# Machine Learning

# Offensive Tweet Detection

*November 16, 2018*

Diogo Daniel Ferreira    76504

Luís Davide Leira    76514

Departamento de Eletrónica, Telecomunicações e Informática
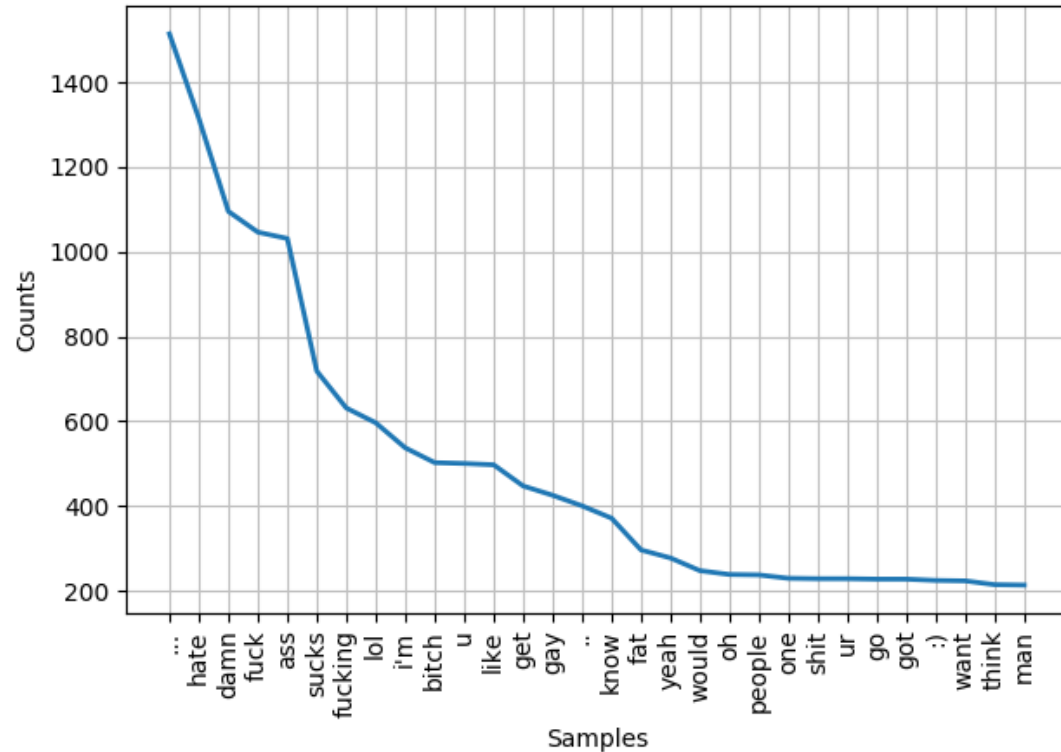
Universidade de Aveiro

# Problem

○ Improve the Twitter implementation of "Safe Search"

○ Don't show offensive tweets from users with "safe search" on

○ Also useful when large websites use Twitter widget to show tweets related with a topic

○ Twitter can suspend or delete accounts with successive violations of safe tweets
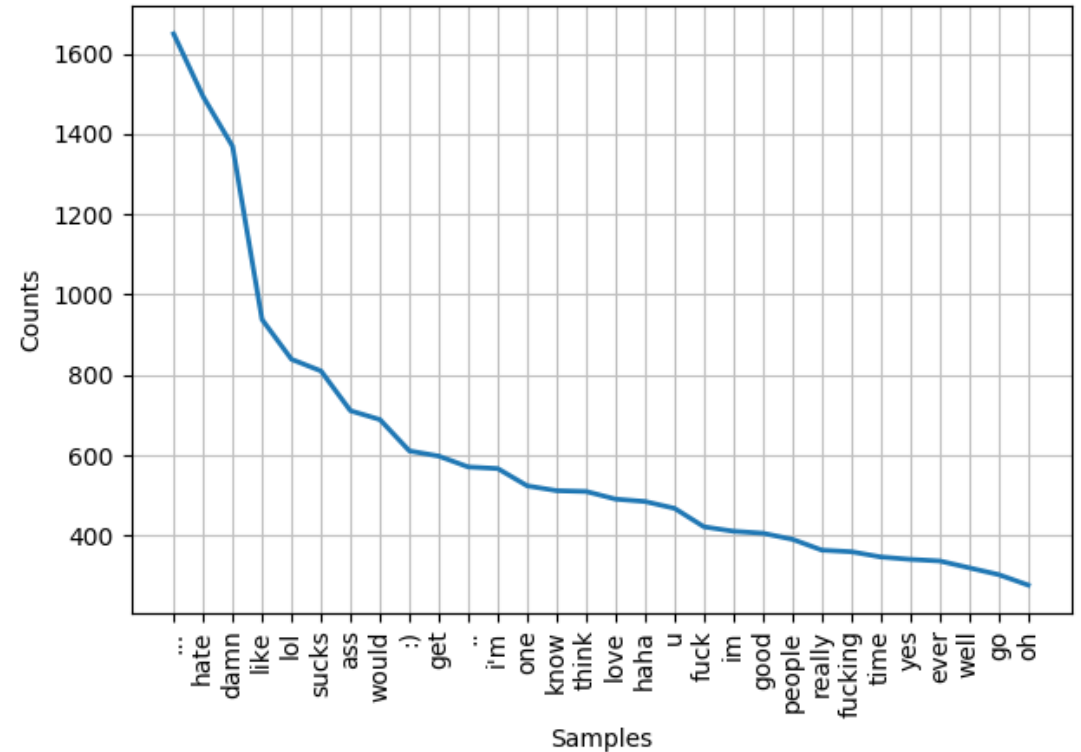
# Dataset

