# Supervised Machine Learning Techniques to Classify Software Requirements

**Submitted by:**       **Faiza Latif**

**Registration No:**    **01-241161-005**

**Supervisor Name:**    **Dr.Muhammad Asfand-e-Yar**

**March,2018**

# Abstract

Requirements are the functionalities that are discovered before building any product. A systematic approach through which the software engineer collects requirements from diverse sources and implements them into the software development processes is called Requirement engineering. Requirements engineering contains a set of activities for discovering, analyzing, documenting, validating and maintaining a set of requirements for a system. Functional Requirements (FRs) and Non-Functional Requirements (NFRs) are two basic types of requirements in Software Requirement Specification documents. The classification of these requirements is an important task as it provides an ease for the team manager and the software development first. The NFRs grabs less attention from the development team. Some of the NFR categories are very important to consider while developing the software. This research study proposes a technique to classify the requirements in to FRs and NFRs with the help of Machine Learning techniques. The NFRs are defined in the requirement document but in some cases the NFRs are not clearly mentioned. Therefore, using Machine Learning the requirements are classified and system will automatically identify the categories of NFRss, which are evaluated using accuracy, precision and F-Measure. The accuracy explains that whether the NFR is accurately classified in the document, the precision explains whether the mentioned requirements are properly placed in the classified field. F-measure or F-score is the weighted average of precision and recall. Furthermore, it also classifies the NFRs into sub categories. Different ML approaches and classification algorithms will be used in the study.

# Contents