



UNIVERSITAT DE
BARCELONA



Machine Learning and Causal Inference approaches for systemic multi-disease associations in UK Biobank

Alejandro González Álvarez

BCN - AIM

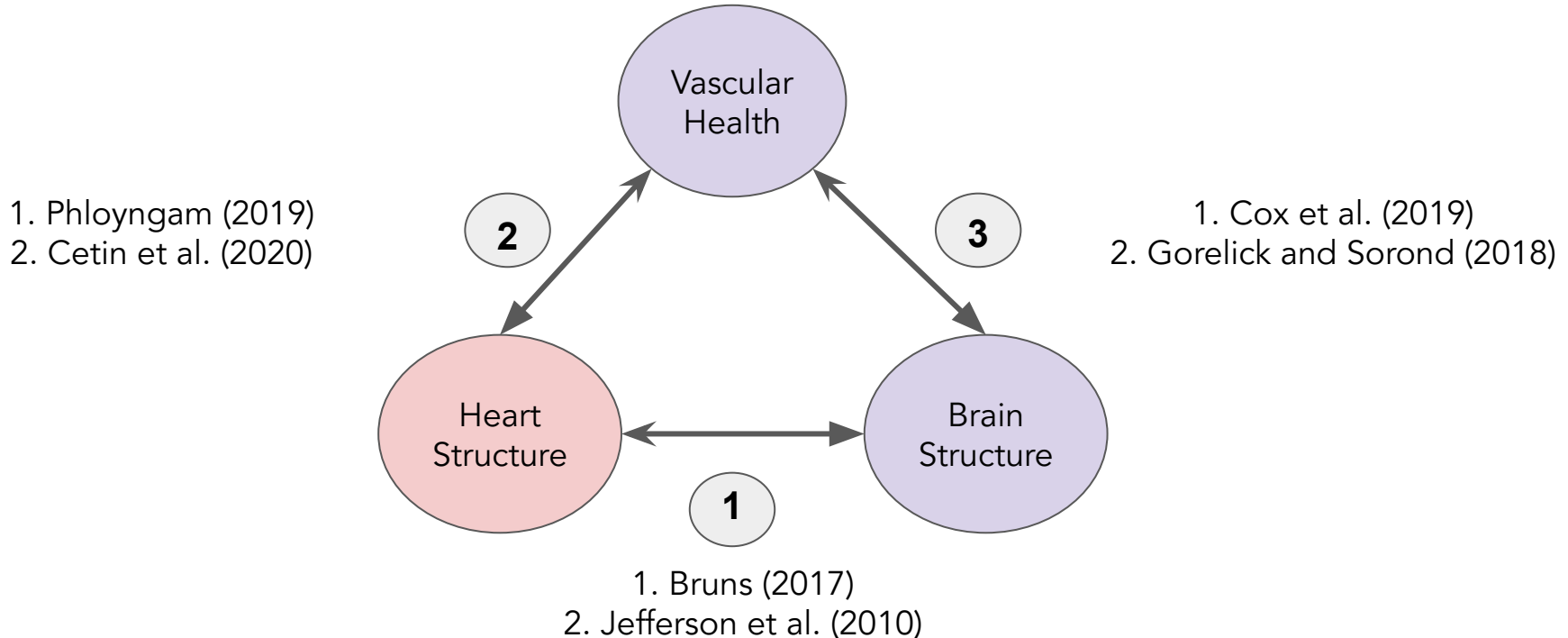
1. Problem statement

- Neurocardiology: specialty which was born with the goal of studying and understanding the pathophysiological interplay of the nervous and cardiovascular systems.
- Cardiology vs neurology. Classical diagnosis and treatment approaches of illnesses have been treated as differentiated and isolated specialties.
- Multi-organ disease association. Many observational studies have supported the clinical relevance of multi-organ disease association.
- These associations have been established largely on the basis of epidemiological data, due to insufficient knowledge on the underlying pathophysiologic mechanisms.



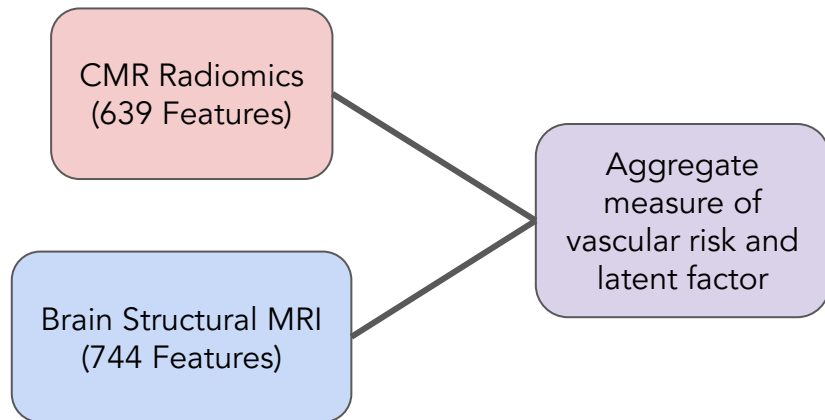
2. Related work

Plenty of studies have looked at the relationships between heart structure, brain structure and vascular health individually

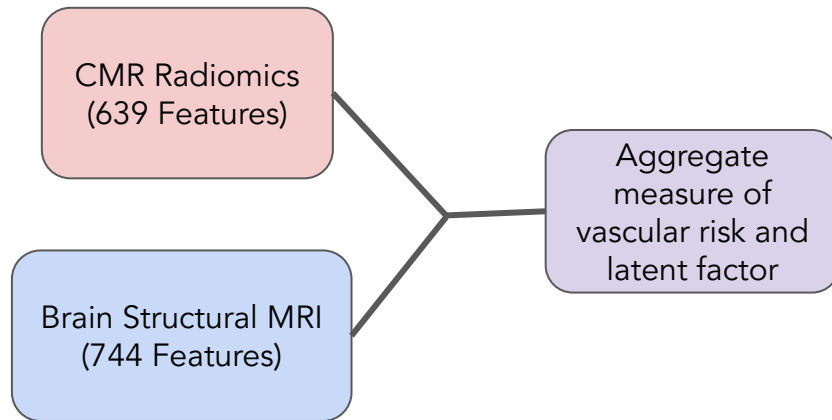


2.1 Our approach

2.1.1 Traditional ML with Separate Models

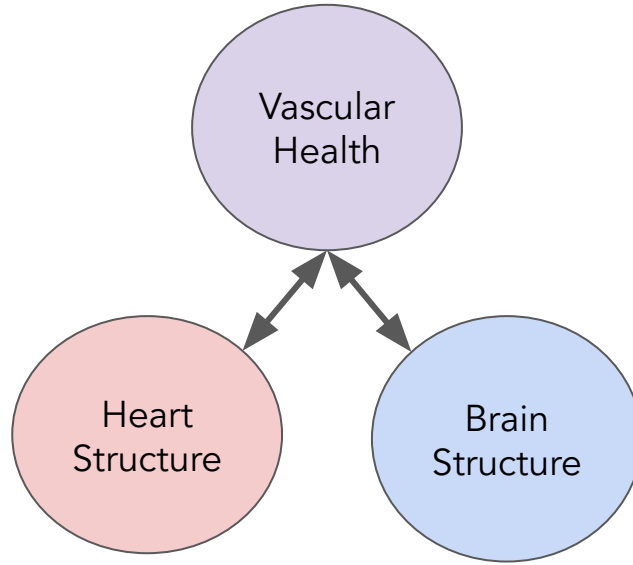


2.1.2 Traditional ML with Joint Model



- In previous work, connections and associations were studied independently. Our approach attempts to study them simultaneously.
- Initial hypothesis: brain MRI indices and heart radiomics are independent and provide unique information, therefore together they will improve performance.