- 20001 tweets

- Human-labelled data in two categories:
  - 1 (offensive)
  - 0 (non offensive)

- 7822 labelled as offensive tweets and 12179 labelled as non offensive tweets

- Source: https://www.kaggle.com/dataturks/dataset-for-detection-of-cybertrolls/

# Dataset statistics

# Dataset statistics

# Data pre-processing

- Convert all words into lowercase

- Apply tokenizer

- Remove punctuation

- Remove stop words (the, and, or, …)

- Lemmatize

6

# Feature extraction

**Tokens counting**

○ Select the N most used tokens and transform them into features

○ Count the number of feature tokens that each tweet has

○ Normalize the counting

The pre-processed tokens can be grouped in a contiguous sequence: n-gram

We have chosen to treat each one as a word. Example:

○ 1-gram: ["We", "Love", "AI"]

○ 2-gram: ["We Love", "Love AI"]

# Feature extraction

**TF-IDF weighting**

- Select the N most used tokens and transform them into features
- Calculate the TF-IDF for the feature tokens that each tweet has
  - TF-IDF is a numerical statistic that reflects the importance of a token in a tweet
    - TF – term frequency
    - IDF – inverse document frequency

# Classifiers

Initially, 11 classifiers were put to test:

- Logistic Regression
- SGD Classifier
- Linear SVC
- NuSVC
- SVC
- Gaussian Naive Bayes

- Decision Tree
- Random Forest
- AdaBoost
- XGBoost
- Multi-layer Perceptrons

# Training & Testing approach

Dataset division: 80% train and cross-validation (5-fold cross-validation), 20% test

Parameters to be taken into consideration:

- Maximum number of tokens to create features
- Different N-gram parameters (for the tokens counting)
- Token counting and TF-IDF weighting combination
- Classifiers hyperparameters

Too many combinations!

# Training & Testing approach

**First phase:**

- Test all classifiers with the default hyperparameters
- Test with tokens counting, TF-IDF and both combined
- Test with the following number of most used tokens:
    - 100, 1000, 5000, 10000, 25000
- Test with 1-gram, 2-gram and both combined (only applied to tokens counting)
- Perform 5-fold validation and store the confusion matrix
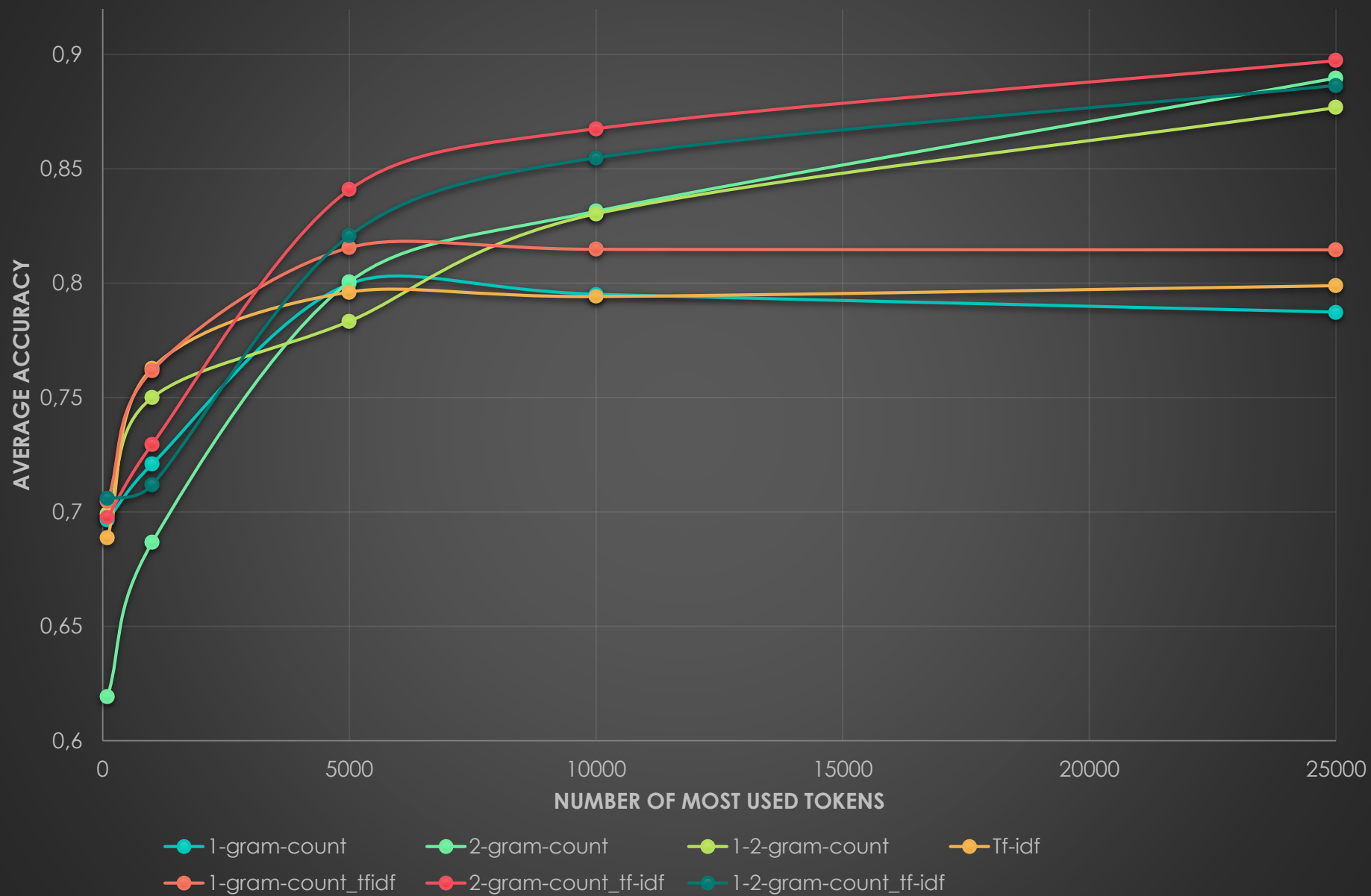
# Training & Testing approach

