

International Hellenic University

Natural Language Processing and Text Mining

Academic Year 2022-2023

Assignment: "Relevance Prediction"

A. Introduction. The area of text analytics (or text mining) includes techniques from multitude scientific areas (A.I., statistics, linguistics), and it has a wide range of applications (security, marketing, information retrieval, opinion mining). While structured data are "ready" for analysis and evaluation, unstructured text requires transformation in order to uncover the underlying information. The transformation of unstructured text into a structured set of data is not a straight forward task and text analytics offers a wide variety of tools to tackle with the idioms, ambiguities and irregularities of natural language.

In this assignment, you are going to work in teams (2 persons per team). The task that must be performed is to predict the relevance of a result with respect to a query. Each query is represented by a search term or terms and the result is a specific product. Based on the relevance between the query and the answer, a relevance score is assigned. The higher the score the better the relevance of the answer. The relevance is a real number in the range [1,3] (in steps of 0.5), where 1 denotes minimum relevance and 3 is the maximum relevance. Your job is to be able to predict the relevance score for an unknown combination between a query and a result. The evaluation measure used is the **Root-Mean-Square Error (RMSE)**.

B. Datasets. There are three datasets available. The first one, **train.csv**, contains pairs of queries and answers and also contains the relevance score. A screenshot of the dataset is shown in Figure 1. As you can observe, a product may be the best answer for different queries (search terms).

id	product_uid	product_title	search_term	relevance
2	100001	Simpson Strong-Tie 12-Gauge Angle	angle bracket	3
3	100001	Simpson Strong-Tie 12-Gauge Angle	l bracket	2.5
9	100002	BEHR Premium Textured DeckOver 1-gal. #SC-141 Tugboat Wood and Concrete Coating	deck over	3

Figure 1. Screenshot of **train.csv**.

The second dataset is **product_descriptions.csv**, and for each different product it contains a textual description of the product. A screenshot is shown in Figure 2.

```
[{"product uid": "100001", "product description": "Not only do angles make joints stronger, they also provide more consistent, straight corners. Simpson Strong-Tie offers a wide variety of angles in various sizes and thicknesses to handle light-duty jobs or projects where a structural connection is needed. Some can be bent (skewed) to match the project. For outdoor projects or those where moisture is present, use our ZMAX zinc-coated connectors, which provide extra resistance against corrosion (look for a \"Z\" at the end of the model number).Versatile connector for various 90 connections and home repair projectsStronger than angled nailing or screw fastening aloneHelp ensure joints are consistently straight and strongDimensions: 3 in. x 3 in. x 1-1/2 in.Made from 12-Gauge steelGalvanized for extra corrosion resistanceInstall with 10d common nails or #9 x 1-1/2 in. Strong-Drive SD screws"}]
```

Figure 2. Screenshot of **product_descriptions.csv**.

```
[{"product uid": "100001", "name": "Bullet01", "value": "Versatile connector for various 90 connections and home repair projects"}, {"product uid": "100001", "name": "Bullet02", "value": "Stronger than angled nailing or screw fastening alone"}, {"product uid": "100001", "name": "Bullet03", "value": "Help ensure joints are consistently straight and strong"}, {"product uid": "100001", "name": "Bullet04", "value": "Dimensions: 3 in. x 3 in. x 1-1/2 in."}, {"product uid": "100001", "name": "Bullet05", "value": "Made from 12-Gauge steel"}, {"product uid": "100001", "name": "Bullet06", "value": "Galvanized for extra corrosion resistance"}, {"product uid": "100001", "name": "Bullet07", "value": "Install with 10d common nails or #9 x 1-1/2 in. Strong-Drive SD screws"}, {"product uid": "100001", "name": "Gauge", "value": "12"}, {"product uid": "100001", "name": "Material", "value": "Galvanized Steel"}, {"product uid": "100001", "name": "MFG Brand Name", "value": "Simpson Strong-Tie"}, {"product uid": "100001", "name": "Number of Pieces", "value": "1"}, {"product uid": "100001", "name": "Product Depth (in.)", "value": "1.5"}, {"product uid": "100001", "name": "Product Height (in.)", "value": "3"}, {"product uid": "100001", "name": "Product Weight (lb.)", "value": "0.26"}, {"product uid": "100001", "name": "Product Width (in.)", "value": "3"}, {"product uid": "100002", "name": "Application Method", "value": "Brush,Roller,Spray"}, {"product uid": "100002", "name": "Assembled Depth (in.)", "value": "6.63 in"}, {"product uid": "100002", "name": "Assembled Height (in.)", "value": "7.76 in"}, {"product uid": "100002", "name": "Assembled Width (in.)", "value": "6.63 in"}]
```

Figure 3. Screenshot of **attributes.csv**.

The third file, **attributes.csv**, contains additional attributes for some of the products (some products may not contain additional attributes). A screenshot is shown in Figure 3.

C. Deliverables. You are requested to write code in any programming language (preferably in Python) in order to solve the relevance prediction problem. You are free to use supervised or unsupervised learning techniques, use feature extraction techniques for text, apply text preprocessing, and in general you are free to apply any technique aiming at the **minimization of RMSE**. In addition to the source code, you must prepare a technical report, no more than 5 pages, explaining your solution and the techniques you have tested, sets of features that have shown the best performance, feature engineering, etc. You should upload your solutions using the corresponding assignment link in the IHU elearning platform. The deadline for the submission is set to **June 18, 2023**. The assignment grade will contribute to **30%** of your total course grade.

Try to get the most out of this assignment, by applying techniques that you have learned or other methods that you may find in books and related material. Also, have fun with the project!

The NLP&TM Team