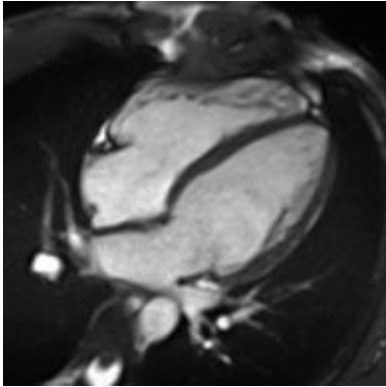
2.1.3 Causal Inference techniques



- If our initial hypothesis fails and performance does not improve, it will mean that both datasets do not provide unique information and as a result they might be similar and somehow correlated.
- Therefore, the relationship between them will be high. Causal Inference techniques pretend to find this link between them.

3. Datasets description

2065 UK Biobank Patients and 1416 variables



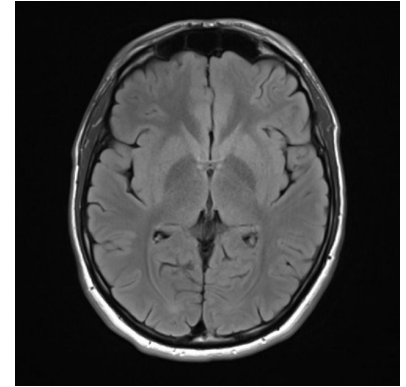
3.1 Heart CMR Radiomics

Heart imaging derived data that quantifies various changes in heart structures



3.2 VRFs

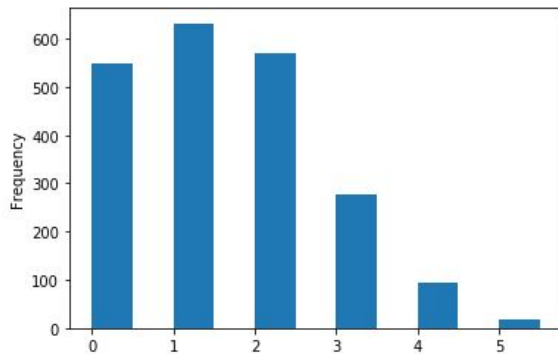
Cardiovascular risk factors that capture how well the heart works



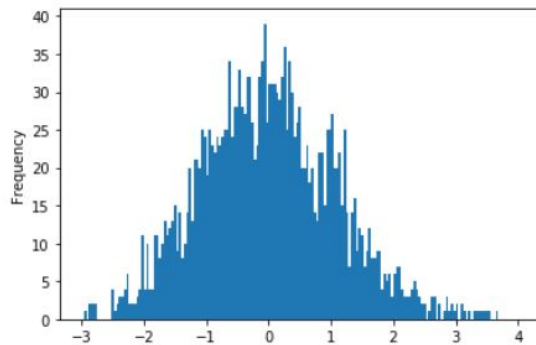
3.3 Brain MRI Indices

Brain structural imaging data that contains the structure of various brain regions

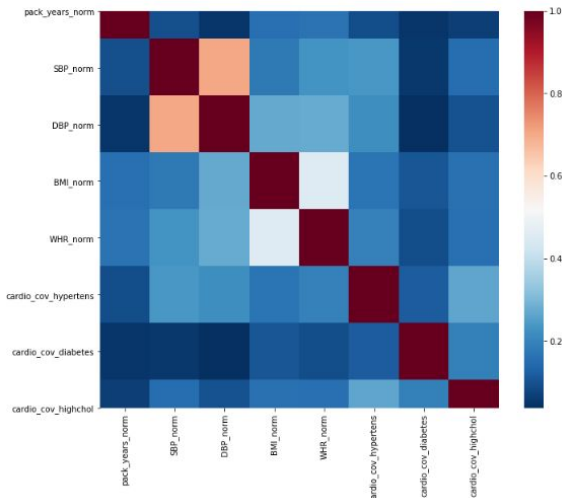
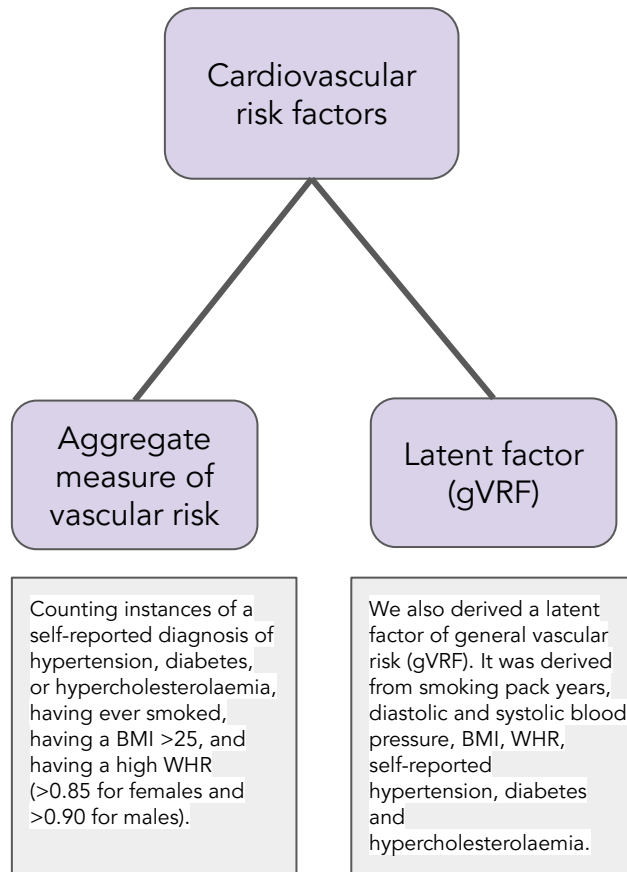
4. Dimensionality reduction



4.1 Aggregate measure



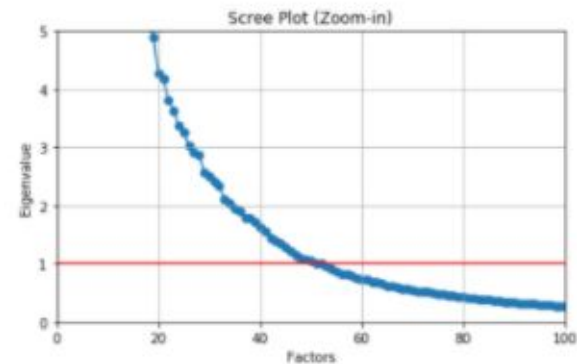
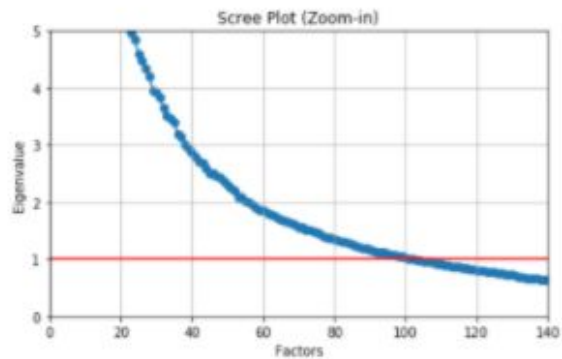
4.2 gVRF



