# Human Activity Recognition Using Smartphones Data

Jens Kobler

February 10, 2025

## Contents

## 1 Main Objective of Analysis

The main objective of this analysis is to correctly classify test data from the human activity recognition dataset. This version of the report aims at delivering baseline results. Future version could include improvement recommendations stated in the last chapter.

This project was inspired by the ”IBM Machine Learning Professional Certificate”.

## 2 Desription of Data

Description from Website (Link):

The Human Activity Recognition Trondheim (HARTH) dataset is a professionally-annotated dataset containing 22 subjects wearing two 3-axial accelerometers for around 2 hours in a free-living setting. The sensors were attached to the right thigh and lower back. The professional recordings and annotations provide a promising benchmark dataset for researchers to develop innovative machine learning approaches for precise HAR in free living.The provided sampling rate is 50Hz. Video recordings of a chest-mounted camera were used to annotate the performed activities frame-by-frame.

Columns:

- timestamp: date and time of recorded sample

- back_x: acceleration of back sensor in x-direction (down) in the unit g

- back_y: acceleration of back sensor in y-direction (left) in the unit g

- back_z: acceleration of back sensor in z-direction (forward) in the unit g

- thigh_x: acceleration of thigh sensor in x-direction (down) in the unit g

- thigh_y: acceleration of thigh sensor in y-direction (right) in the unit g

- thigh_z: acceleration of thigh sensor in z-direction (backward) in the unit g

- label: annotated activity code

Activity code from original dataset: (1) walking, (2) running, (3) shuffling, (4) stairs (ascending), (5) stairs (descending), (6) standing, (7) sitting, (8) lying, (13) cycling (sit), (14) cycling (stand), (130) cycling (sit, inactive), (140) cycling (stand, inactive).

The file with the filename *S006.csv* is selected from the 22 different files. It is used to create baseline results. In Figure 1 the class distribution of the dataset is displayed. As visible the dataset is heavily unbalanced. Class 5 represents sitting (class labels are encoded (again); see next chapter).
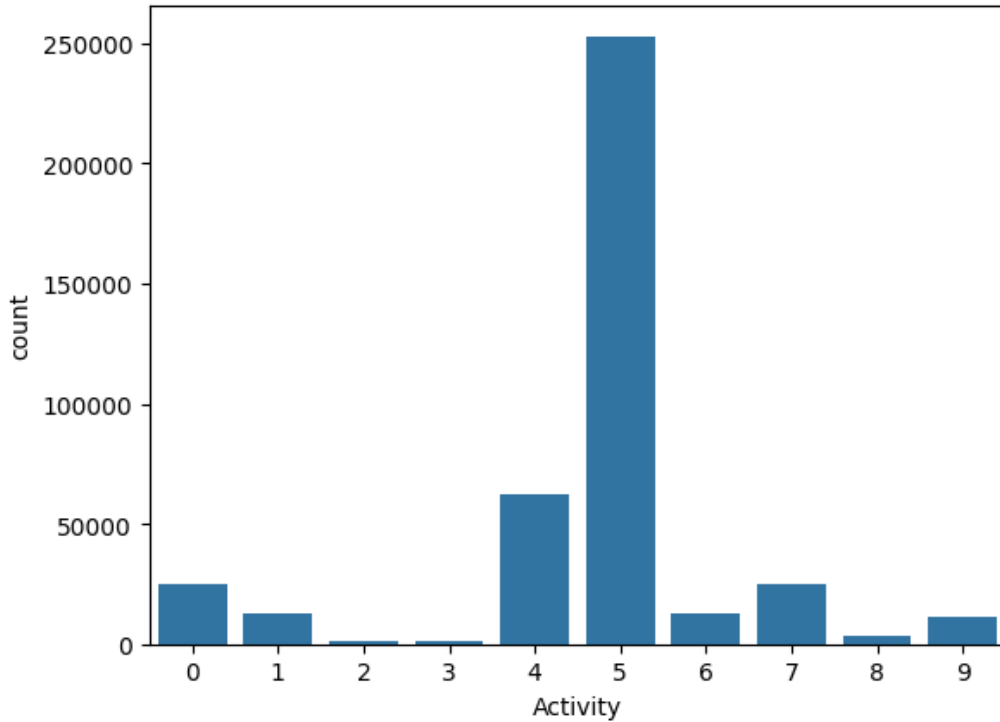


Figure 1: Class distribution.

