

Classification of Strong and Weak Passwords and Generation of Strong Passwords

Abstract

This project focuses on developing a machine learning-based approach for classifying strong and weak passwords, and also generating strong passwords. The initial phase involved experimenting with two datasets (ROCKYOU and PWNED), where features were created and the clustering techniques were employed to generate labels for password strength. Despite challenges with the datasets and the results not being satisfactory, the second phase utilized the PWLDS dataset, leading to impressive outcomes. A neural network classifier achieved high accuracy in password classification, and a Variational Autoencoder (VAE) was implemented to generate strong passwords. This report covers the complete process from dataset preprocessing to training models and evaluating their performance.

Introduction

In this digital era, maintaining personal and commercial data safe depends much on password security. Since breakable passwords are easy to guess, systems that can correctly select and generate strong passwords are quite crucial. This project investigates how to generate better passwords by means of machine learning by means of password strength sorting.

The project was divided into two main stages:

1. The first stage was about applying clustering techniques to produce labels and create enhanced features using the first datasets (ROCKYOU and PWNED).
2. Turning now to the PWLDS dataset, which already featured class names, features were developed and a neural network classifier trained to predict password strength. VAE was deployed and password generation then came from this variational autoencoder model.

Datasets

1. ROCKYOU and PWNED (Pond Dataset)

The study initially examined two datasets, known as Pwned and Rockyou. The PWNED dataset included the hacked copies of the passwords, hash values as well as the breach count. More elements could then be included: the length of the password, the count of alphabets, count of numerals, count of special characters, uppercase and lowercase letters, and their corresponding ratios. Unsupervised learning techniques like clustering had to be applied to label the passwords into weak and strong groups as no labels already existed.

2. PWLDS Dataset

The issues arising from the initial datasets led the PWLDS dataset to be chosen for the second phase. It had roughly five million passwords, and every one of them was assigned a one of five security strength label:

0: Very Weak

1: Weak

2: Normal

3: Strong

4: Very Strong

This dataset was far superior since it was ideal for the models requiring supervised learning since it included labelled data.

Stage 1: Experimentation with Clustering

Challenges

The PWNED dataset lacked predefined labels, which was the primary issue in the first stage. Clustering algorithms were applied in order to overcome this issue of labelling:

1. **K-Means Clustering:** Weak, average, and strong labels were created using the K-Means Clustering technique, but the results were not what anticipated.
2. **Agglomerative Clustering:** Another effort that also failed rather poorly was agglomerative clustering.
3. **Gaussian Mixture Models:** Though they gave us some ideas, Gaussian Mixture Models were too difficult to compute and unsuitable for this study.

Despite these setbacks, K-Means was selected as the best option for generating the initial password labels, even though the results were not ideal.

Stage 2: PWLDS Dataset - Neural Network Classifier

Dataset Preparation

The second stage drew on the PWLDS dataset with preset class labels. The dataset was preprocessed to extract features such as:

- Password length
- Counts of alphabets, numerics, special characters, uppercase and lowercase letters
- Repeated characters
- Uppercase to lowercase ratio

These elements then included the strength values (0–4) to produce the dataset utilized for neural network training.

Model Development

A neural network classifier (using PyTorch) was trained on the preprocessed dataset. The model aimed to predict the strength of a password based on the extracted features.

Performance Evaluation

The neural network classifier performed well, achieving:

- **94% validation accuracy**
- **93.7% testing accuracy**

Strong success across all classes was shown by the model's good F1, recall, and accuracy scores as well. Given both the training and validation curves kept near to the same point, it appeared the model was not overfitting.

Stage 3: Variational Autoencoder for Password Generation

Strong passwords in the last phase of the research were produced using a variational autoencoder (VAE). The VAE was trained using the PWLDS dataset, then learnt to generate passwords fitting into the "strong" and "very strong" classifications.

Results

Strong passwords made by the VAE indicate that it may be feasible to generate safe passwords automatically. This generation approach considered the issue of creating difficult to guess strong passwords.

Discussion, Conclusion and Future Work

The project progressed through significant experimentation:

1. **Clustering (Stage 1):** Despite the challenges with the unsupervised methods, this stage helped lay the groundwork for the subsequent phases by generating initial labels for password strength.
2. **Neural Network Classifier (Stage 2):** The neural network showed strong performance with the PWLDS dataset, providing a reliable method for classifying passwords based on their strength.
3. **Variational Autoencoder (Stage 3):** The VAE demonstrated its potential in password generation, offering a means of creating strong passwords with high complexity.

Turning now to the PWLDS dataset marked a turning point in the study that resulted in more dependable and practical models following the first stage's unsatisfactory outcome.

This project has solved the issues of classification and password creation. The VAE revealed that the neural network predictor could create secure passwords; it was also rather accurate. Our next focus will be:

1. Improving the VAE's ability to generate even stronger passwords.
2. Exploring additional generative models for password creation.
3. Expanding the dataset and refining the neural network architecture.

Overall, this system has the potential to enhance password security in various applications by automatically classifying and generating secure passwords.