# Project 2 *

February 22, 2016

# 1 Dataset and Problem Statement

## 1.1 Part A

# 2 Modeling Text Data and Feature Extraction

## 2.1 Part A

## 2.2 Part C

10 Most significant features with TFICF scores:

Class:comp.sys.ibm.pc.hardware
('scsi', 0.42809947405535098),
('drive', 0.31974715073739174),
('use', 0.22321146724973656),
('mb', 0.1968842917792481),
('ide', 0.18884045322355955),
('card', 0.15550544489414134),
('disk', 0.1477153861736914),
('control', 0.13410160556882411),
('dos', 0.12739571821010173),
('jumper', 0.11502197860943709)

Class:comp.sys.mac.hardware
('mac', 0.29656004455958551),
('use', 0.2251638723660454),
('scsi', 0.1890126678045923),
('appl', 0.18559985263916312),
('drive', 0.17315418917163491),
('mb', 0.17113312542860176),
('simm', 0.1637222404222877),
('problem', 0.14905214086203006),
('quadra', 0.13825360012263685),
('nubus', 0.12476544401311132)

---

Class:misc.forsale
('dos', 0.23204826904651923),
('new', 0.19885758974308329),
('sale', 0.18905428351306636),
('offer', 0.17954408282558959),
('use', 0.17596750747049816),
('includ', 0.17096030197337017),
('ship', 0.15768604272172954),
('price', 0.14163238406162046),
('wolverin', 0.1270563291946013),
('sell', 0.11874230178903535)

Class:soc.religion.christian
('god', 0.37385757805445796),
('christian', 0.27290603508238487),
('jesus', 0.23172842957647949),
('church', 0.20108044082283888),
('christ', 0.16786399940330063),
('peopl', 0.14539203709627244),
('say', 0.14475435272304318),
('bibl', 0.13177856565674298),
('believ', 0.12944992776554082),
('think', 0.12307308403324817)

# 3    Feature Selection

## 3.1    Part D

On applying LSI to the TFIDF matrix with k=50, each document was mapped to a 50 dimensional vector.

# 4    Learning Algorithms

## 4.1    Part E

# 5    Multi-class Classification

## 5.1    Part A