# DATA SCIENCE

## 1.1 MRI AND ALZHEIMER

The `https://www.oasis-brains.org/`OASIS provides free datasets of MRI scans done on demented and non-demented people. They provide two free datasets containing demographic data accompanied with MRI scans. In this section we explore the demographic datasets OASIS-1 and OASIS-2[1, 2].

### 1.1.1 *Initial data exploration*

First we look at the structure of the data. The longitudinal set contains 373 scans and 15 rows an the cross-sectional set contains 436 scans and 12 rows. We plot the columns and data-types in 1. The length of the bar indicates the non-empty elopements. We observe that the longitudinal contains fewer empty elements than the cross-sectional set, hence we choose to start with this dataset for our regression and classification.

Subsequently, we check if there are easy to find correlations between the numerical columns. This is done by creating a correlation matrix heatmap seen in fig2. We want to attempt to predict the Clinical Dementia Rating or if a person is demented from this dataset, hence we look for correlations with the CDR rating. In both datasets it can be observed that there is an anti-correlation between
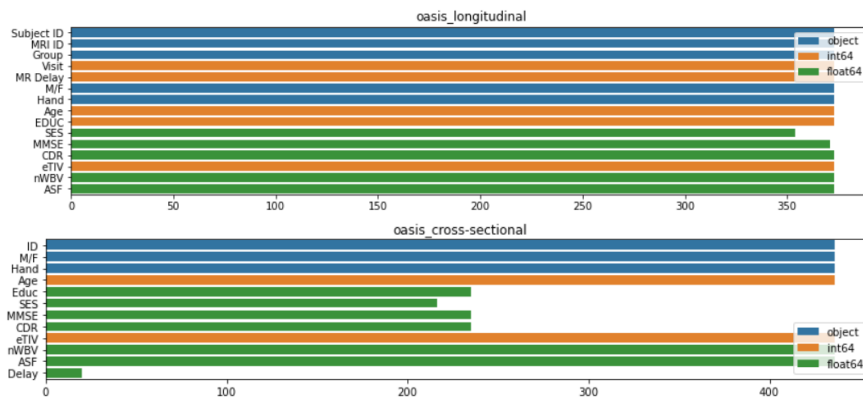


*Figure 1: Colour shows the datatype, barlength shows non-empty elements*

*(a) Correlation for Longitudal*          *(b) Correlation for Cross-sectional*
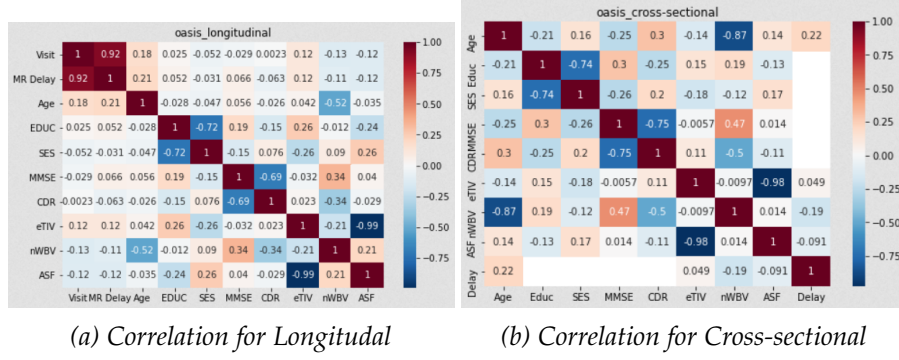
*Figure 2: Dark red is highly correlated, white is no correlation, and blue is anti-correlated.*

Mini Mental State Examination(MMSE) and Clinical Dementia Rating(CDR). We also observe small correlation with Normalize Whole Brain Volume(nWBV) and Education.

We visualize the correlations observed in fig3. The scatterplot in fig3a shows the correlation between CDR and MMSE. In fig3b and fig3c we visualise the correlation for nBWV and Education respectively.



*(a) Scatterplot showing the relation between MMSE and CDR.*



*(b) Density plots for nBWV.*          *(c) Density plots for Educ.*
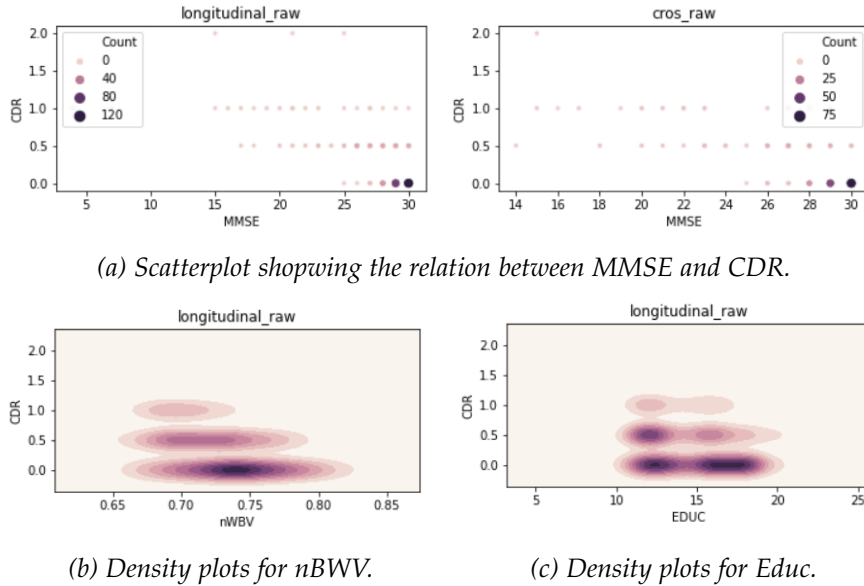
*Figure 3: Visualization of correlations*

