Context in a Digital Environment

The screenshot displays the 'Historical Corpus Analysis Portal' interface. It features a 'Documents Distribution' line graph showing frequency over time (1901-1954). Below the graph, there's a 'Test Classification' section with a 'Keywords' list. The keywords include terms like 'autonomia amministrativa', 'autonomia regionale', 'autonomia provinciale', and 'libertà'. The interface also shows a snippet of text from a document dated 1920-01-17.

Autonomia 1919 - 1922

This section compares n-grams from Alcide De Gasperi's corpus with Google Ngram Viewer. On the left, under 'Alcide De Gasperi', a list of n-grams is shown: 'autonomia amministrativa', 'autonomia regionale', 'speciale autonomia', 'autonomia trentina', 'progetto di autonomia', 'autonomia della provincia', 'nostra autonomia culturale', 'ampia autonomia provinciale', 'bisogno di autonomia', 'parliamo di autonomia', 'ordinamenti di autonomia', 'rispetto reciproco', 'aprile 1920', 'decentramento amministrativo', and 'problema'. On the right, under 'Google Ngram Viewer', a list of n-grams is shown: 'autonomia amministrativa', 'completa autonomia', 'piena autonomia', 'autonomia politica', 'larga autonomia amministrativa', 'autonomia economica', 'autonomia nazionale', 'autonomia locale', 'indipendenza regionale', 'speciale autonomia', 'libertà', 'necessaria autonomia', 'autogoverno', 'Trentino', 'riconoscimento', 'rivendicazione', and 'necessità'.

We extracted from De Gasperi’s corpus all the n-grams and compared them to Google n-grams after ranking both lists according to several measures. This allows us to compare the co-occurrences of a given target word (or multiword expression).