# Report about Statistical Learning Theory (SLT)

Binary classification is a fundamental problem in machine learning and statistics, where the goal is to categorize elements of a given dataset into one of two distinct classes.

Given a dataset $\{D\} = \{(x\_i, y\_i)\}_{\{i=1\}}^{n}$), where each ($x_i$ in $\{R\}^m$) is an m-dimensional feature vector and ($y_i$ in $\{-1, 1\}$) is a binary label, the task is to learn a function (f: $\{R\}^m \Rightarrow \{0, 1\}$) that accurately predicts the label (y) for new, unseen feature vectors (x).

The objective is to find a classifier (f(x)) that minimizes the expected prediction error on unseen data. This is often formalized using a loss function (L(y, f(x))), such as the zero-one loss. The goal is to minimize the empirical risk.

To solve the binary classification problem, we choose a hypothesis class, which is a set of candidate functions. Common choices include linear classifiers (e.g., logistic regression), decision trees, support vector machines (SVMs), and neural networks.

## How SLT offer math basic framework to solve the problem of binary classification in Machine Learning?

Statistical Learning Theory (SLT) provides a mathematical framework for addressing the problem of binary classification in machine learning by offering a structured approach to understanding, analyzing, and improving learning algorithms. Here are the key ways SLT contributes to solving binary classification problems:

1. Defining the Learning Problem. SLT starts by formalizing the learning problem. In binary classification, the goal is to find a function that maps input data to binary labels. This involves selecting a hypothesis from a hypothesis space that will generalize well to unseen data.

2. Empirical Risk Minimization (ERM). SLT introduces the concept of Empirical Risk Minimization as a strategy for learning. The empirical risk is the average loss over the training dataset, and ERM seeks to minimize this risk. By minimizing empirical risk, SLT provides a clear mathematical objective for training a classifier.

3. Generalization and Overfitting. A major concern in binary classification is overfitting, where a model performs well on training data but poorly on unseen data. SLT addresses this by providing theoretical bounds on the generalization error, which is the difference between empirical risk and true risk (expected error on new data).

4. Structural Risk Minimization (SRM). To mitigate overfitting, SLT proposes Structural Risk Minimization, which balances empirical risk with model complexity. SRM guides the selection of models that generalize better by minimizing both empirical risk and complexity.