### 1.1.2   Regression and Classification

Next we are going to see if we can predict the CDR using two machine learning techniques; regression and classification. We start out with training a XGBregressor using the features: Male/Female, Age, Education, nBWV, and MMSE. We train and test the regressor on differ-
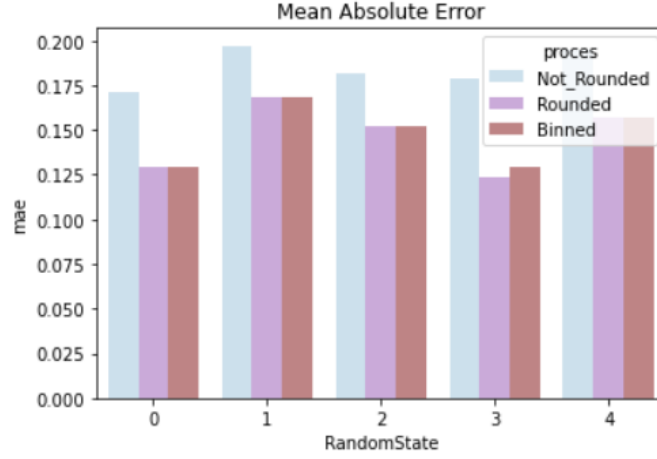
*Figure 4: Comparing the MAE for not rounded, rounded and binned regressor results*

ent data splits. The regressor returns a float indicating the predicted CDR, however CDR rating is bracketed in 0, 0.5, 1, and 2. Therefor we apply three post-processes, i.e, no post-processing (Not rounded), rounding the data, and binning the data. In fig4 we compare the mean absolute error(mae) of these different processes. We conclude that using rounded and binned is superior to not rounding when comparing the mae.



*(a) Confusion matrix for RS:0*



*(b) Confusion matrix for RS:1*



*(c) Confusion matrix for RS:2*



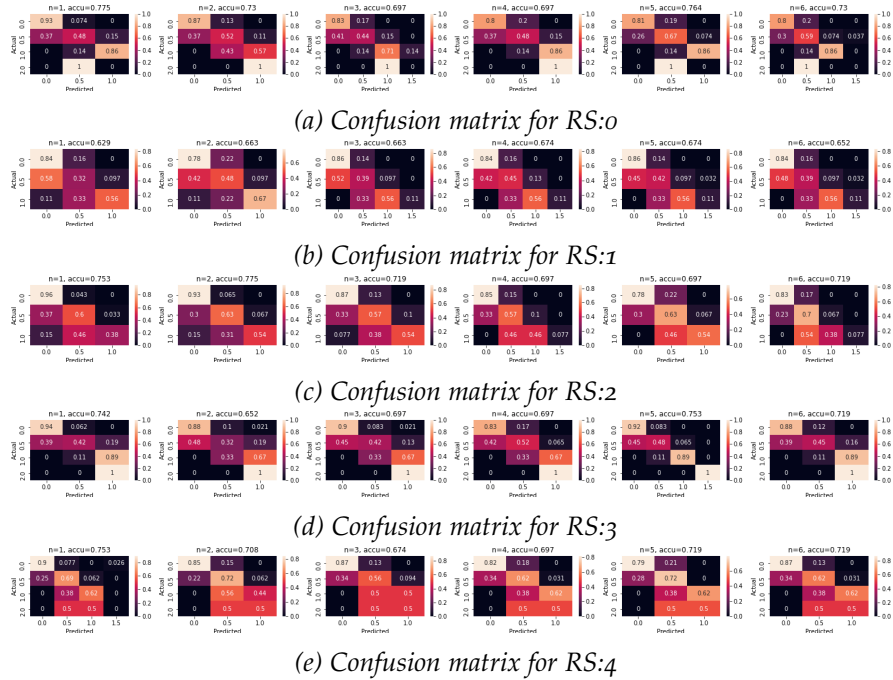*(d) Confusion matrix for RS:3*



*(e) Confusion matrix for RS:4*

*Figure 5: Confusion matrices for binned XGBregression.*

In fig5 the binned result confusion matrices have been plotted for different data-splits and max depth. It can be observed that the con-

