

---

# example

---

*example description*

*Author: Anton Zhitomirsky*

**Contents**

<b>1</b>	<b>TF-IDF</b>	<b>2</b>
1.1	Motivating Problem . . . . .	2
1.2	Solution . . . . .	2
1.3	Introduction . . . . .	2

# 1 TF-IDF

## 1.1 Motivating Problem

If we have a search in Google for “low carb breakfast” they get shown generic results for breakfast. Some of these terms shouldn’t have equal weighting to each-other.

The problem is that the search is using basic string/keyword matching. It treats all words as equally important. Less relevant results are ranked as highly as more relevant ones.

## 1.2 Solution

Rank more important words as more important.

## 1.3 Introduction

This is pre-neural networks.

- **Term Frequency (TF):** Measures how often a term occurs in a document. The more often a term appears in a document, the more important it is for that document.
- **Inverse Document Frequency (IDF):** Measures how common or rare a term is across all documents in the corpus. Terms that appear in many different documents are less significant than those that appear in a smaller number of documents.

This means that a search for “low-carb breakfast” will prioritize recipes where “low-carb” is a significant term, rather than returning recipes with the more common “breakfast” term.

## TF-IDF

### Term Frequency (TF):

Upweights words  $w$  that are more important to  $d$

Frequency of  $w$  occurring together with  $d$

$$TF_{w,d} = \frac{\text{count}(w, d)}{\sum_{w'} \text{count}(w', d)}$$

### Inverse Document Frequency (IDF):

Downweights words that appear everywhere

Size of document collection

$$IDF_{w,D} = \log \frac{|D|}{|\{d \in D : w \in d\}|}$$

Number of documents in  $D$  that contain  $w$

$$TF\text{-}IDF_{w,d,D} = TF_{w,d} IDF_{w,D}$$

7