**Second phase:**

○ Use the four best classifiers from the first phase

○ Perform 5-fold cross validation on those four classifiers with different hyperparameters

○ Select the best hyperparameters for each one of them and get the confusion matrix on the test data
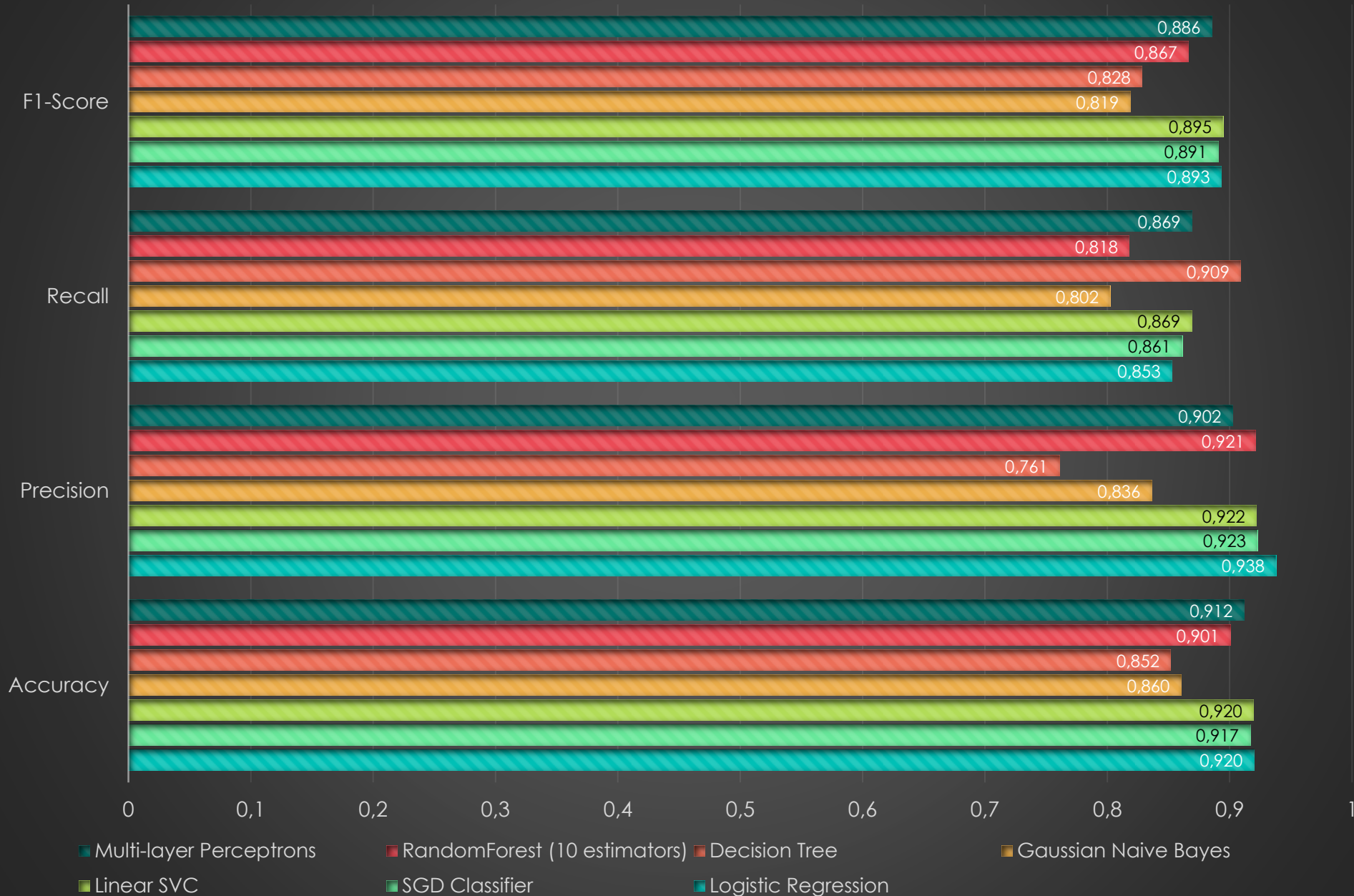
**Third phase:**

○ Test the four best classifiers with the best hyperparameters from the second phase