fusion is rather dependent on the data splits, that is, it is very dependent on the supplied data.
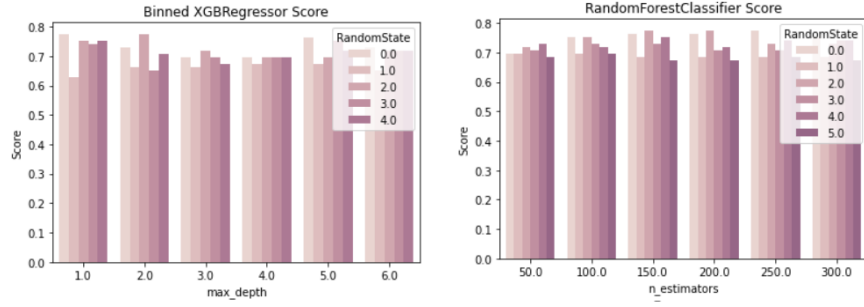
Next we look at a classifier instead of a regressor. In contrast to a regressor, which returns a continues float, a classifier classifies. Hence, we convert the CDR to categories. In this case we'll use the Random Forest Classifier with different number of estimators and compare the results. The confusion matrices for the RFC are plotted in fig6.
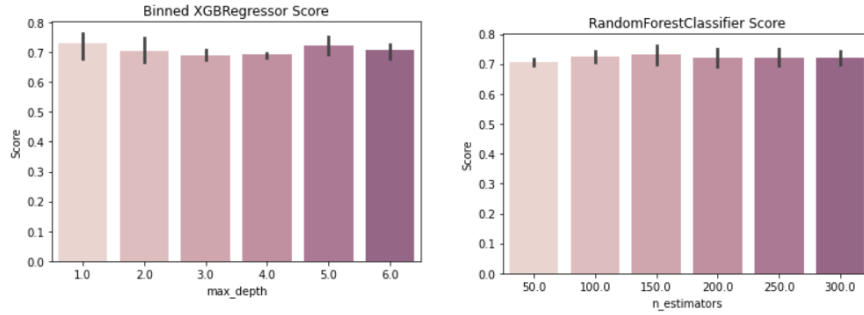


*(a) Confusion matrix for n_estimators=50 with average error=0.777*



*(b) Confusion matrix for n_estimators=100 with average error=0.785*



*(c) Confusion matrix for n_estimators=150 with average error=0.787*



*(d) Confusion matrix for n_estimators=200 with average error=0.79*



*(e) Confusion matrix for n_estimators=250 with average error=0.788*



*(f) Confusion matrix for n_estimators=300 with average error=0.788*

*Figure 6: Confusion matrices for RandomForestClassifier.*

We observe that in general most confusion is with neighboring categories except in splitting using random state 4 and 5. Again this is more split than model dependent. In all splits there are only one, two ore three rows with a CDR of 2, hence they are hard to predict.
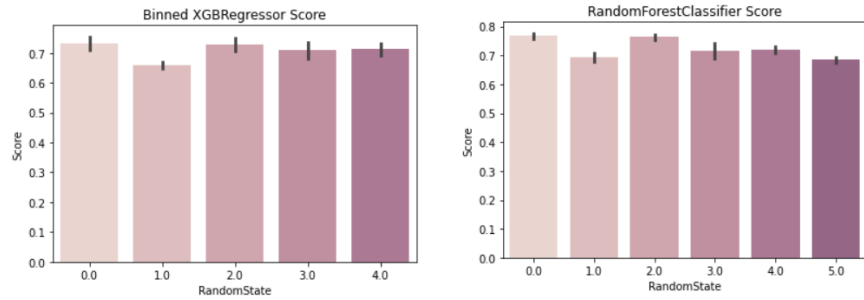
In fig7 we plot the score for the different models and Random states. Both models have similar results with a score around 0.7. Tuning the two variables explored here to obtain higher average scores also resulted in a larger confidence intervals.

*(a) XGBregression score for different max depths on different data splits.*

*(b) Random forest classifier score for different n on different data-splits*

*(c) XGBregression score for different max depths on different data splits.*

*(d) Random forest classifier score for different n on different data-splits*

*(e) XGBregression score for different random state splits.*

*(f) Random forest classifier score for different random state splits*

*Figure 7: Comparing regression and classifier for different n and RS*

1.1.3   *Notes for the future*

There are a few takeaways for future projects. Considering the data
used it would be great to have a larger sample size, especially for
categories with few samples.

Looking at the data representation and exploration, I have the fol-
lowing notes

- Having more consistency in explored model variables.

- Having more consistency in representation of the data.

- Having more colour consistency.

- Creating confusion matrices averaged on different splittings.

- Creating barplots with barcolours dependent on their values.

- Creating a heatmap showing scores depending on RS and vari-
  ables.

- Showing feature importance for the different models.

