

Traveling Subject法を用いた 複数プロトコルデータセット のHarmonize法

Multi-site data

- 多施設多疾患データプロジェクトの増加(HCP, ADNI, SRPB, CRHD, DecNef)
- コホートデータ(思春期においてcortical thicknessとageに負の相関が見られる)

メーカー、スキャナー、撮像プロトコル、ヘッドコイル、位相方向など



疾患による違い、発達による変化などは機種間差と交絡し、検出しにくくなる

Traveling Subjects

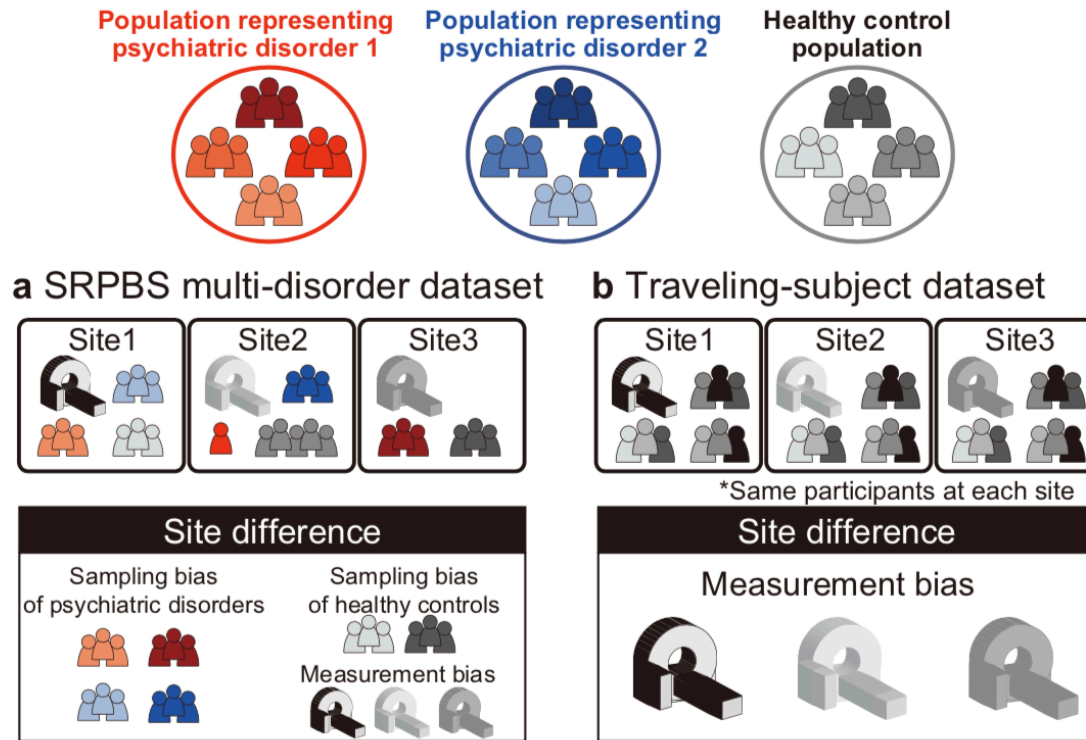
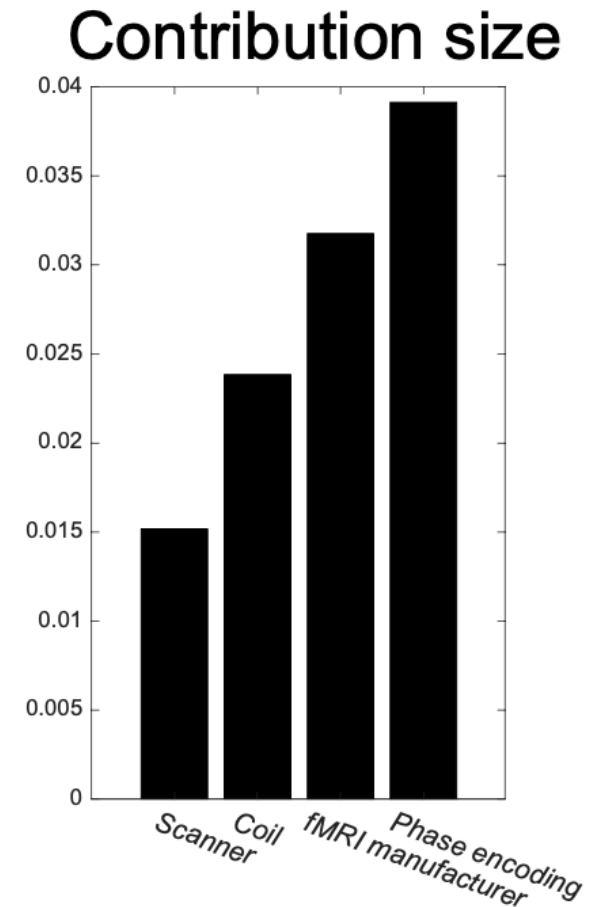


Figure 1: Schematic examples illustrating the two main datasets.