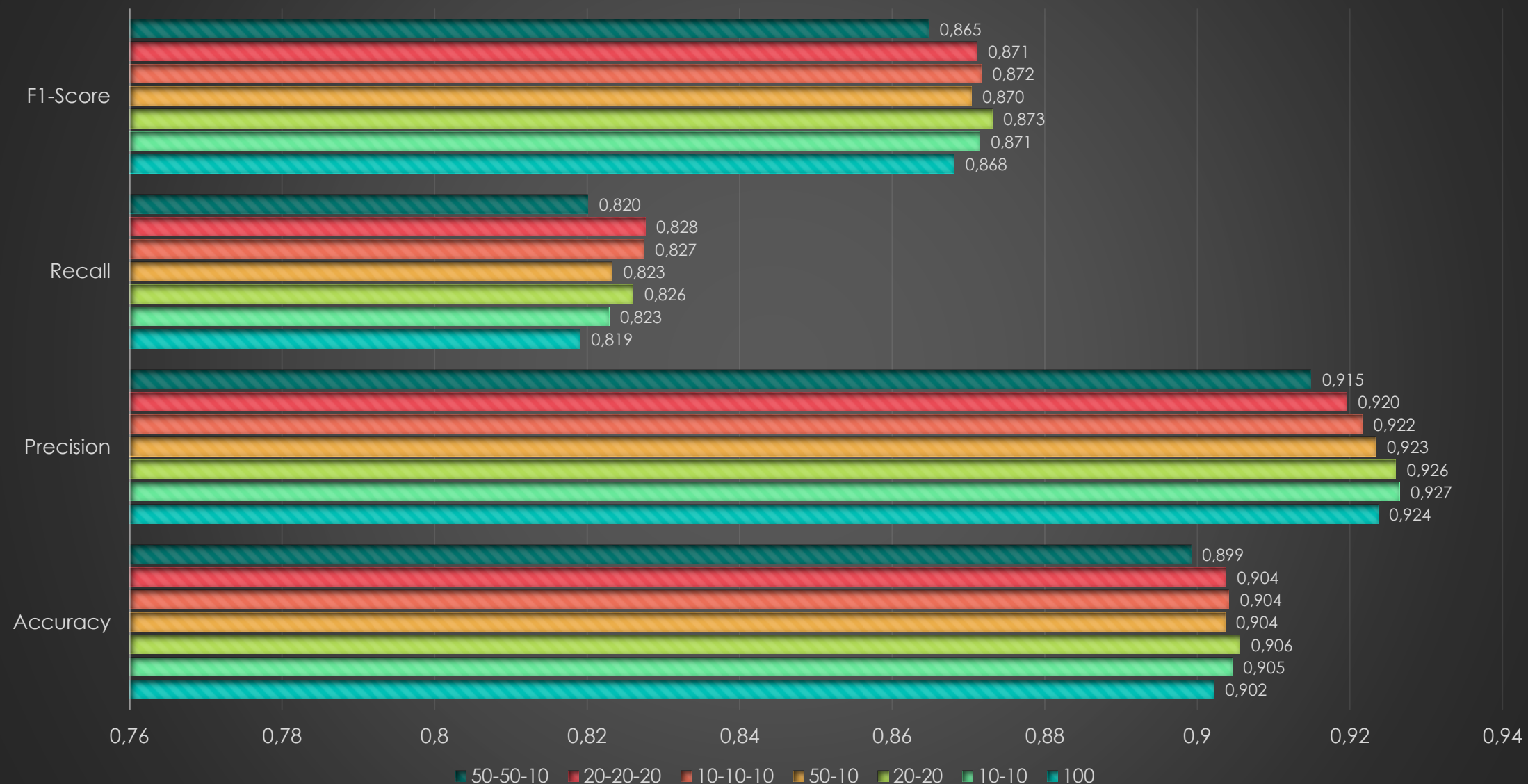Average accuracy with different number of most used tokens

13

Results with 2-gram count and tf-idf features

14

# Multi-layer Perceptrons (MLP)

- Activation function: ReLU

- Solver: Adam

- 500 epochs

- Early Stopping activated

- Different network configurations tested:

  - 100, 10-10, 20-20, 50-10, 10-10-10, 20-20-20, 50-50-10

- Different lambda values tested:

  - 0, 0.0001, 0.001, 0.01, 0.1, 1
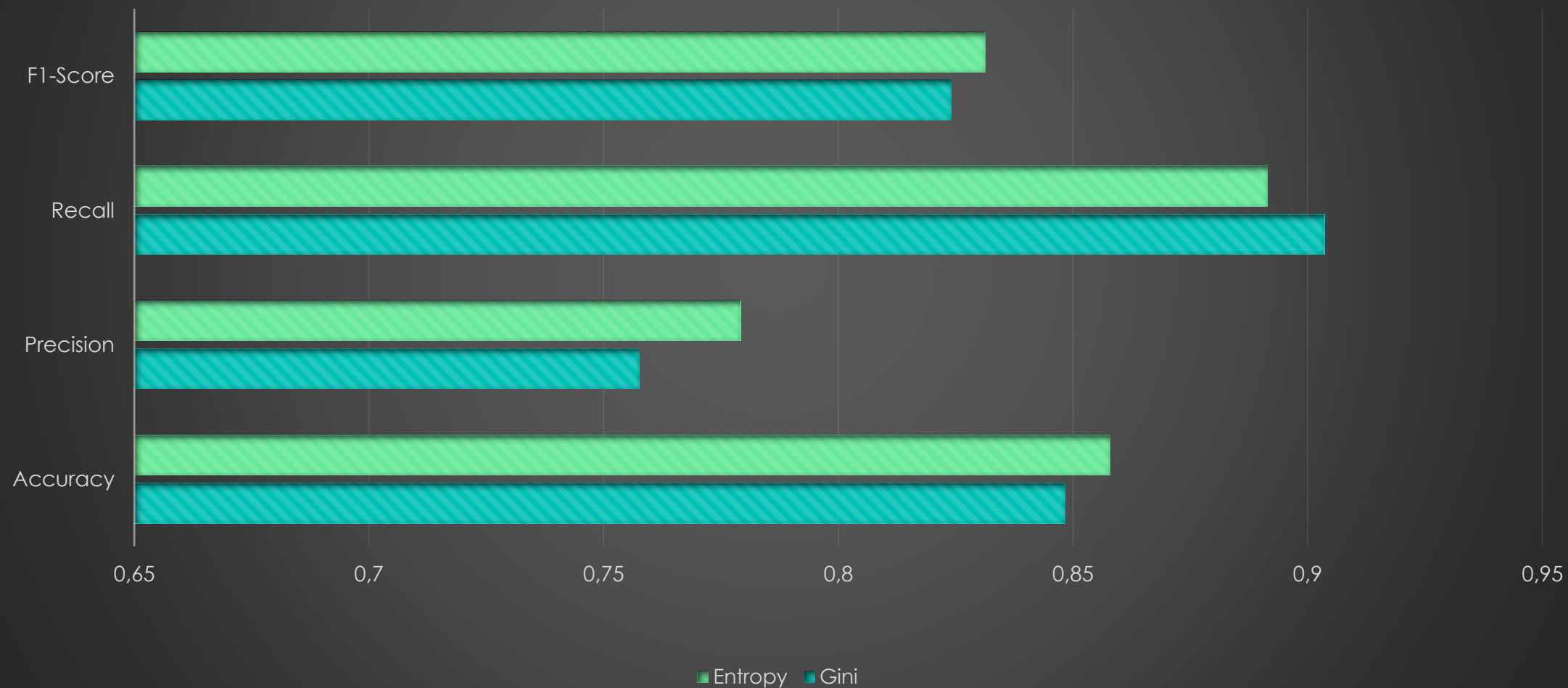
Results with different network configurations