### 1.1.4 Code snippets

*Code: Data exploration*

```python
def data_analysis_plot(data):
    #
    length = data.shape[1]
    lengthplot = length/5
    plt.figure(figsize=(14,lengthplot))

    ax = sns.barplot(y=data.columns,x=data.notnull().sum(),
    hue=data.dtypes,dodge=False)
    ax.set_title(data.name)

    plt.figure(figsize=(8,5))
    ax2 = sns.heatmap(data.corr(),annot=True,cmap='RdBu_r')
    ax2.set_title(data.name)

    plt.figure(figsize=(8,5))
    ax3 = sns.heatmap(data.corr()[["CDR"]],annot=True,cmap='
    RdBu_r')
    ax3.set_title(data.name)

def data_info(data):
    print("Shape: ")
    print(data.shape)
    print("Columns: ")
    print(data.columns)

sets = [longitudinal_raw,cros_raw]
for datasets in sets:
    print(datasets.name)
    data_info(datasets)
    data_analysis_plot(datasets)
```

*Code: Data Scatterplot MMSE-CDR*

```python
fused = longitudinal_raw.groupby(["MMSE","CDR"]).count()
fused["Count"]=fused["MRI ID"]
fused.head()
fused2 = cros_raw.groupby(["MMSE","CDR"]).count()
fused2["Count"]=fused2["ID"]
fused.head()
figz, axs = plt.subplots(ncols=2,figsize=(12,2.5))
sns.scatterplot(x=fused.index.get_level_values("MMSE"),y=fused
    .index.get_level_values("CDR"),data=fused,size="Count",
    hue="Count",ax=axs[0]).set_title("longitudinal_raw")
sns.scatterplot(x=fused2.index.get_level_values("MMSE"),y=
    fused2.index.get_level_values("CDR"),data=fused2,size="
    Count", hue="Count",ax=axs[1]).set_title("cros_raw")
```

*Code: Regression*

```python
def multi_plot(y_test,prediction,n,vari):
    accu = np.round(accuracy_score(y_test, prediction),3)
    conf = pd.crosstab(y_test, prediction, rownames=['Actual'
    ], colnames=['Predicted'],normalize='index')
```

```
4      #conf = confusion_matrix(y_test,prediction,normalize='true
       ')
5      sns.heatmap(conf, annot=True,ax=multi_ax[n-1]).set(xlabel=
       'Predicted', ylabel='Actual',title=vari+"="+str(n)+", "+"
       score="+accu.astype('str'))
6      #plot_confusion_matrix(model,X_test, y_test,normalize='
       true')
7      return accu
8
9  def test_model2(X,y,n,rs):
10      X_train, X_test, y_train, y_test = train_test_split(X, y,
       random_state=rs)
11      model = XGBRegressor(random_state=0,max_depth=n)
12      model.fit(X_train,y_train)
13      prediction = model.predict(X_test)
14      rounded_pred = np.absolute(np.round(2*prediction)/2)
15      mae = mean_absolute_error(rounded_pred,y_test)
16      score = multi_plot(y_test.astype(str),rounded_pred.astype(
       str),n,vari="n")
17      return score
18
19 #Create plots
20 l=7
21
22 reg_score = pd.DataFrame(columns=["max_depth","RandomState","
       Score"])
23
24 for i in range(5):
25      rs = i
26      multi_fig, multi_ax = plt.subplots(ncols=l-1,figsize=(l
       *3,2))
27      multi_fig.tight_layout()
28      for n in range(1,l):
29          score = test_model2(X,y,n,rs)
30          reg_score.loc[str(n)+str(rs)]=[n,rs,score]
31      #output = multi_fig
32      name = 'PP_dataS_MRIalz_RS'+str(i)+'.png'
33      multi_fig.savefig(name)
34      print(r"\begin{subfigure}{\textwidth}\includegraphics[
       width=\textwidth]{Pictures/datas/"+name+
35 "}\caption{Confusion matrix for RS:"+str(i)+"}\end{subfigure}"
       )
```

*Code: Classification*

```
1  #GradientBoostingClassifier,RandomForestClassifier,
2  def classie(X,y,n,rs):
3      X_train, X_test, y_train, y_test = train_test_split(X, y,
       random_state=rs)
4      nn = n*50
5      model = RandomForestClassifier(n_estimators=nn)
6      model.fit(X_train,y_train)
7      prediction = model.predict(X_test)
8      multi_plot(y_test,prediction,rs,vari="RS")
9      accu = accuracy_score(y_test, prediction)
10     return accu
11
12
```

```
13
14  X = longitudinal[features]
15  #y_g = longitudinal["Group"]
16  y_g = longitudinal["CDR"].astype('str')
17
18  l=7
19  nplots = 6
20
21  classi_error = pd.DataFrame(columns=["n_estimators","
        RandomState","Score"])
22  for n in range(1,l):
23      multi_fig, multi_ax = plt.subplots(ncols=nplots,figsize=(l
        *3,2))
24      #multi_fig.tight_layout()
25      multi_fig.subplots_adjust(left=None, bottom=None, right=
        None, top=None, wspace=1, hspace=None)
26      errtot = 0
27      for rs in range(nplots):
28          c_score = classie(X,y_g,n,rs)
29          errtot+= c_score
30          classi_error.loc[str(n)+str(rs)]=[n*50,rs,c_score]
31      err = errtot/nplots
32      #print(err)
33      name = 'PP_dataS_MRIalz_RFCn='+str(n)+'.png'
34      multi_fig.savefig(name)
35      print(r"\begin{subfigure}{\textwidth}\includegraphics[
        width=\textwidth]{Pictures/datas/"+name+
36  r"}\caption{Confusion matrix for n{\textunderscore}estimators=
        "+str(n*50)+" with average error="+str(round(err,3))+"}\
        end{subfigure}")
```