(a) The SRPBS multi-disorder dataset includes patients with psychiatric disorders and healthy controls. The number of patients and scanner types differed among sites. Thus, site differences consist of sampling bias and measurement bias. (b) The traveling-subject dataset includes only healthy controls, and the participants were the same across all sites. Thus, site differences consist of measurement bias only. SRPBS: Strategic Research Program for Brain Sciences.



GLM

I-of-K binary coding scheme

$$\text{Connectivity} = \mathbf{x}_m^T \mathbf{m} + \mathbf{x}_{s_{hc}}^T \mathbf{s}_{hc} + \mathbf{x}_{s_{mdd}}^T \mathbf{s}_{mdd} + \mathbf{x}_{s_{scz}}^T \mathbf{s}_{scz} + \mathbf{x}_d^T \mathbf{d} + \mathbf{x}_p^T \mathbf{p} + \text{const} + e,$$

$$\text{such that } \sum_j^9 p_j = 0, \sum_k^{12} m_k = 0, \sum_k^6 s_{hc k} = 0, \sum_k^3 s_{mdd k} = 0, \sum_k^3 s_{scz k} = 0, d_1(\text{HC}) = 0,$$

m : measurement bias

S_{hc}: Sampling bias for HC

S_{mdd}: Sampling bias for MDD

S_{scz} : Sampling bias for Scz

d:disorder factor

p: participant factor

Const : mean activities across sites across participants

e: $e \sim \mathbf{N}(0, \gamma^{-1})$ noise

Yamashita et al., 2019

$$\text{ROI thickness} = \mathbf{X}_m^T \mathbf{m} + \mathbf{X}_p^T \mathbf{p} + \text{const} + e$$

この式にTS01.csv, TS02.CSVを代入して計算を行います。

Data

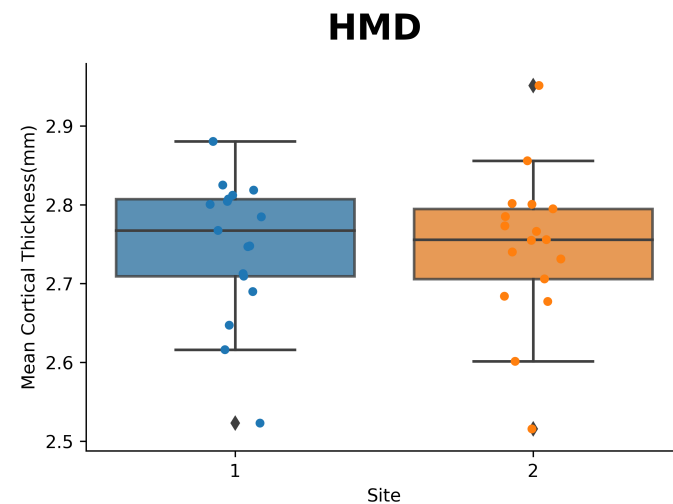
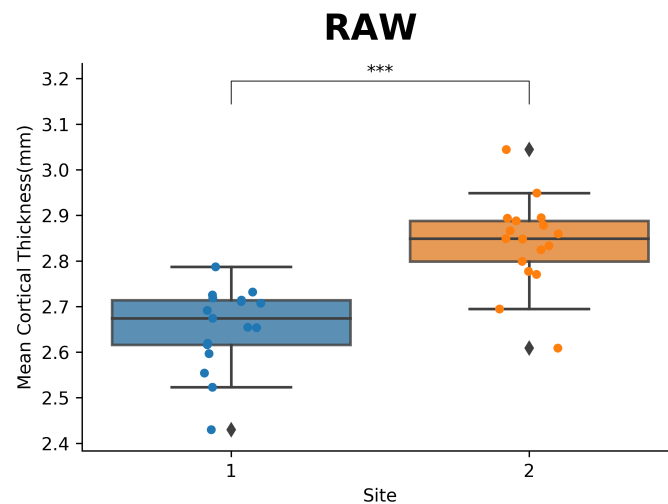
国際脳参画の思春期コホート男女40名

