# Introduction to Machine Learning
## NPFL 054

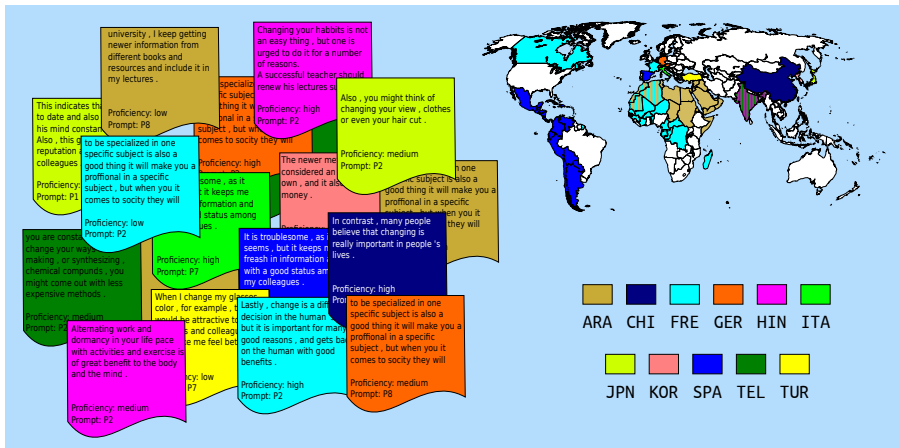http://ufal.mff.cuni.cz/course/npfl054

**Barbora Hladká**                     **Martin Holub**

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

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

# NLI

Identifying the native language (L1) of a writer based on a sample of their writing in a second language (L2)

**Our data**

- **L1s**: Arabic (ARA), Chinese (ZHO), French(FRA), German (DEU) Hindi (HIN), Italian (ITA), Japanese (JPN), Korean (KOR), Spanish (SPA), Telugu (TEL), Turkish (TUR)
- **L2**: English
- **Real-world objects**: For each L1, 1,000 texts in L2 from The ETS Corpus of Non-Native Written English (former TOEFL11), i.e. *Train ∪ DevTest*
- **Target class:** L1

*More detailed info is available at the course website.*

# NLI

**Topic**
Most advertisements make products seem much better than they really are

**Sample text**
now a days the publisity is the best way to promoved a produt and if you wanth to sale a product you should bring some information that makes , that the people who is seeing the advertisements make sure that the product very good and in the future this person could buy it .

**L1 = Spanish**

# tf-idf

**T**erm **F**requency-**I**nverse **D**ocument **F**requency

- How important a word is to a document $D$ in a collection $C$ ($|C| = N$)?
- term frequency
  $\mathrm{tf}_{t,D}$ = the number of times a term $t$ occurs in $D$ (other possibilities exist)
- document frequency
  $\mathrm{df}_{t,D}$ = the number of documents in $C$ in which a term $t$ occurs, i.e.,
  $|\{D \in C : t \in D\}|$
- inverse document frequency
  $\mathrm{idf}_{t,D} = \log N/\mathrm{df}_{t,D}$

$$\mathrm{tfidf}_{t,D,C} = \mathrm{tf}_{t,D} \cdot \mathrm{idf}_{t,C}$$

Other variants of $\mathrm{tf}_{t,D}$ and $\mathrm{idf}_{t,D}$ exist.