*Code: Full code*

```
1   # %% [code]
2   # This Python 3 environment comes with many helpful analytics
        libraries installed
3   # It is defined by the kaggle/python Docker image: https://
        github.com/kaggle/docker-python
4   # For example, here's several helpful packages to load
5
6   import numpy as np # linear algebra
7   import pandas as pd # data processing, CSV file I/O (e.g. pd.
        read_csv)
8   import seaborn as sns
9   import matplotlib.pyplot as plt
10
11  # Input data files are available in the read-only "../input/"
        directory
12  # For example, running this (by clicking run or pressing Shift
        +Enter) will list all files under the input directory
13
14  import os
15  for dirname, _, filenames in os.walk('/kaggle/input'):
16      for filename in filenames:
17          print(os.path.join(dirname, filename))
18
```

```python
19  # You can write up to 5GB to the current directory (/kaggle/
        working/) that gets preserved as output when you create a
        version using "Save & Run All"
20  # You can also write temporary files to /kaggle/temp/, but
        they won't be saved outside of the current session
21
22  print("Setup done")
23
24  # %% [code]
25  long_path = "/kaggle/input/mri-and-alzheimers/
        oasis_longitudinal.csv"
26  longitudinal_raw = pd.read_csv(long_path)
27  longitudinal_raw.name = "oasis_longitudinal"
28
29  cross_path = "/kaggle/input/mri-and-alzheimers/oasis_cross-
        sectional.csv"
30  cros_raw = pd.read_csv(cross_path)
31  cros_raw.name = "oasis_cross-sectional"
32  print("Datasets done")
33
34  # %% [code]
35  def data_analysis_plot(data):
36      #
37      length = data.shape[1]
38      lengthplot = length/5
39      plt.figure(figsize=(14,lengthplot))
40
41      ax = sns.barplot(y=data.columns,x=data.notnull().sum(),
        hue=data.dtypes,dodge=False)
42      ax.set_title(data.name)
43
44      plt.figure(figsize=(8,5))
45      ax2 = sns.heatmap(data.corr(),annot=True,cmap='RdBu_r')
46      ax2.set_title(data.name)
47
48      plt.figure(figsize=(8,5))
49      ax3 = sns.heatmap(data.corr()[["CDR"]],annot=True,cmap='
        RdBu_r')
50      ax3.set_title(data.name)
51
52  def data_info(data):
53      print("Shape: ")
54      print(data.shape)
55      print("Columns: ")
56      print(data.columns)
57
58  sets = [longitudinal_raw,cros_raw]
59  for datasets in sets:
60      print(datasets.name)
61      data_info(datasets)
62      data_analysis_plot(datasets)
63
64
65  # %% [code]
66  meaning = {"SES":"Socioeconomic Status","MMSE":"Mini Mental
        State Examination","CDR":"Clinical Dementia Rating",
```

```
67          "eTIV":"Estimated  Total  intracranial Voluma","nWBV"
        :"Normalize Whole Brain Volume","ASF":"Atlas Scaling
        Factor"}
68
69
70
71
72
73 # %% [code]
74 def obj_cols(data):
75     s = (data.dtypes == 'object')
76     object_cols = list(s[s].index)
77     print("Dataset: "+data.name)
78     for i in object_cols:
79         if ("ID" in i) != True:
80             print(i+" unique values:")
81             print(data[i].value_counts())
82     return object_cols
83
84 for data in sets:
85     obj_cols(data)
86
87 # %% [code]
88 #Find dupliacte IDs
89 def non_uniq(data):
90     print(data.value_counts()[data.value_counts().values > 1])
91
92 for i in [longitudinal_raw["Subject ID"],longitudinal_raw["MRI
        ID"],cros_raw["ID"]]:
93     non_uniq(i)
94 cros_raw = cros_raw.rename(columns={"Educ":"EDUC"})
95 cros_raw.columns
96
97 # %% [code]
98 corcols = ["MMSE","nWBV","EDUC"]
99 for i in corcols:
100    fused = longitudinal_raw.groupby([i,"CDR"]).count()
101    fused["Count"]=fused["MRI ID"]
102    fused.head()
103    fused2 = cros_raw.groupby([i,"CDR"]).count()
104    fused2["Count"]=fused2["ID"]
105    fused.head()
106    figz, axs = plt.subplots(ncols=2,figsize=(12,2.5))
107    sns.scatterplot(x=fused.index.get_level_values(i),y=fused.
        index.get_level_values("CDR"),data=fused,size="Count", hue
        ="Count",ax=axs[0]).set_title("longitudinal_raw")
108    sns.scatterplot(x=fused2.index.get_level_values(i),y=
        fused2.index.get_level_values("CDR"),data=fused2,size="
        Count", hue="Count",ax=axs[1]).set(title="cros_raw")
109
110 # %% [code]