FreeSurfer6.0で前処理したCortical thicknessデータ (実際の数値からランダムに%増減させています)

- 1回目：11歳ごろ(Data01)

- 2回目：16歳ごろ(Data02)

	Manufacture	Scanner	Coil	protocol	TS
Site1(Data01)	Philips	Achieva	SENSE_HEAD_8	SRPB	TS01
Site2(data02)	SIEMENS	Prisma	32Ch Head	CRHD	TS02



Data

TS dataを用いて(TS01.csv, TS02.csv)
サイトごとの係数を求める

GML_HMD-ROI-Tutorial2020 .ipynb



Output files
TS_code.csv
HMD_thickness.csv
site_coef_thickAvg.csv



得られた係数を使って
データ(Data01.csv, Data02.csv)を補正する

APPLY-Tutorial2020 .ipynb



Output files
HMD_Data01.csv
HMD_Data02.csv
thick_mean_raw.png
thick_mean_hmd.png

Python code

GML_HMD-ROI-Tutorial2020 .ipynbを開いてください。

```
import sys
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.pipeline import Pipeline
from sklearn.decomposition import PCA
from sklearn.linear_model import Ridge
from sklearn.model_selection import GridSearchCV
import datetime
from sklearn import metrics as mtr
from sklearn.metrics import mean_squared_error, auc, roc_curve
import itertools
import numpy.matlib
from math import sqrt
```

モジュールを読み込み

Python code

```
# 上で作成したTSデータとコーディングデータをコピーする
ROI = TS_data.copy()
ROI=ROI.drop(['SubjID'], axis='columns')
site=TS_code.drop(['site_plot'], axis='columns')

n_samples, n_features = ROI.shape
print("%d samples, %d features" % (n_samples, n_features))
print(" ")

# 切片を求める
mean=np.mean(ROI, axis=0)
print(mean.shape)

mean_roi=np.matlib.repmat(mean,n_samples,1)
print(mean.shape,site.shape)

ROI_modified=ROI-mean_roi

# Ridge回帰を定義
clf=Ridge(alpha=0.01,fit_intercept=False)

clf.fit(site,ROI_modified)
```

Python code

```
# site effectの係数を取り出す
coef_site1 = numpy.matlib.repmat(clf.coef_[:,17],n_samples,1)
coef_site2 = numpy.matlib.repmat(clf.coef_[:,18],n_samples,1)

site1 = numpy.matlib.repmat(site['site1'],n_features,1).T
site2 = numpy.matlib.repmat(site['site2'],n_features,1).T

# site effectを求める
site1_effect=coef_site1*site1
site2_effect=coef_site2*site2

# site effectを取り除く
ROI_HMD=ROI-site1_effect-site2_effect
ROI_HMD=pd.concat([TS_data['SubjID'],ROI_HMD], axis=1)
ROI_HMD.to_csv('HMD_thickness.csv', index=False)

# site effectの係数を書き出して保存する
coef=pd.DataFrame(np.vstack((clf.coef_[:,17],
                             clf.coef_[:,18]
                             )).T,
                  columns=['site1', 'site2'])
coef.to_csv('site_coef_thickAvg.csv', index=False)
```

Apply to data

APPLY-Tutorial2020 .ipynbを開いてください

```
1 #site1#
2 #まずデータ1を補正する
3
4 data01=pd.read_csv('Data01.csv', skipinitialspace=True)
5 coef=pd.read_csv('site_coef_thickAvg.csv', skipinitialspace=True)
6
7 #SubjID列を削除し、ROIデータのためのデータフレームを作成する
8 ROI_data=data01.drop(['SubjID'],axis='columns')
9 n_samples,n_feature=ROI_data.shape #新しいデータフレームのサイズを確認
10
11 coef_site1 = nb.repmat(coef['site1'],n_samples,1)
12 print(coef_site1.shape)
13
14 #site1によるバイアスを求める
15 site1 = np.ones_like(ROI_data)
16 site1_effect=coef_site1*site1
17
18 #site1によるバイアスを取り除く
19 HMD_ROI_data=ROI_data-site1_effect
20 HMD_ROI=pd.concat([data01['SubjID'],HMD_ROI_data], axis=1)
21 HMD_ROI.to_csv('HMD_Data01.csv', index=False)
```

site2のデータについて同じ手順で行う

Visualization