Recall with different network configurations and lambda values
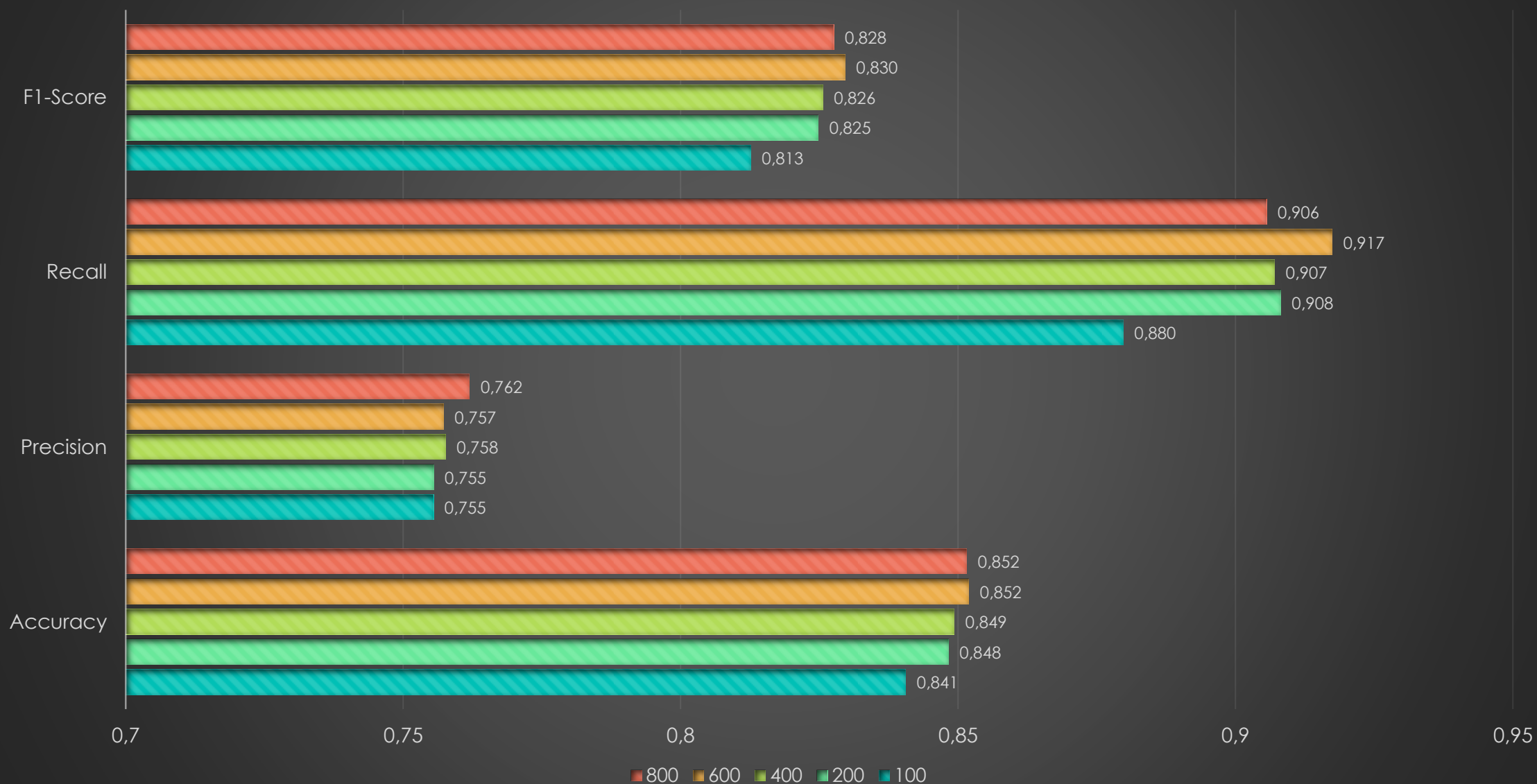
# Decision Tree

- Criterion: measures the quality of a split
  - Gini
  - Entropy
- Maximum depth of decision tree
  - 100
  - 200
  - 400
  - 600
  - No maximum depth (~800)