111 for i in corcols:
112    cm = sns.cubehelix_palette(light=1, as_cmap=True)
113    figz2, axss = plt.subplots(ncols=2,figsize=(12,2.5))
114    sns.kdeplot(longitudinal_raw[i], longitudinal_raw["CDR"],
        shade=True,cmap=cm ,ax=axss[0]).set_title("
        longitudinal_raw")
```

```
115     sns.kdeplot(cros_raw[i], cros_raw["CDR"],shade=True,cmap=
        cm ,ax=axss[1]).set_title("cros_raw")
116
117 # %% [code]
118 longitudinal = longitudinal_raw.copy()
119 longitudinal["M/F"]=longitudinal["M/F"].map({"M":0,"F":1})
120 print(longitudinal.head())
121 print(longitudinal.shape)
122 longitudinal = longitudinal.dropna()
123 print(longitudinal.shape)
124
125 # %% [code]
126 print("Import ML modules")
127 from sklearn.model_selection import train_test_split
128 from sklearn.ensemble import RandomForestRegressor
129 from sklearn.metrics import mean_absolute_error
130 from sklearn.metrics import accuracy_score
131 from xgboost import XGBRegressor
132 print("Importing done")
133
134 # %% [code]
135 from sklearn.model_selection import cross_val_score
136 from sklearn.pipeline import Pipeline
137
138 featuress=["M/F","Age","EDUC","nWBV","MMSE","eTIV"]
139 features=featuress[0:5]
140 print("Features used:")
141 print(features)
142 X = longitudinal[features]
143 y = longitudinal["CDR"]
144
145 def get_score(line):
146     # Multiply by -1 since sklearn calculates *negative* MAE
147     scores = -1 * cross_val_score(line, X, y,
148                                   cv=5,
149                                   scoring='
        neg_mean_absolute_error')
150     print("Average MAE score:", scores.mean())
151     return scores.mean()
152
153 def runmodels(X,y):
154     rf_model = RandomForestRegressor(random_state=1,
        n_estimators=100)
155     XGB_model = XGBRegressor(random_state=0)
156     models = [rf_model,XGB_model]
157     for model in models:
158         my_pipeline = Pipeline(steps=[
159                 ('model', model)
160                 ])
161         get_score(my_pipeline)
162
163 runmodels(X,y)
164
165 # %% [code]
166 prediction = [0,0.4,1.1,1.4,2,2.5]
167 binned_pred = pd.cut(prediction,[-1,0.25,0.75,1.5,10],labels
        =[0,0.5,1,2])
```

```python
168  print(binned_pred)
169  print(binned_pred.astype('str'))
170
171  # %% [code]
172  def show_accuracy(y_test, prediction, vari, n):
173      score = accuracy_score(y_test, prediction)
174      print(score)
175      conf = pd.crosstab(y_test, prediction, rownames=['Actual'
         ], colnames=['Predicted'], normalize='index')
176      #conf = confusion_matrix(y_test, prediction, normalize='true
         ')
177      figy, axconf = plt.subplots(figsize=(5,2))
178      tit = vari+": "+str(n)+", score:"+str(round(score,3))
179      axconf = sns.heatmap(conf, annot=True).set(xlabel='
         Predicted', ylabel='Actual', title=tit)
180      #plot_confusion_matrix(model, X_test, y_test, normalize='
         true')
181      return accuracy_score
182
183  def test_model(X,y,n):
184      X_train, X_test, y_train, y_test = train_test_split(X, y,
         random_state=n)
185      XGB_model = XGBRegressor(random_state=0,max_depth=5 )
186      XGB_model.fit(X_train, y_train)
187      prediction = XGB_model.predict(X_test)
188      mae1 = mean_absolute_error(prediction, y_test)
189      print("Not rounded: ",mae1)
190      rounded_pred = np.absolute(np.round(2*prediction)/2)
191      mae2 = mean_absolute_error(rounded_pred, y_test)
192      print("Rounded: ",mae2)
193      binned_pred = pd.cut(prediction,[-1,0.25,0.75,1.5,10],
         labels=[0,0.5,1,2],right=False)
194      mae3 = mean_absolute_error(binned_pred, y_test)
195      print("Binned: ",mae3)
196      vari = "Split random state"
197      show_accuracy(y_test.astype(str),binned_pred.astype(str),
         vari,n)
198      return [mae1,mae2,mae3]
199
200  regr_mae = pd.DataFrame(columns=["RandomState","proces","mae"
         ])
201  proc=["Not_Rounded","Rounded","Binned"]
202  for n in range(5):
203      result=test_model(X,y,n)
204      for i in range(len(result)):
205          regr_mae.loc[str(n)+str(i)]=[n]+[proc[i]]+[result[i]]
206  print(regr_mae)
207  plt.show()
208  sns.barplot(x="RandomState", y="mae",data=regr_mae,hue="proces
