

Assignment 1: Linear Algebra and Probability

Submission: Tuesday August 16th

Maximum 2 students per group

Prof. Fabio A. González
Machine Learning - 2011-II
Maestría en Ing. de Sistemas y Computación

- Let $D = \{d_1, \dots, d_n\}$ be a set of documents and $T = \{t_1, \dots, t_m\}$ a set of terms (words). Let $TD = (TD_{i,j})_{i=1\dots m, j=1\dots n}$ be a matrix such that $TD_{i,j}$ corresponds to the number of times the term t_i appears in the document d_j . Assume a process where a document d_j is randomly chosen with uniform probability and then a term t_i , present in d_j , is randomly chosen with a probability proportional to the frequency of t_i in d_j .
 - How do you transform the matrix TD to obtain a matrix $P(T, D)$, such that $P(T, D)_{i,j} = P(t_i, d_j)$ (the joint probability of term t_i and document d_j).
 - How do you obtain the $P(T|D)$ matrix?
 - How do you obtain the $P(D|T)$ matrix?
- Let l_i be the length, number of characters, of term t_i , and let $L = (l_1, \dots, l_m)$ be a column vector.
 - Calculate $E[l]$, the expected value of the random variable l corresponding to the length of a randomly chosen term.
 - Calculate $\text{Var}(l)$, the variance of l .
 - Show that $\text{Var}(al + b) = a^2\text{Var}(l)$, where a and b are arbitrary constants.
- Find an expression for $P(T|D)$ that exclusively uses $P(D|T)$, $P(T, D)$ and, optionally, constant matrices/vectors.
- Find an expression for the matrix $COV = (\text{Cov}(t_i, t_j))_{i,j=1\dots m}$, the covariance of the random variables corresponding to terms.
- (optional) What is the meaning of the Eigenvector of COV corresponding to the largest Eigenvalue?

Note: In all the cases use standard matrix and scalar operations: transposition, matrix multiplication (*), matrix elementwise operations (+, -, .*, ./), matrix-scalar operations, etc.