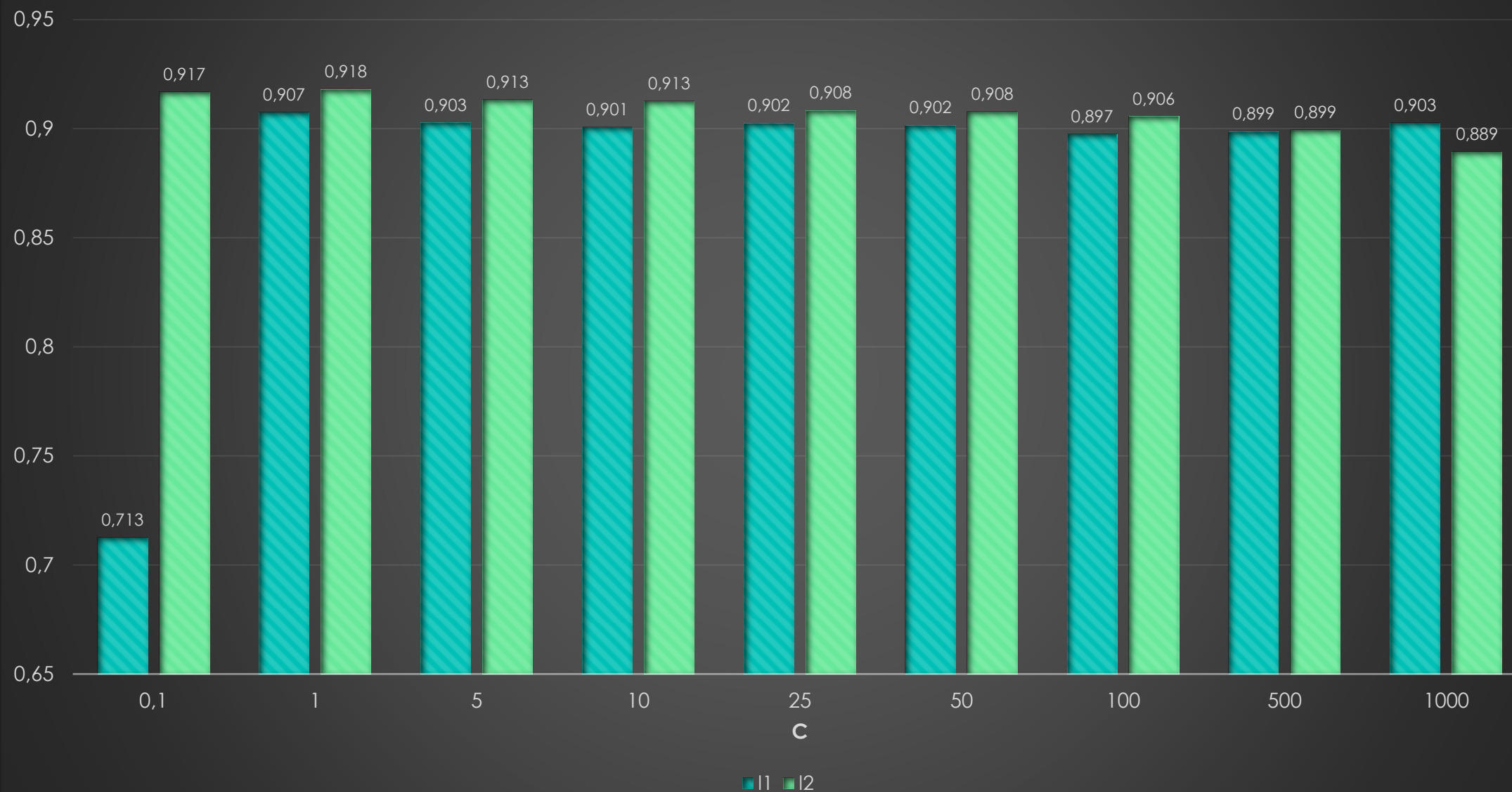**Results score when using different criterion for splitting nodes in tree**
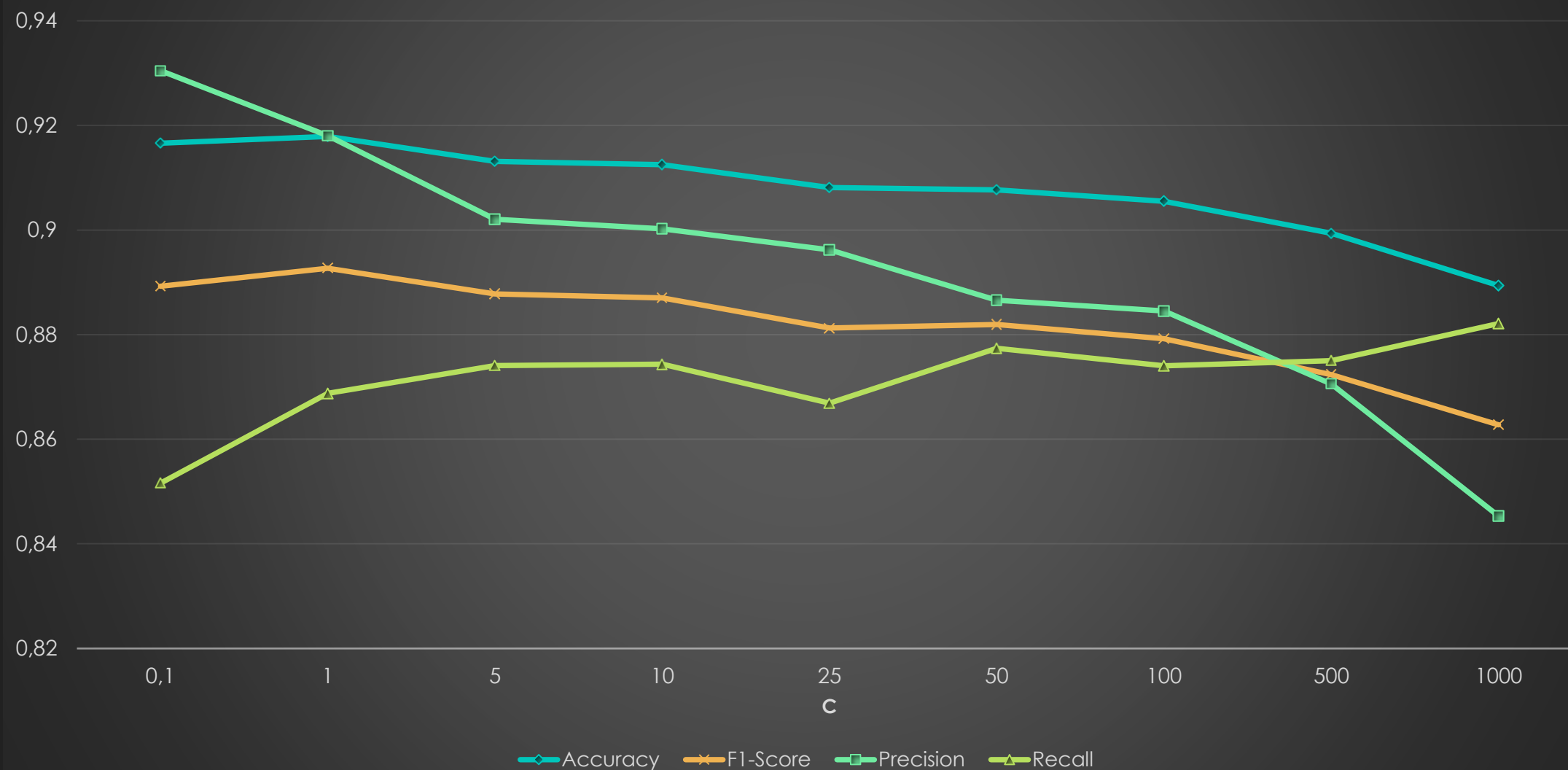