         ", palette=sns.light_palette("green")).set(title="Mean
         Absolute Error")
209
210  # %% [code]
211  plt.show()
212  #clors = sns.color_palette("RdYlGn", 5)
213  clors = sns.color_palette("cubehelix_r",6)
```

```
214 sns.barplot(x="RandomState", y="mae",data=regr_mae,hue="proces
        ", palette=clors).set(title="Mean Absolute Error")
215
216 # %% [code]
217 from sklearn.preprocessing import FunctionTransformer
218 XGB_model = XGBRegressor(random_state=0)
219 trans = FunctionTransformer(np.round, validate=True)
220 le_line = Pipeline(steps=[('Custom transformation',trans),
221                 ('model', XGB_model)]
222                 )
223
224 scores = -1 * cross_val_score(le_line, X, y,
225                                     cv=5,
226                                     scoring='
        neg_mean_absolute_error')
227 print(scores)
228 print("Average MAE score:", scores.mean())
229
230 # %% [code]
231 print("Doing Classifiers")
232 from sklearn.ensemble import GradientBoostingClassifier,
        RandomForestClassifier
233 from sklearn.metrics import confusion_matrix
234 from sklearn.metrics import plot_confusion_matrix
235
236
237 def classifiers(longitudinal):
238     featuress=["M/F","Age","EDUC","nWBV","MMSE","eTIV"]
239     features=featuress[0:5]
240     print(features)
241     X = longitudinal[features]
242     y = longitudinal["CDR"]
243     y_trans = y.copy().astype(str)
244     X_train, X_test, y_train, y_test = train_test_split(X,
        y_trans, random_state=1)
245
246     #GBC = GradientBoostingClassifier(n_estimators=200)
247     #RFC = RandomForestClassifier(n_estimators = 200)
248     #classifiers = [GBC,RFC]
249
250
251
252     for n in range(2,3):
253         vari = "n_estimators"
254         esti = n*50
255         model = RandomForestClassifier(n_estimators=esti)
256         model.fit(X_train, y_train)
257         prediction = model.predict(X_test)
258         show_accuracy(y_test, prediction, vari,n)
259         print(model.feature_importances_)
260
261
262
263 classifiers(longitudinal)
264
265
266 # %% [code]
```

```python
267 def multi_plot(y_test, prediction, n, vari):
268     accu = np.round(accuracy_score(y_test, prediction),3)
269     conf = pd.crosstab(y_test, prediction, rownames=['Actual'
        ], colnames=['Predicted'], normalize='index')
270     #conf = confusion_matrix(y_test, prediction, normalize='true
        ')
271     sns.heatmap(conf, annot=True, ax=multi_ax[n-1]).set(xlabel=
        'Predicted', ylabel='Actual', title=vari+"="+str(n)+", "+"
        score="+accu.astype('str'))
272     #plot_confusion_matrix(model, X_test, y_test, normalize='
        true')
273     return accu
274
275 def test_model2(X, y, n, rs):
276     X_train, X_test, y_train, y_test = train_test_split(X, y,
        random_state=rs)
277     model = XGBRegressor(random_state=0, max_depth=n)
278     model.fit(X_train, y_train)
279     prediction = model.predict(X_test)
280     rounded_pred = np.absolute(np.round(2*prediction)/2)
281     mae = mean_absolute_error(rounded_pred, y_test)
282     score = multi_plot(y_test.astype(str), rounded_pred.astype(
        str), n, vari="max depth")
283     return score
284
285 #Create plots
286 l=7
287
288 reg_score = pd.DataFrame(columns=["max_depth","RandomState","
        Score"])
289
290 for i in range(5):
291     rs = i
292     multi_fig, multi_ax = plt.subplots(ncols=l-1, figsize=(l
        *3,2))
293     multi_fig.tight_layout()
294     for n in range(1,l):
295         score = test_model2(X, y, n, rs)
296         reg_score.loc[str(n)+str(rs)]=[n, rs, score]
297     #output = multi_fig
298     name = 'PP_dataS_MRIalz_RS'+str(i)+'.png'
299     multi_fig.savefig(name)
300     print(r"\begin{subfigure}{\textwidth}\includegraphics[
        width=\textwidth]{Pictures/datas/"+name+
301 "}\caption{Confusion matrix for RS:"+str(i)+"}\end{subfigure}"
        )
302
303
304 # %% [code]
305 sns.barplot(x=reg_score["max_depth"], y=reg_score["Score"],
306             hue=reg_score["RandomState"], palette=sns.
        cubehelix_palette(n_colors=10)).set(title="Binned
        XGBRegressor Score")
307
308 # %% [code]
309 sns.barplot(x=reg_score["RandomState"], y=reg_score["Score"],
```

```
310              hue=reg_score["max_depth"],palette=sns.