Loadings	
DBP_norm	0.723294
SBP_norm	0.686705
WHR_norm	0.490945
BMI_norm	0.457422
cardio_cov_hypertens	0.380753
cardio_cov_highchol	0.283791
pack_years_norm	0.186169
cardio_cov_diabetes	0.155718

Brain Structural MRI

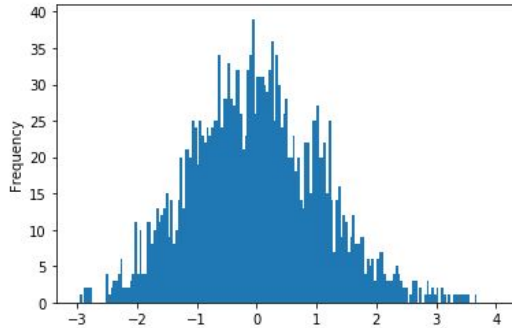
Heart CMR Radiomics



Features	# Variables	KMO Score	Bartlett's test	Scree test
gVRF	8	0,6418	(2719.71, 0)	3
Brain MRI Indices	744	0.9526	(inf, nan)	100
Heart CMR Radiomics	639	0.9781	(inf, nan)	50

5. Machine Learning approaches

5.1 Predicting gVRF

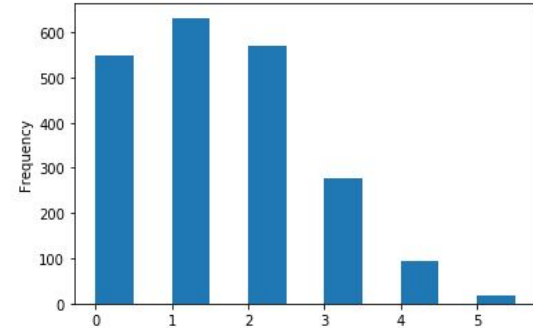


`from sklearn.linear_model import LinearRegression`

Evaluation of performance

- R²: coefficient of determination, regression score function.
- MAPE: mean absolute percentage error.
- MAE: mean of the absolute value of the errors.
- MSE: mean of the squared errors.
- RMSE: square root of the mean of the squared errors.

5.2 Predicting Aggregate measure



`from sklearn.ensemble import RandomForest`

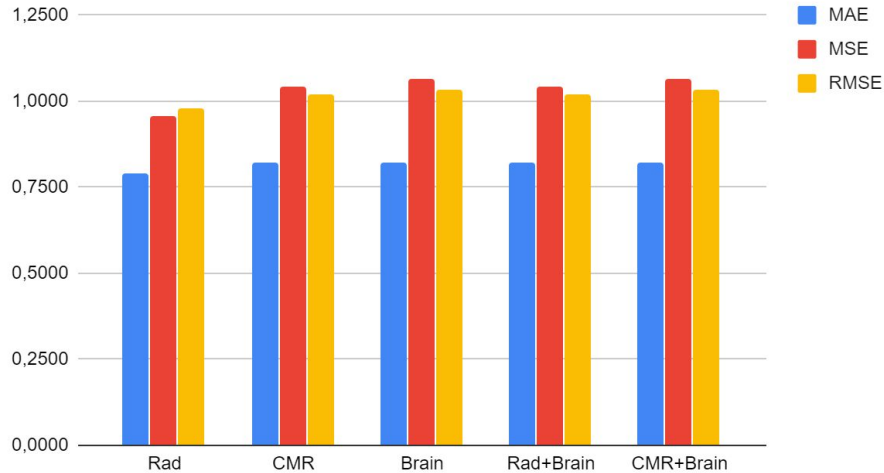
Evaluation of performance

- Accuracy.
- Confusion matrix.
- ROC curve.
- AUC.

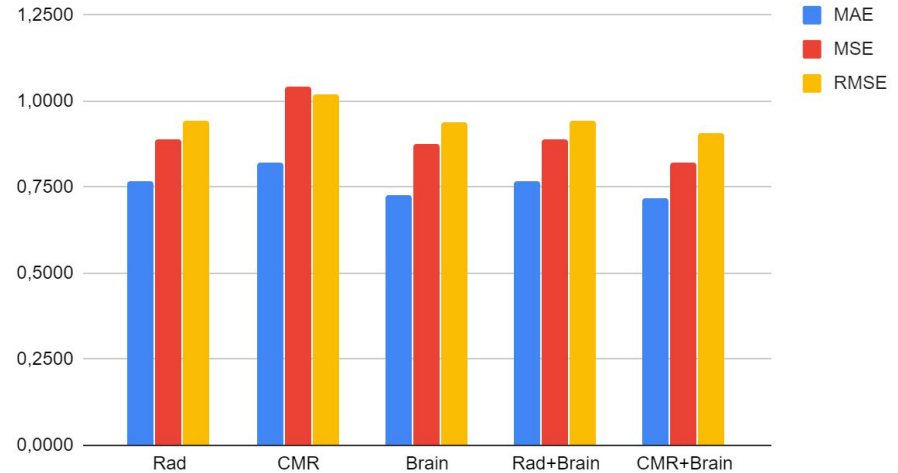
5.1 Predicting gVRF

Best results are obtained with heart radiomics in case of FA, and brain MRI indices and its combination with cardio CMR in case of SelectKbest feature selection

