# Introduction to Machine Learning, Fall 2014 - Exercise session II

Rodion "rodde" Efremov

November 3, 2014

## Problem 1 (3 points)

Consider a document-term matrix, where $tf_{ij}$ is the number of times that the $i^{th}$ word (term) appeares in the $j^{th}$ document, and let $m$ be the total number of documents in the collection. Consider the variable transformation that is defined by

$$tf'_{ij} = tf_{ij} \log \frac{m}{df_i}, \tag{1}$$

where $df_i$ is the number of documents in which the $i^{th}$ term appears, which is known as the *document frequency* of the term. This transformation is known as the *inverse document frequency* transformation.

  (a) What is the effect of this transformation if a term occurs in only one document? In every document?

  (b) What is the overall effect and what might be the purpose of this transformation?

  (c) Can you think of other (non-document) data in which this transformation might be useful?

### (a)

If the the term occurs in only one document, we have $df_i = 1$, and, thus, $tf'_{ij} = tf_{ij} \log m$. If $df_i = m$, we have that $tf'_{ij} = tf_{ij} \log \frac{m}{m} = 0$.

### (b)

The inverse document frequency aims to minimize the effect of over-emphasized terms such as "the", "in", and so on, and to emphasize the words that are not used very often, and, thus, allow better chance of differentiating between documents.

# Problem 3 (3 points)

Proximity is typically defined between a pair of objects.

(a) Give two ways in which you might define the 'proximity' among a set of (more than two) objects (i.e. a single measure of how similar an arbitrary number of items are all to one another)

(b) How might you define the distance between two sets of points in Euclidian space?

(c) How might you define the proximity between two sets of data objects? (Make no assumptions about the data objects, except that a proximity measure is defined between any pair of objects.)

## (a)

If $P(\mathbf{x}, \mathbf{y})$ is a "normal" function giving the proximity of data objects $\mathbf{x}$ and $\mathbf{y}$, and $A$ is an arbitrary set of data objects (which can be indexed like $A_1, A_2, \ldots, A_n$), I would try

$$P(A) = \left( \sum_{1 \leq i < j \leq n} P(A_i, A_j) \right) \binom{n}{2}^{-1},$$

or namely the average proximity over all pairs of data objects from $A$. For another one, I would choose a small $\varepsilon > 0$, and define $P(A)$ as

$$\prod_{1 \leq i < j \leq n} \left( \max \left( P(A_i, A_j), \varepsilon \right) \right)^{\frac{1}{\binom{n}{2}}} \geq \varepsilon,$$

or namely the geometric mean over all pairs of data objects from $A$. The $\varepsilon$ is chosen to be positive so that a single proximity value of 0 does not dominate the entire proximity of a set and bring it down to 0. (Above, it is assumed that $P(\mathbf{x}, \mathbf{y}) \geq 0$ for all data objects $\mathbf{x}, \mathbf{y}$.)

## (b)

I would do the way topologists do: if $A$ and $B$ are two (finite) sets of points in Euclidian space, then

$$P(A, B) = \min_{(\mathbf{x}, \mathbf{y}) \in A \times B} d(\mathbf{x}, \mathbf{y}).$$

## (c)

Suppose we are given two sets of data objects, $A$ and $B$. Now, the easiest way to define proximity between them is

$$P(A, B) = \sum_{\mathbf{x} \in A} \sum_{\mathbf{y} \in B} \frac{P(\mathbf{x}, \mathbf{y})}{|A||B|}.$$