

preschool_data_wrangling

Organize data for MRS project

Import data from old tracking sheet and reduce to kids with MRS & relevant columns

```
import numpy as np
import pandas as pd

df = pd.read_csv("/Users/meaghan/OneDrive_UCal/Preschool_data/preschool_mri_data_sheet.csv")

mrs = df.loc[df["spec"] == '1', ['study_code', 'subj_id', 'spec_id', 'date_scan', 'hand', 'female', 'mri_
```

Recode t1_quality from text to numeric (1=good, 2=medium, 3=bad) & convert date_scan to datetime format

```
mrs = mrs.replace({"good": 1, "medium": 2, "bad": 3, "good ": 1, "medium ": 2, "bad ": 3})

mrs['date_scan_dt'] = pd.to_datetime(mrs['date_scan'])
mrs['date_scan'] = mrs['date_scan_dt'].dt.date
```

Then import data from REDCap export (2019-present), rename variables to match mrs dataframe and reduce to kids with MRS & relevant columns & convert date format to datetime format

```
df2 = pd.read_csv("/Users/meaghan/OneDrive_UCal/Preschool_data/redcap_report_mri_data_oct1921.csv")

df2 = df2.rename(columns={"meta_subj_id": "subj_id", "mri_studycode": "study_code", "mri_date": "date_s

mrs2 = df2.loc[df2["spec"] == 1, ['study_code', 'subj_id', 'spec_id', 'date_scan', 'mri_age_y', 'spec',

mrs2['date_scan_dt'] = pd.to_datetime(mrs2['date_scan'], yearfirst=True)
mrs2['date_scan'] = mrs2['date_scan_dt'].dt.date
```

Merge mrs and mrs2 dataframes, clean up, and create new separate columns indicating yes/no for acquisition of ACG and LAG voxels, count total data sets per voxel and how many subjects have MRS data for each voxel

```
mrs_all = mrs.append(mrs2)

mrs_all = mrs_all.replace({"ACG": 1, "LAG": 2, "LAG, ACG": 3, "LAG and ACG": 3})

mrs_all['spec_location'].value_counts()
```

```
## 2.0    324
## 1.0    131
## 3.0      3
## Name: spec_location, dtype: int64
```

```

mrs_all.loc[(mrs_all['spec_location'] == 1) | (mrs_all['spec_location'] == 3), 'spec_acg'] = 1
mrs_all.loc[(mrs_all['spec_location'] == 2), 'spec_acg'] = 0

mrs_all.loc[(mrs_all['spec_location'] == 2) | (mrs_all['spec_location'] == 3), 'spec_lag'] = 1
mrs_all.loc[(mrs_all['spec_location'] == 1), 'spec_lag'] = 0

mrs_all['spec_acg'].value_counts()

```

```

## 0.0    324
## 1.0    134
## Name: spec_acg, dtype: int64

```

```

mrs_all['spec_lag'].value_counts()

```

```

## 1.0    327
## 0.0    131
## Name: spec_lag, dtype: int64

```

```

mrs_all.groupby('spec_acg')['subj_id'].nunique()

```

```

## spec_acg
## 0.0     111
## 1.0     130
## Name: subj_id, dtype: int64

```

```

mrs_all.groupby('spec_lag')['subj_id'].nunique()

```

```

## spec_lag
## 0.0     130
## 1.0     111
## Name: subj_id, dtype: int64

```

Save mrs_all dataframe to .csv

```

mrs_all.to_csv("/Users/meaghan/OneDrive_UCal/Preschool_data/MRS/data/mrs_data_summary_Oct2021.csv", index=False)

```

Create separate data frames for ACG and LAG data and write to .csv

```

mrs_acg = mrs_all.query('spec_acg == 1')
mrs_acg.to_csv("/Users/meaghan/OneDrive_UCal/Preschool_data/MRS/data/mrs_data_ACG_Oct2021.csv", index=False)

mrs_lag = mrs_all.query('spec_lag == 1')
mrs_lag.to_csv("/Users/meaghan/OneDrive_UCal/Preschool_data/MRS/data/mrs_data_LAG_Oct2021.csv", index=False)

```

Age range per voxel location

```

mrs_acg['mri_age_y'].min()

```

```

## 2.3381

```

```
mrs_acg['mri_age_y'].max()
```

```
## 8.0246
```

```
mrs_lag['mri_age_y'].min()
```

```
## 2.4887
```

```
mrs_lag['mri_age_y'].max()
```

```
## 10.44
```