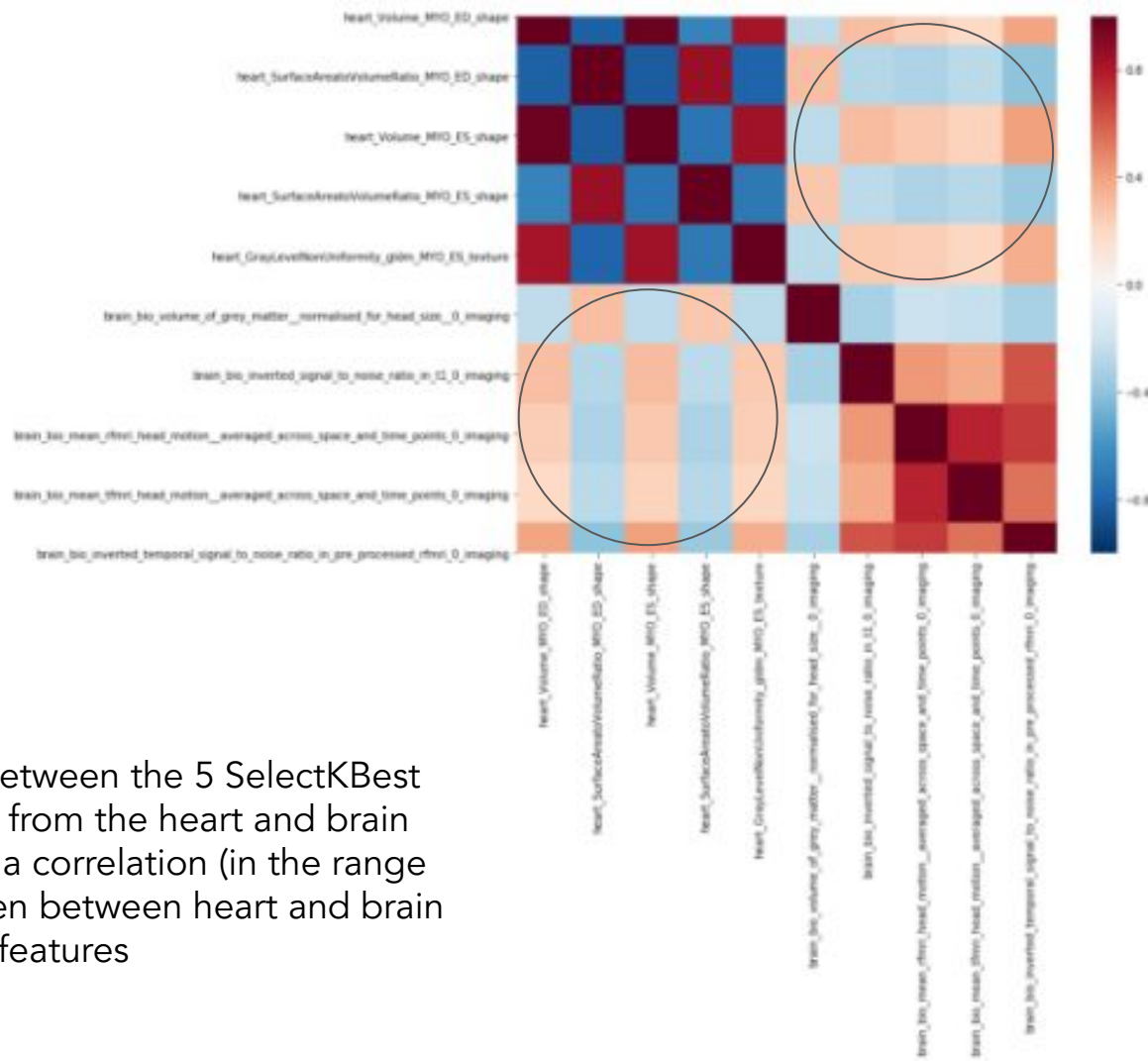
FA Errors



Kbest Errors



However, overall, the combination of these datasets does not seem to improve at all, or improve very little our performance metrics, which may lead us to reject our initial hypothesis

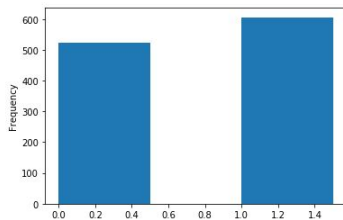


5.2 Predicting aggregate measure of vascular risk

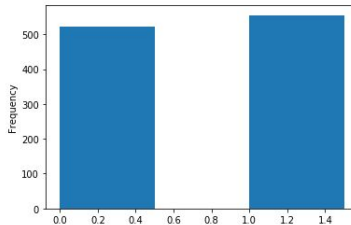
- Multiclass classification problem ❌
- Binarize classes:
 - Random oversampling
 - Random undersampling
 - SMOTE-Tomek



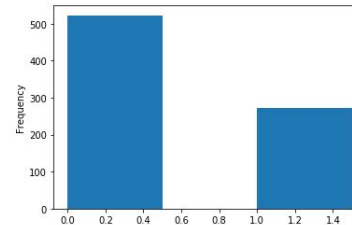
0vs1



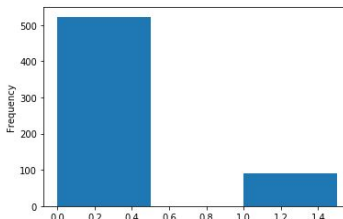
0vs2



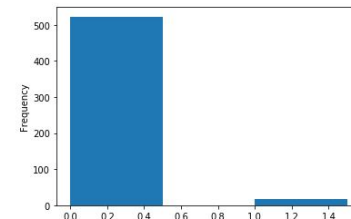
0vs3



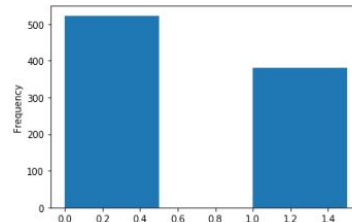
0vs4



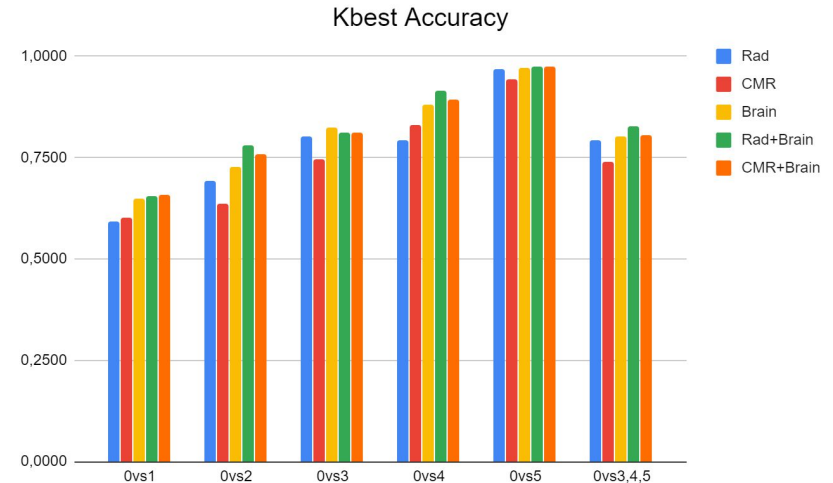
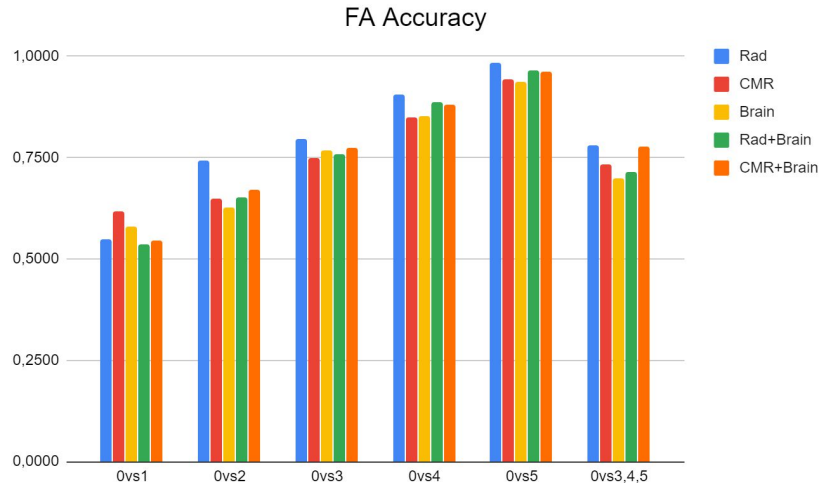
0vs5



