# Final Poject FDS: Stroke Prediction

Dipartimento di Ingegneria Informatica, Informatica e Statistica

Andrea Potì, Amedeo Ranaldi,
Onur Ergun, Giulio D'Erasmo

Master Degree in Data Science

A. A. 2021-2022

SAPIENZA
Università di Roma

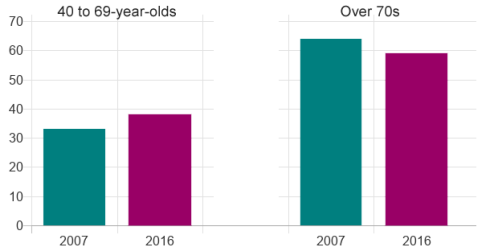# Table of contents

# Task and Motivation

## Task

Predicting the probability of a person having a stroke in the future.

## Motivation

Each year over 13 million people will have a stroke and 5.5 million people dies as a result. Our motivation is to save lives by predicting if an individual will have a stroke or not.

**First-time strokes are happening earlier in life in England**
Percentage by age group



Source: Public Health England

BBC

SAPIENZA
Università di Roma

# Task and Motivation

## Task

Predicting the probability of a person having a stroke in the future.

## Motivation

Each year over 13 million people will have a stroke and 5.5 million people dies as a result. Our motivation is to save lives by predicting if an individual will have a stroke or not.

## Related work

Stroke risk prediction using machine learning: a prospective cohort study of 0.5 million Chinese adults (Journal of the American Medical Informatics Association, 28(8), 2021, 1719–1727, 9 May 2021)

## Tentative material and methods

For the investigate the task and resolve it we choose to use and
compare the following **models** studied in the course and implemented
in the Assignment:

- Naive Bayes
- Logistic Regression
- Gaussian Discriminant Analysis (GDA)

# Benchmark Analysis

- The Stroke Prediction **dataset** that we will use has been taken from the following link :
  https://www.kaggle.com/fedesoriano/stroke-prediction-dataset

**Attribute Information**

1) id: unique identifier
2) gender: "Male", "Female" or "Other"
3) age: age of the patient
4) hypertension: 0 if the patient doesn't have hypertension, 1 if the patient has hypertension
5) heart_disease: 0 if the patient doesn't have any heart diseases, 1 if the patient has a heart disease
6) ever_married: "No" or "Yes"
7) work_type: "children", "Govt_jov", "Never_worked", "Private" or "Self-employed"
8) Residence_type: "Rural" or "Urban"
9) avg_glucose_level: average glucose level in blood
10) bmi: body mass index
11) smoking_status: "formerly smoked", "never smoked", "smokes" or "Unknown"*
12) stroke: 1 if the patient had a stroke or 0 if not
*Note: "Unknown" in smoking_status means that the information is unavailable for this patient

## Benchmark Analysis

- The Stroke Prediction **dataset** that we will use has been taken
  from the following link :
  https://www.kaggle.com/fedesoriano/stroke-prediction-dataset

- The **metric** we will to compare the models and the precision we
  obtain are: the loglikelooh history for the logistic regression and
  accuracy score for the Naive Bayes