

Introduction to Machine Learning

NPFL 054

<http://ufal.mff.cuni.cz/course/npfl054>

Barbora Hladká

Martin Holub

{Hladka | Holub}@ufal.mff.cuni.cz

Charles University,
Faculty of Mathematics and Physics,
Institute of Formal and Applied Linguistics

Term Frequency-Inverse Document Frequency

- How important a word is to a document D in a collection C ($|C| = N$)?
- term frequency
 $tf_{t,D}$ = the number of times a term t occurs in D
- document frequency
 $df_{t,D}$ = the number of documents in C in which a term t occurs, i.e.,
 $|\{D \in C : t \in D\}|$
- inverse document frequency, term importance
 $idf_{t,D} = \log \frac{N}{df_{t,D}}$

$$tfidf_{t,D,C} = tf_{t,D} \cdot idf_{t,C}$$