# 3 Data Cleaning and Feature Engineering

1. delete timestamp column

2. scale features with minmax scaler

3. encode target labels

# 4 Variations of Supervised Learning Models

The dataset is split into train and test set with test size 0.2 in a stratified way.

## 4.1 Logistic Regression (LR)

The standard logistic regression model of the sklearn library is used with solver liblinear.

## 4.2 K Nearest Neighbor (KNN)

The standard K Nearest Neighbor model of the sklearn librar is used with different k $\in$ [5, 10, 15, 20].

## 4.3 Support Vector Machine (SVC)

The standard support vector machine of the sklearn library is used.

# 5 Final Result and reasoning

In Figure 2 the heatmaps for the different models are displayed. The logistic regression model is highly biased to class 4, 5 and 7. The LR model is not able to deal with the unbalanced dataset.
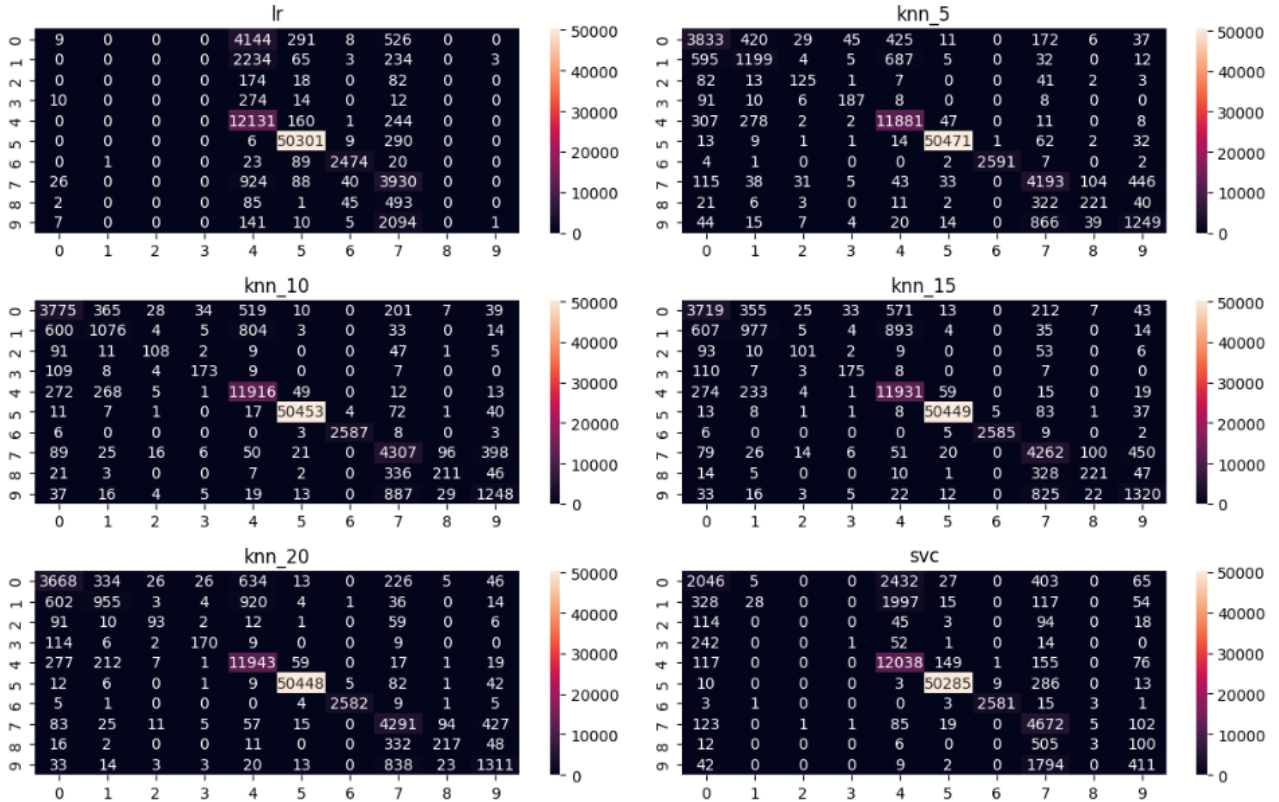