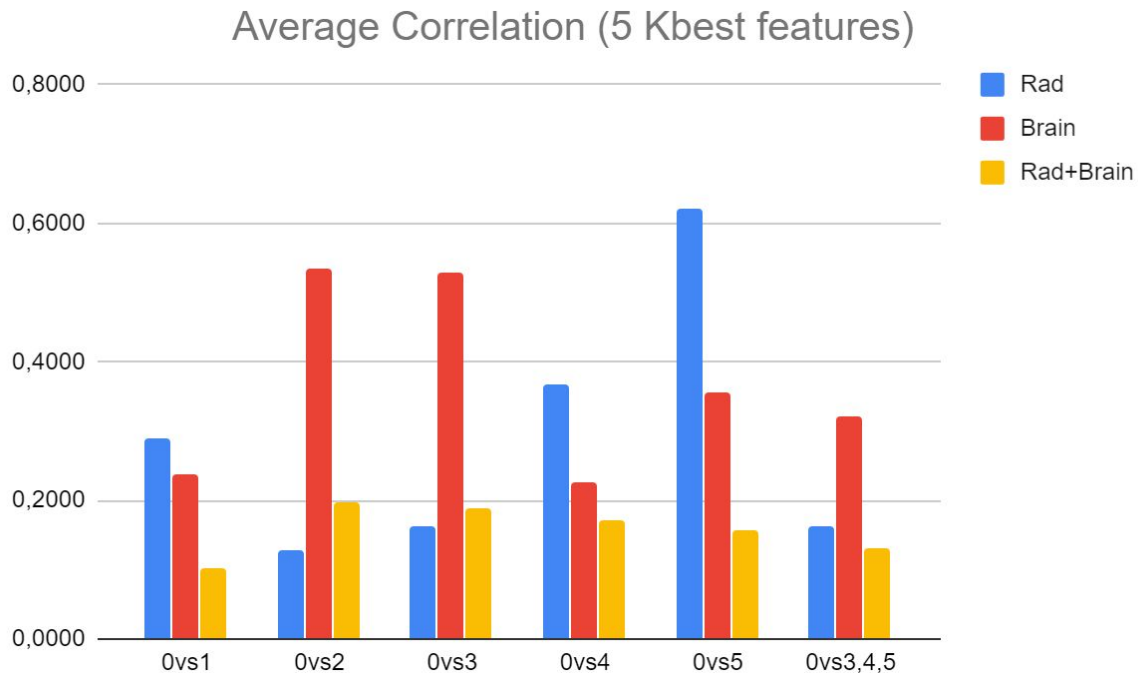
0vs3,4,5



1. The main trend is that model accuracy and AUC improves as VRFs burden increases.
2. Combined effects of risk factors seem to be better detected by:
 - Heart CMR radiomics when using FA as the dimensionality reduction technique.
 - Brain MRI indices when using SelectKBest as the dimensionality reduction technique.
3. Again, as it happened when predicting gVRF, the combination of heart and brain datasets does not seem to improve at all, or improve very little our performance metrics, which may lead us to reject our initial hypothesis.



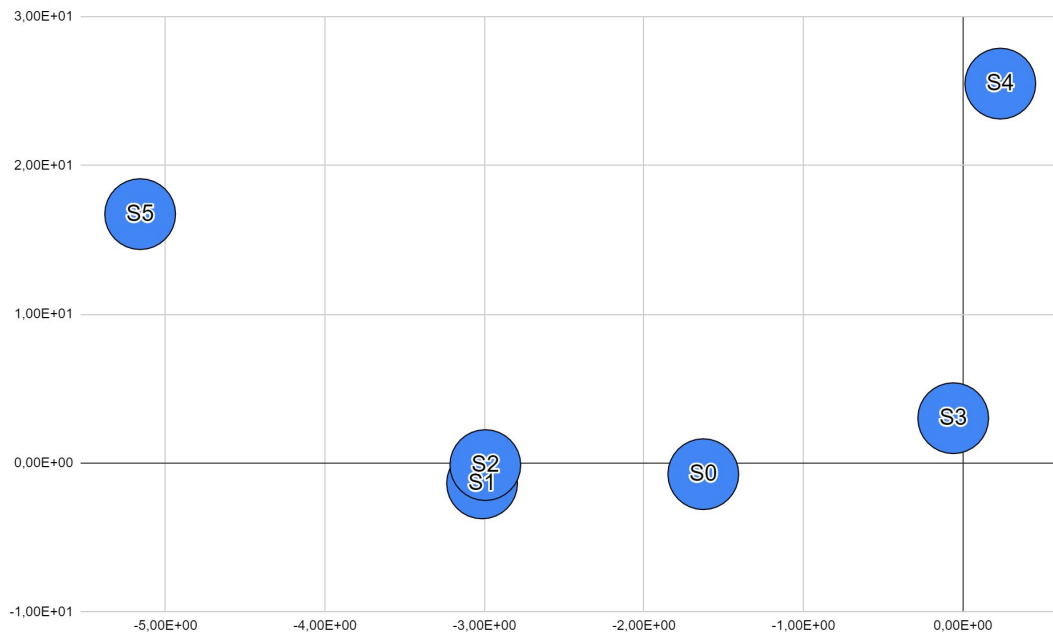
The average correlation between the top 5 SelectKBest features from each dataset seems to be very low and not significant, so we cannot infer yet that these datasets do not provide unique information



5.3 K-means clustering

- The multiclass classification task did not work.
- Model accuracy and AUC improved as VRFs burden increased.
- These results reinforce the hypothesis mentioned earlier that different classes in our aggregate measure might be increasingly different as the number of vascular risk factors increases.
- Patients with zero VRFs, the most healthy ones, and patients with five VRFs, the most at risk, seemed to be the most differentiated ones.
- To measure the distance between classes we ran the K-means clustering algorithm and computed the centroids for each class.

- S5, patients with an aggregate score of five, and S0, patients with an aggregate score of zero, are the most distant ones, reinforcing our hypothesis that these two classes differ the most.
- S1, S2 and S3 are the closest aggregate measures to S0, showing why these comparisons obtained the lowest performance metrics and why they are the most similar classes from our target variable.

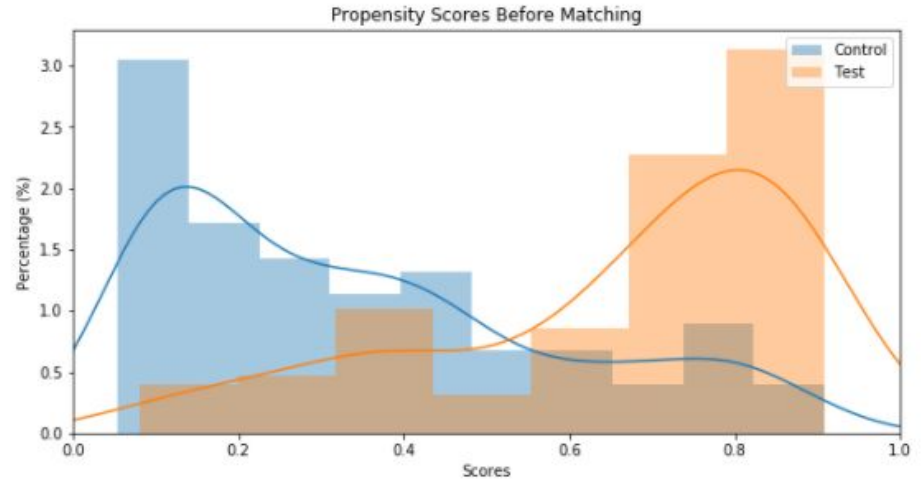


6. Propensity score matching analysis

Individuals with a total of 4 and 5 VRFs were each matched on sex and age to a single participant who had no VRFs