Entropy ■ Gini

# Linear SVC

- Penalty: norm used in the penalization
  - l1
  - l2
- C: penalty parameter C of the error term
  - 0.1
  - 1
  - 5
  - …
  - 1000

Accuracy between l1 and l2 penalty with different C values
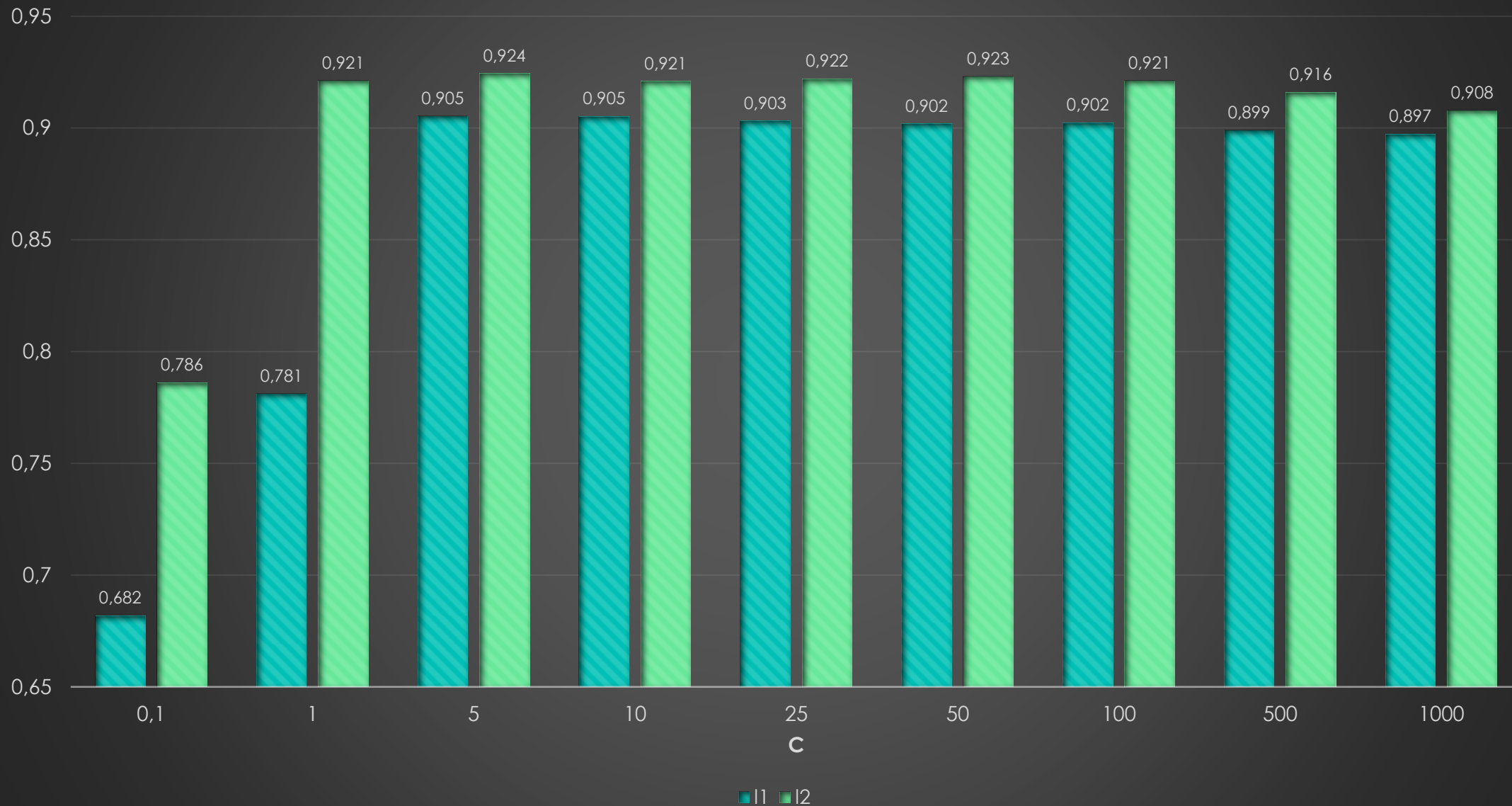
l2 penalty results with different C values