Figure 2: Heatmaps.

In Figure 3 the results for different metrics and models is illustrated. The K-Nearest-Neighbor algorithm achieves the highest results. The chocie between the different k is not having a big difference. Pick KNN algorithm for classification.
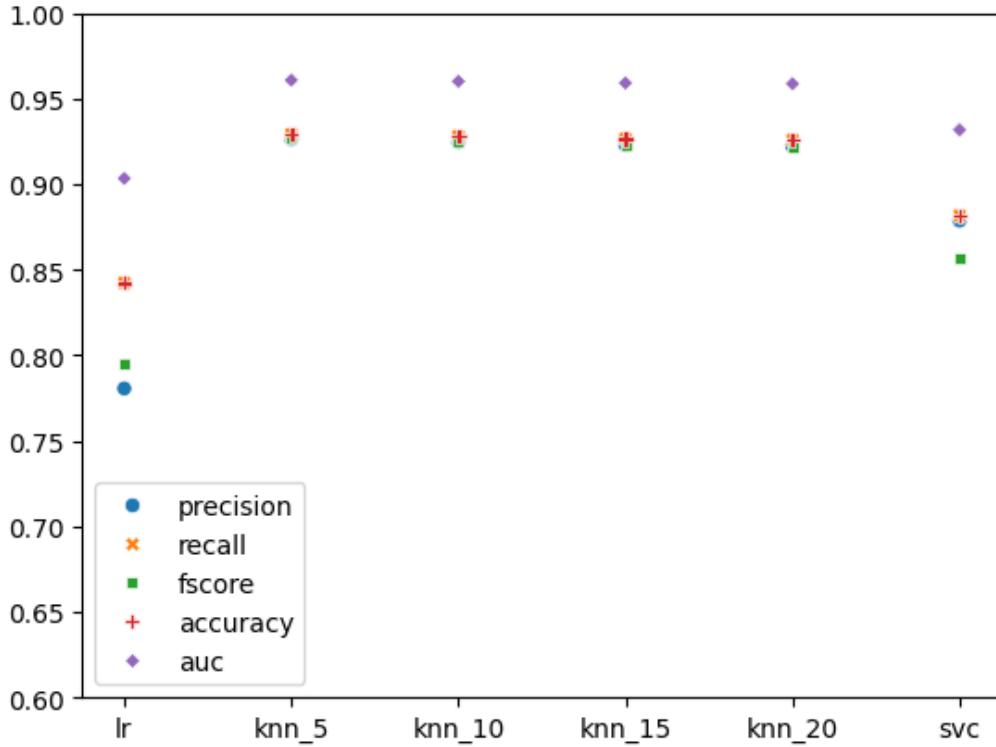
Figure 3: Results.

# 6   Flaws in model and Revisiting Analysis

- Only one out of 22 tables were used. Concatenate all tables from all participants.

- One could balance the classes. For example one can take the minimum of all classes as a number of samples.

- One could try different unsupervised (clustering) methods for example KMeans.

- One could train a deep neural network.

- Try different hyperparameters.

- One could try Principal Component Analysis (PCA), to reduce the number of features. However, all features seem necessary considering the physics.

- One could visualize the features space and the labeled samples with TSNE or PCA...