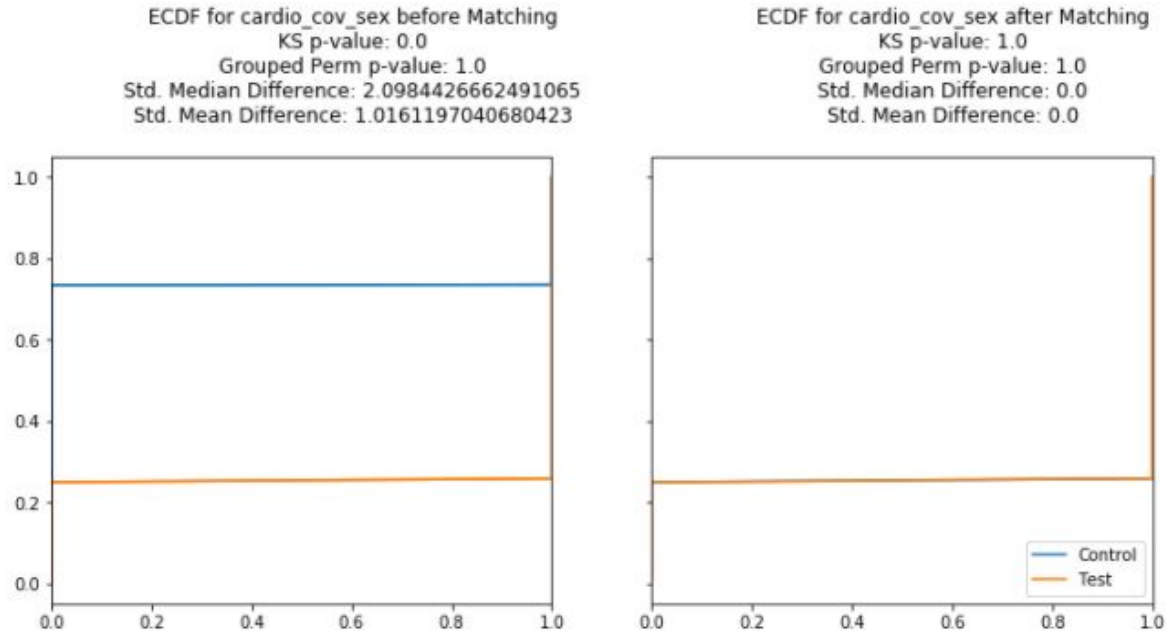
Fitting Models on Balanced Samples: 100\100
Average Accuracy: 75.9%

The average accuracy of our 100 models is 75.9%, suggesting that there's separability within our data and justifying the need for the matching procedure

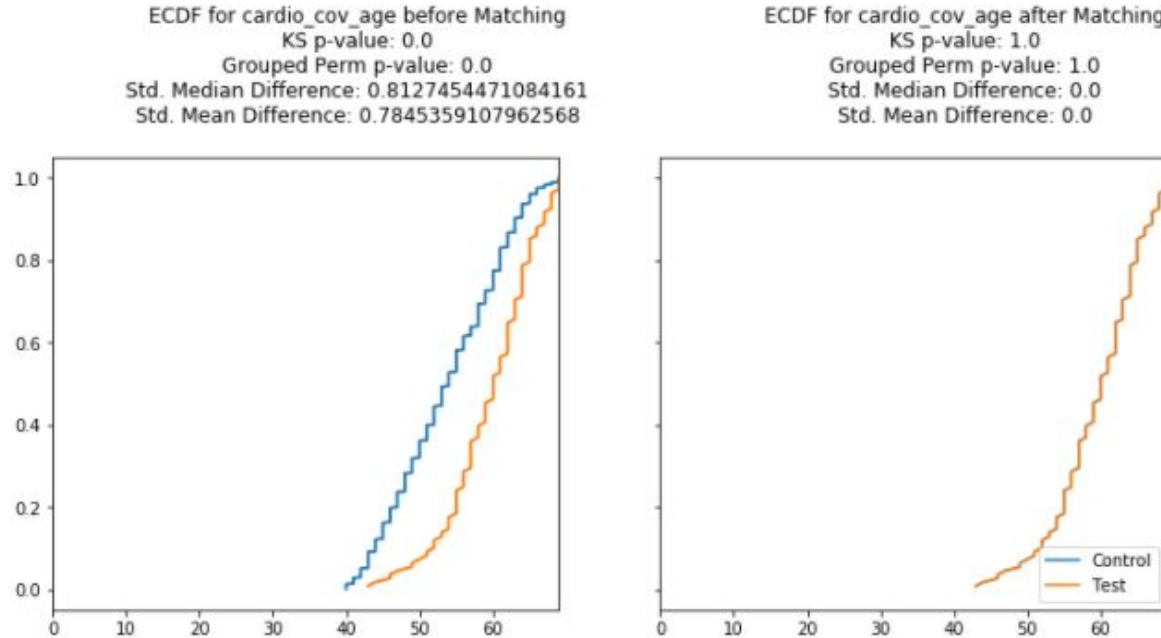


The plot above demonstrates the separability present in our data. Test profiles (patients with four and five VRFs) have a much higher propensity, or estimated probability of defaulting than the control group (patients with zero VRFs), given the features we isolated in the data

- Kolmogorov-Smirnov Goodness of fit Test (KS-test). This test statistic is calculated on 1000 permuted samples of the data, generating an empirical p-value.
- Chi-Square Distance. Similarly this distance metric is calculated on 1000 permuted samples.
- Standardized mean and median differences

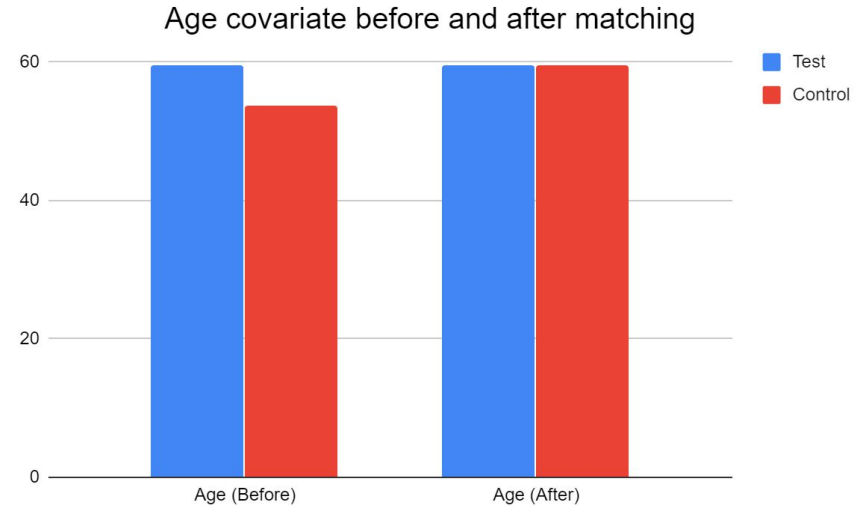
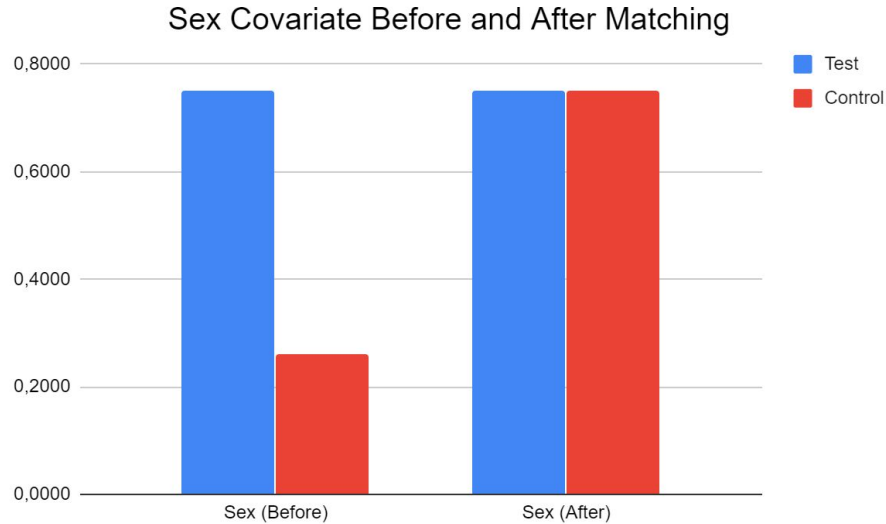