# Logistic Regression

- Penalty: norm used in the penalization
    - l1
    - l2

- C: penalty parameter C of the error term
    - 0.1
    - 1
    - 5
    - …
    - 1000

Accuracy between l1 and l2 penalty with different C values
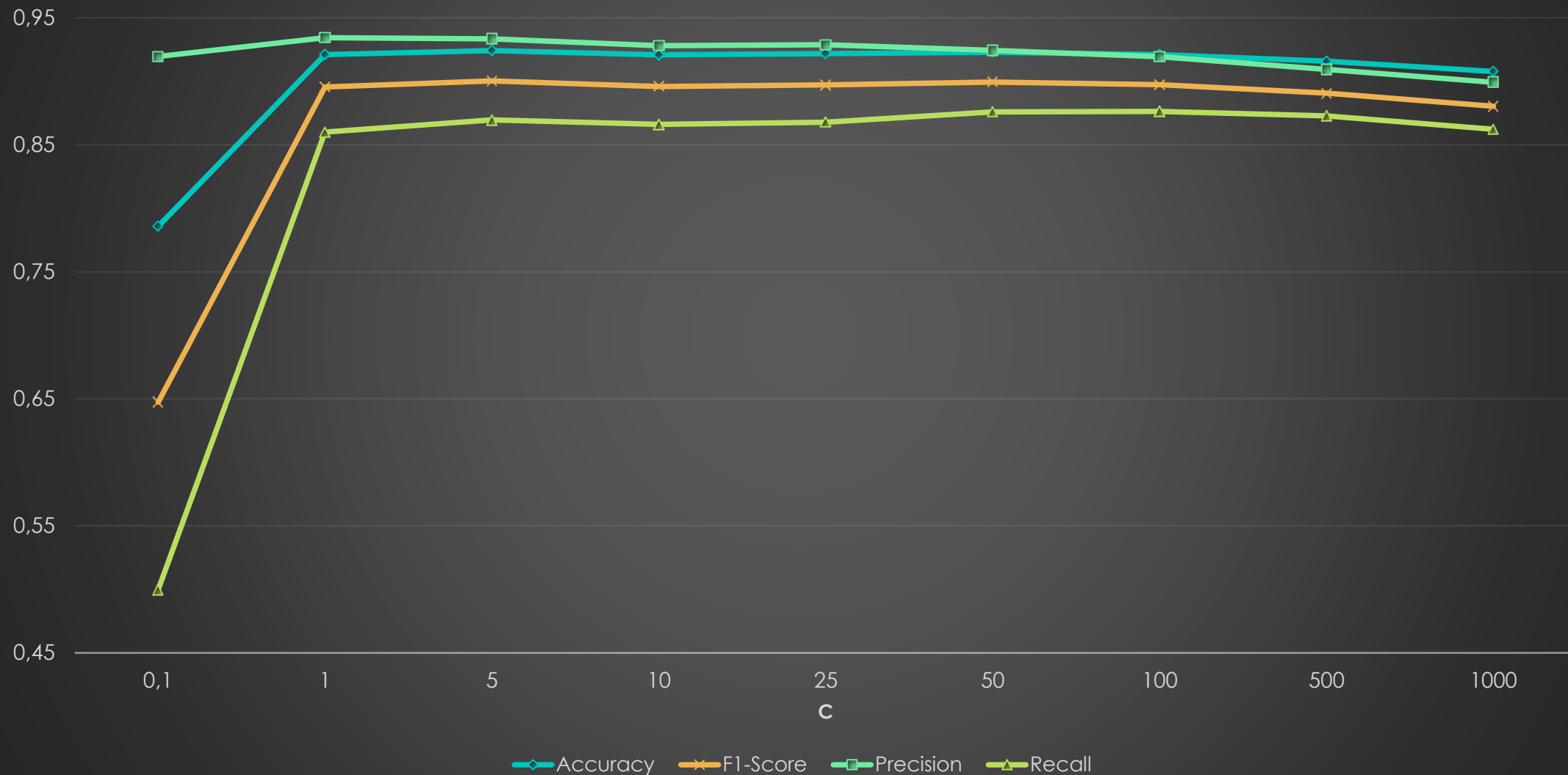
l2 penalty results with different C values

# Final results

| Classifier | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 0.956 | 0.940 | 0.946 | 0.943 |
| Linear SVC | 0.952 | 0.924 | 0.952 | 0.938 |
| Decision Tree | 0.871 | 0.764 | 0.959 | 0.850 |
| Multi-layer Perceptrons | 0.944 | 0.943 | 0.907 | 0.925 |

*Data division: 80% train and 20% test*