         cubehelix_palette(n_colors=10)).set(title="Binned
         XGBRegressor Score")
311
312 # %% [code]
313 #GradientBoostingClassifier,RandomForestClassifier,
314 def classie(X,y,n,rs):
315      X_train, X_test, y_train, y_test = train_test_split(X, y,
         random_state=rs)
316      nn = n*50
317      model = RandomForestClassifier(n_estimators=nn)
318      model.fit(X_train,y_train)
319      prediction = model.predict(X_test)
320      multi_plot(y_test,prediction,rs,vari="RS")
321      accu = accuracy_score(y_test, prediction)
322      return accu
323
324
325
326 X = longitudinal[features]
327 #y_g = longitudinal["Group"]
328 y_g = longitudinal["CDR"].astype('str')
329
330 l=7
331 nplots = 6
332
333 classi_error = pd.DataFrame(columns=["n_estimators","
         RandomState","Score"])
334 for n in range(1,l):
335      multi_fig, multi_ax = plt.subplots(ncols=nplots,figsize=(l
         *3,2))
336      #multi_fig.tight_layout()
337      multi_fig.subplots_adjust(left=None, bottom=None, right=
         None, top=None, wspace=0.1, hspace=None)
338      errtot = 0
339      for rs in range(nplots):
340          c_score = classie(X,y_g,n,rs)
341          errtot+= c_score
342          classi_error.loc[str(n)+str(rs)]=[n*50,rs,c_score]
343      err = errtot/nplots
344      #print(err)
345      name = 'PP_dataS_MRIalz_RFCn='+str(n)+'.png'
346      multi_fig.savefig(name)
347      print(r"\begin{subfigure}{\textwidth}\includegraphics[
         width=\textwidth]{Pictures/datas/"+name+
348 r"}\caption{Confusion matrix for n{\textunderscore}estimators=
         "+str(n*50)+" with average error="+str(round(err,3))+"}\
         end{subfigure}")
349
350 # %% [code]
351 classi_error.head()
352 sns.barplot(x=classi_error["n_estimators"],y=classi_error["
         Score"],
353             hue=classi_error["RandomState"],palette=sns.
         cubehelix_palette(n_colors=10)).set(title="
         RandomForestClassifier Score")
354
```

```
355  # %% [code]
356  print(reg_score.groupby("max_depth").agg({"Score":["mean","std
         "]}))
357  print(reg_score.groupby("RandomState").agg({"Score":["mean","
         std"]}))
358  print(classi_error.groupby("n_estimators").agg({"Score":["mean
         ","std"]}))
359  print(classi_error.groupby("RandomState").agg({"Score":["mean"
         ,"std"]}))
360
361  sns.barplot(reg_score["max_depth"],reg_score["Score"],palette=
         sns.cubehelix_palette(n_colors=10)).set(title="Binned
         XGBRegressor Score")
362  plt.show()
363  sns.barplot(classi_error["n_estimators"],classi_error["Score"
         ],palette=sns.cubehelix_palette(n_colors=10)).set(title="
         RandomForestClassifier Score")
364
365
366  # %% [code]
367  sns.barplot(reg_score["RandomState"],reg_score["Score"],
         palette=sns.cubehelix_palette(n_colors=10)).set(title="
         Binned XGBRegressor Score")
368  plt.show()
369  sns.barplot(classi_error["RandomState"],classi_error["Score"],
         palette=sns.cubehelix_palette(n_colors=10)).set(title="
         RandomForestClassifier Score")
370
371
372  # %% [code]
373  v = reg_score["Score"].values
374  colors=plt.cm.plasma((v-v.min())/(v.max()-v.min()))
375  sns.barplot(reg_score["max_depth"],reg_score["Score"],palette=
         colors).set(title="Binned XGBRegressor Score")
376
377  # %% [code]
378  def checksplit(n,X,y):
379      for rs in range(5):
380          X_train, X_test, y_train, y_test = train_test_split(X,
         y, random_state=rs)
381          plt.show()
382          sns.countplot(x=y_train)
383          print(y_train.value_counts())
384
385  checksplit(5,X,y)
386
387  # %% [code]
```

## BIBLIOGRAPHY

[1] D. S. Marcus, T. H. Wang, J. Parker, J. G. Csernansky, J. C. Morris, and R. L. Buckner, "Open access series of imaging studies (oasis): Cross-sectional mri data in young, middle aged, nondemented, and demented older adults," *Journal of Cognitive Neuroscience*, vol. 19, no. 9, pp. 1498–1507, 2007.

[2] D. S. Marcus, A. F. Fotenos, J. G. Csernansky, J. C. Morris, and R. L. Buckner, "Open access series of imaging studies: Longitudinal mri data in nondemented and demented older adults," *Journal of Cognitive Neuroscience*, vol. 22, no. 12, pp. 2677–2684, 2010. PMID: 19929323.