Both lines are very close to each other for both of our covariates (sex and age),
and even indistinguishable after matching



p values from both the KS-test and the grouped permutation of the Chi-Square distance after matching are above 0.05, meaning they were statistically significant

Mean distributions for age and sex before and after matching for both our control and test groups



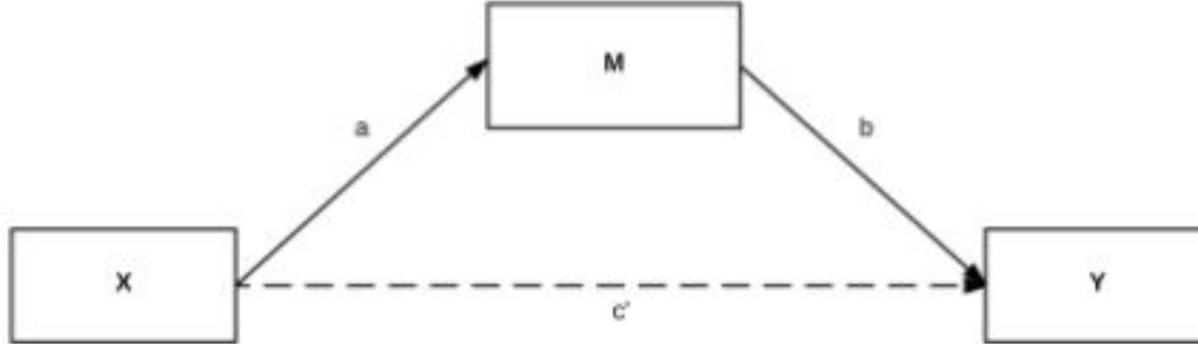
- Before matching, there are more males than females in the test group (patients with four and five VRFs).
- The age mean for the test group (59.58) before matching is higher than the control group (53.79).

*** Summary so far...**

1. Combination of brain and heart structure' datasets did not improve when trying to predict vascular health.
2. We may reject our initial hypothesis, since these datasets may NOT provide unique information, and somehow they are related.
3. Combined effects of risk factors were better predicted by heart CMR radiomics when using FA as the dimensionality reduction technique.
4. Combined effects of risk factors were better detected by brain MRI indices when using SelectKBest as the dimensionality reduction technique.
5. Direct link between VRFs and brain structure?
6. Direct link between VRFs and heart structure?
7. An intermediate factor is playing a mediating role between these associations?

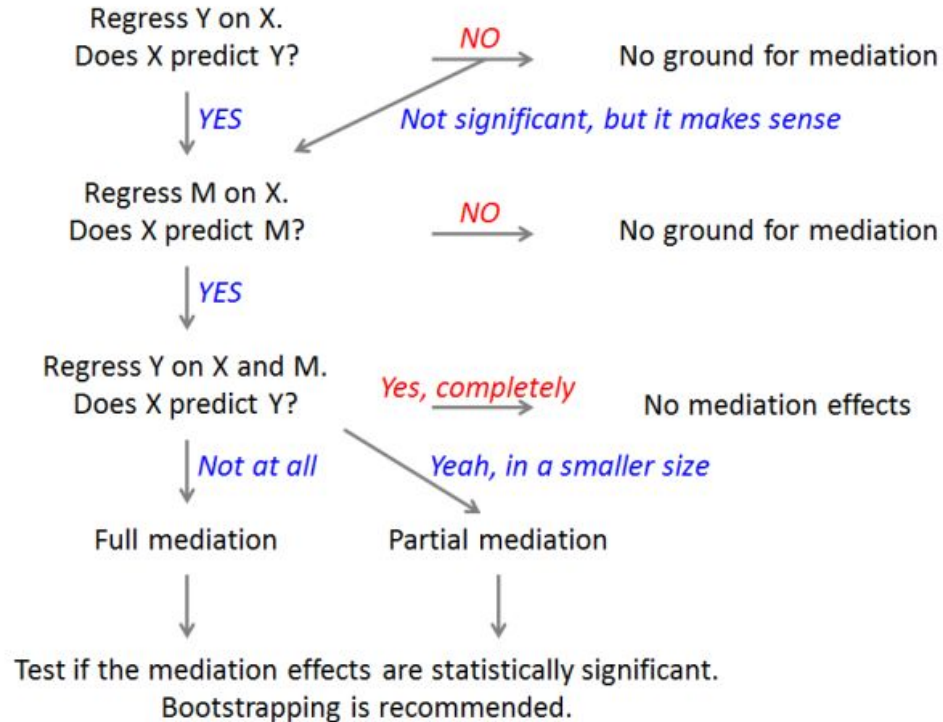
7. Causal Inference techniques (Mediation Analysis)

- a and b reflect the indirect path of the effect of X on the outcome (Y) through the mediator (M).
- c' is the direct effect of X on the outcome after the indirect path has been removed.
- The total effect of X is the combined indirect and direct effects.

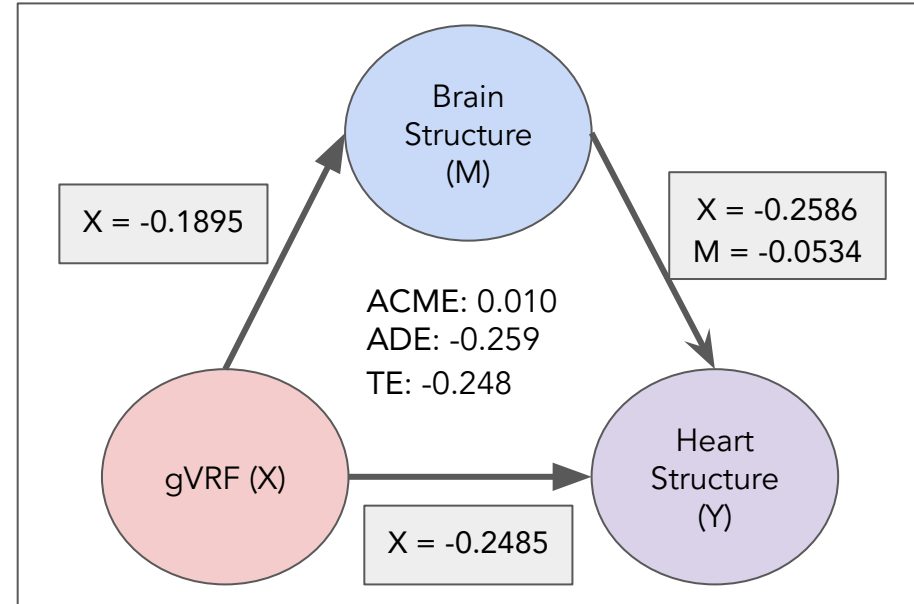
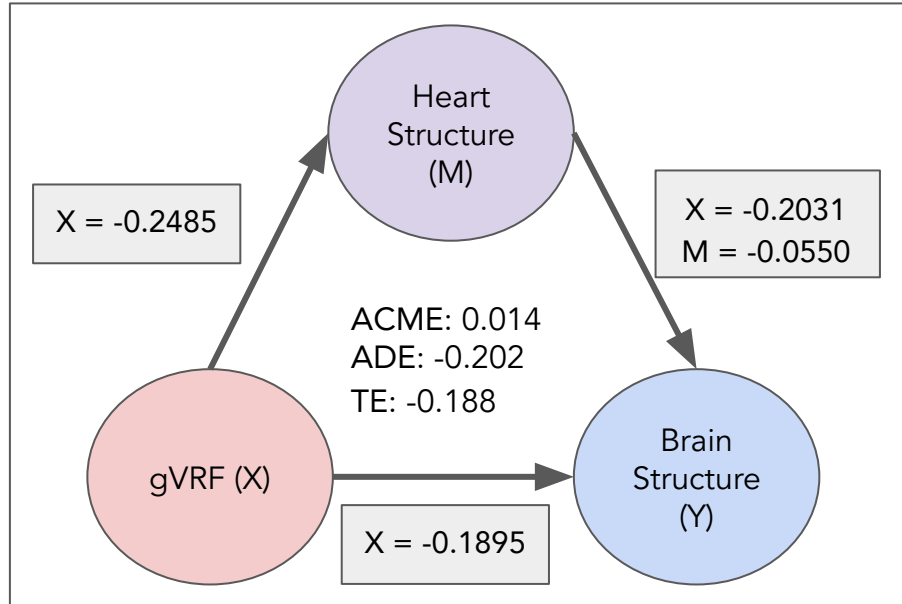


