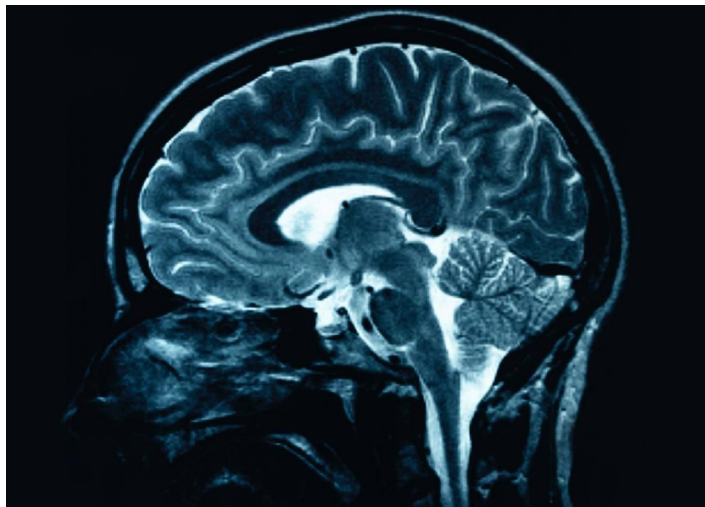


# Practical exercises

Advanced Machine Learning

# Task 1: Predict a person's age from brain image data



MRI image

# MRI

- Magnetic Resonance Imaging (MRI) is a key technology in medical imaging
- Non-invasive + Non-radiative investigation of sensitive organs (brain)
- Switzerland: 45.0 MRI units per 1.000.000 inhabitants (2016)  
West-Africa: 0.22 MRI units per 1.000.000 inhabitants (2018)

# MRI processing

## Raw brain scans are difficult to handle

- 3D brain scans are  $\sim 200 \times 200 \times 200 \sim 10^7$  features/voxels (3D pixels  $\sim 1\text{mm}^3$ )
- 3D structure + individual brain shapes  $\rightarrow$  difficult to recognize disease patterns
- Data is scarce:

	ImageNet	MRI data set for task 1
Image size	224 x 224 x 3	$\sim 200 \times 200 \times 200$
Data set size	1.2 million	$\sim 1200$

# MRI feature extraction

**We are using ~200 anatomical features for this project**

- Informative features derived from image data (with Freesurfer)
- No need to process big images (6 GB) → csv sheet (3 MB)
- No need for image analysis
- Meaningful features are extracted using domain knowledge
  - e.g. cortex volume, left/right hemisphere surface area, white/gray matter volume etc
- => Information loss

# Task 1: Age prediction

- We have modified the derived input data in three ways:
  1. Irrelevant features
  2. Outliers
  3. Perturbations (e.g. missing values, etc.)

# Description of the dataset

# File description

We provide the following files:

- `X_train.csv`, `y_train.csv`: the training set, including the features and labels
- `X_test.csv`: the test set (make predictions based on this file)
- `sample.csv`: a sample submission file in the correct format



# Task description

# Subtask 0: Filling missing values

## Background

There are some missing values in the data, originally they are set to NaN values. Most of the methods cannot handle them automatically. There are different strategies how to impute them: mean, median, most frequent etc.

## Task requirement

We require that students fill missing values in the training and the test set.

# Subtask 1: Outlier Detection

## Background

In the training set, there are some outliers. If the resulting model is not robust enough, it may be sensitive to the outliers. In this case, outliers deletion can be expected to lead to better results.

## Task requirement

We require that students build an outlier detection model to make classification for samples in the training set i.e. whether they are outliers.

# Subtask 2: Feature selection

## Background

To make the task a bit more challenging, we added some manual features to the FreeSurfer-processed dataset.

Feature selection is thus needed:

- Simplifies the models to make them easier to interpret
- Leads to shorter training times and makes the curse of dimensionality more manageable
- Better generalization by reducing overfitting

## Task requirement

We require that students use feature selection methods to label the features as selected features and unselected features.

Here, unselected features includes irrelevant features and redundant features.

# Main task: Age Prediction

## Background

After primary preprocessing and dimensionality reduction, now we finally arrive at the regression task.

## Task requirement

We require that students use suitable regression methods to predict the age of a person from his/her brain data.

# Evaluation metric

## Coefficient of Determination $R^2$

is the proportion of the variance in the dependent variable that is predictable from the independent variable.

$$R^2 := \frac{\sum_i (f_i - \bar{y})^2}{SS_{\text{tot}}} = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

$$SS_{\text{tot}} := \sum_i (y_i - \bar{y})^2$$

$$SS_{\text{res}} := \sum_i (f_i - y_i)^2$$

- Varies between 1 (best) and  $-\infty$  (worst)
- e.g. a regression model that always predicts the empirical mean of the predictor variable has  $R^2 = 0$

## How to compute it in Python:

from sklearn.metrics import r2\_score

```
score = r2_score(y_true, y_pred)
```

# Other considerations

- do NOT use AutoML packages
  - this includes anything that does automatic data cleaning and automatic model selection
- beware of overfitting on the public test data
- describe what you did when you hand in the project
  - keep your implementation for potential review
- do NOT wait until the last day to submit something
  - servers usually get overloaded and crash causing long waiting times

Q&A