- ACME (average causal mediation effect): Total effect minus the direct effect (TE - ADE).
- ADE (average direct effect) : A direct effect of X on Y after taking into account a mediation indirect effect of M ($X + M \rightarrow Y$).
- TE (total effect) (indirect + direct effect): A total effect of X on Y (without M) ($X \rightarrow Y$).

A mediation analysis is comprised of three sets of regressions:
 $X \rightarrow Y$, $X \rightarrow M$, and $X + M \rightarrow Y$

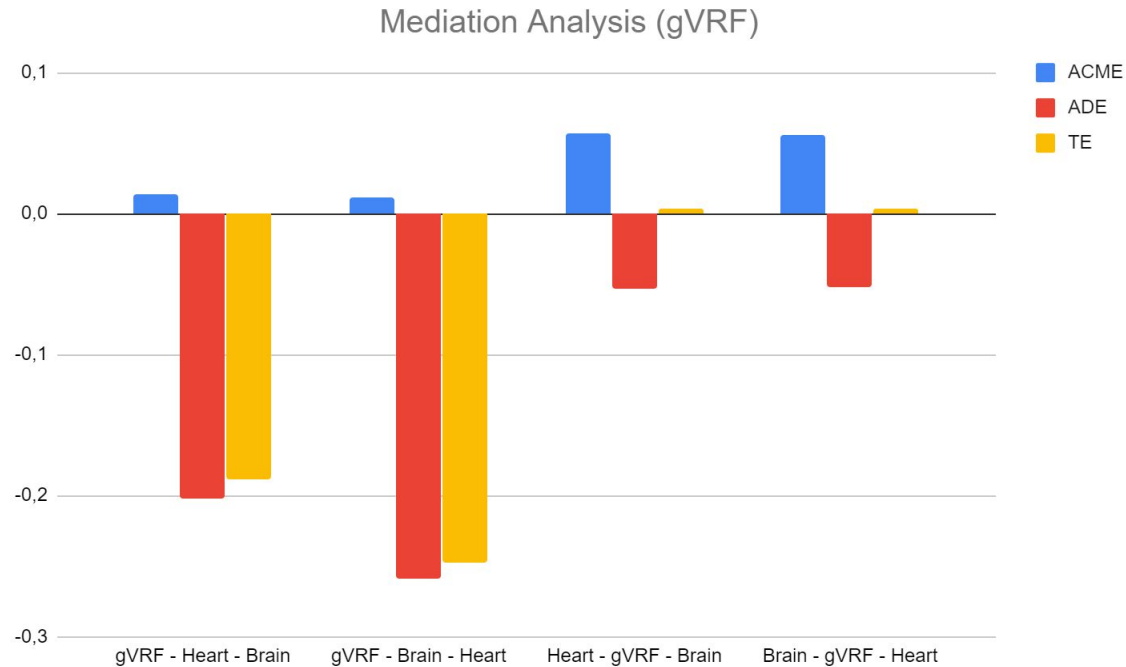


How much of the connection between cardiovascular risk and heart/brain structure can be explained by changes in brain/heart structure?



- VRF → Heart: strongly negatively correlated
- VRF → Brain: strongly negatively correlated
- Heart → Brain: weakly negatively correlated

- Strongest (-) direct effect: VRF → Heart
- Small mediation role of heart and brain (small but significant relationship)



What these factors and negative correlations mean?

TABLE 4.3: Top 10 Heart CMR Radiomics variables' loadings

Variables	Loadings
Heart Inverse Difference glcm RV ES texture	0.891137
Heart Inverse Difference Moment glcm RV ES texture	0.889092
Heart Inverse Difference glcm RV ED texture	0.880894
Heart Inverse Difference Moment glcm RV ED texture	0.880765
Heart Gray Level Non Uniformity Normalized glrlm RV ES texture	0.871713
Heart Large Dependence Low Gray Level Emphasis gldm RV ED texture	0.867649
Heart Gray Level Non Uniformity Normalized glrlm LV ES texture	0.866858
Heart Gray Level Non Uniformity Normalized glrlm RV ED texture	0.858650
Heart Inverse Difference glcm LV ES texture	0.858577
Heart Large Dependence Low Gray Level Emphasis gldm LV ED texture	0.856550

TABLE 4.4: Top 10 Brain MRI Indices variables' loadings

Variables	Loadings
Brain mean l3 in anterior corona radiata on fa skeleton left	0.875992
Brain mean l3 in anterior corona radiata on fa skeleton right	0.870006
Brain mean md in anterior corona radiata on fa skeleton left	0.865722
Brain weighted mean l3 in tract inferior fronto occipital fasciculus right	0.859523
Brain mean md in superior longitudinal fasciculus on fa skeleton left	0.858544
Brain mean l2 in anterior corona radiata on fa skeleton left	0.857454
Brain mean md in anterior corona radiata on fa skeleton right	0.856520
Brain mean l3 in superior corona radiata on fa skeleton left	0.853679
Brain mean md in superior longitudinal fasciculus on fa skeleton right	0.852527
Brain mean l3 in superior longitudinal fasciculus on fa skeleton right	0.851144

8. Discussion and Future Work

1. Combination of features did not improve performance, but we only used a subset of features from each dataset (5 SelectKBest features and 5 latent factors) for comparison purposes.
2. Correlation was not very high between these subsets of features (maybe with more features involved correlation increases).
3. Classes in aggregate measure too similar between them. (Better way to measure vascular health?)
4. Age and sex important role (further analysis needed).
5. VRF \rightarrow Heart: strongly negatively correlated.
6. VRF \rightarrow Brain: strongly negatively correlated.
7. Heart \longleftrightarrow Brain: weakly negatively correlated.
8. Small mediation role of the heart and the brain (which might increase with more features involved).