

Similarity Data based TBI Patients Stratification and Feature Selection

¹Emily Li, ¹Qi Cheng, ¹Alan D. Kaplan, ²Geoffrey Manley

¹Lawrence Livermore National Laboratory

²TRACK-TBI Consortium

May 15, 2019

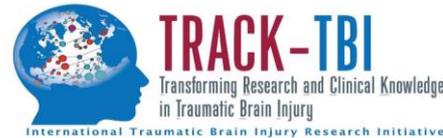


Introduction

- Traumatic brain injury (TBI) is a complex disorder. Though TBI is a major cause of death and disability, it is traditionally stratified based on clinical signs and symptoms, and few targeted treatments exist.
- The Transforming Research and Clinical Knowledge in Traumatic Brain Injury (TRACK-TBI) Pilot multicenter study enrolled 586 acute TBI patients and collected diverse common data elements (TBI-CDEs) across the study population, including imaging, genetics, and clinical outcomes.
 - The most highly-detailed ever collected data in neuroscience, with enormous opportunity for data-rich phenotyping and TBI precision medicine
- The goals of the project include developing better ways to diagnose and provide early prognosis of brain injury and to characterize various subtypes of brain injury that might have different underlying pathophysiology, prognosis, and optimal treatment.
- We conduct multidimensional exploratory analytics and use similarity data based approaches to reveal data-driven patterns in patient outcomes, and identify the key features that characterize these patterns.
 - Provide for diagnostic decision support and precision medicine in neurotrauma

Components of the TRACK-TBI Dataset and Data-Analysis Challenges

- **Data captured across 4 broad domains:**
- **Clinical assessments and demographic information,**
- **Blood biomarkers (genetics and proteomics)**
- **Neuroimaging (Head CT and MRI)**
- **Outcome measures at 3, 6, 12 months**



Missing values

Heterogeneity in data

- Mixed categorical, ordinal, continuous

Multimodality

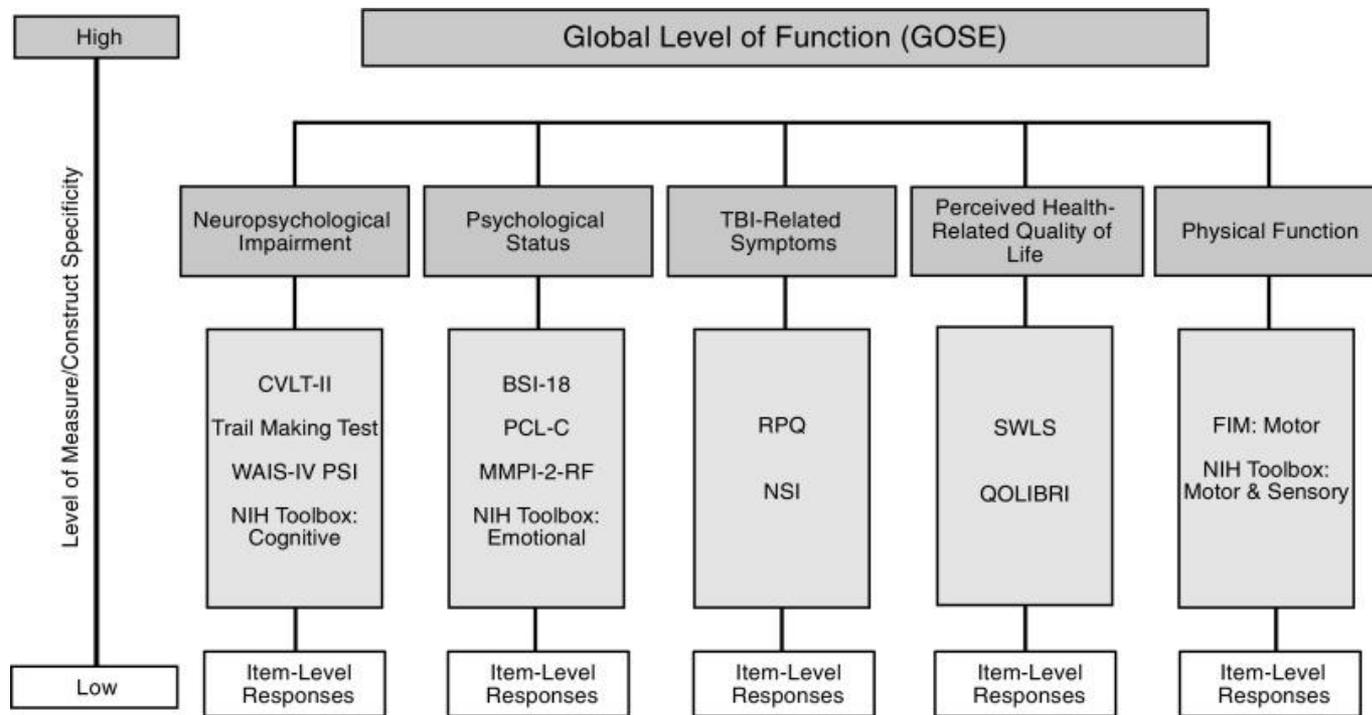
- Scalar and various imaging types

Multidimensional targets (outcomes)

Data Structure

- Primarily working with the 6-month outcome data
- 6-month outcome data breakdown by functional groups

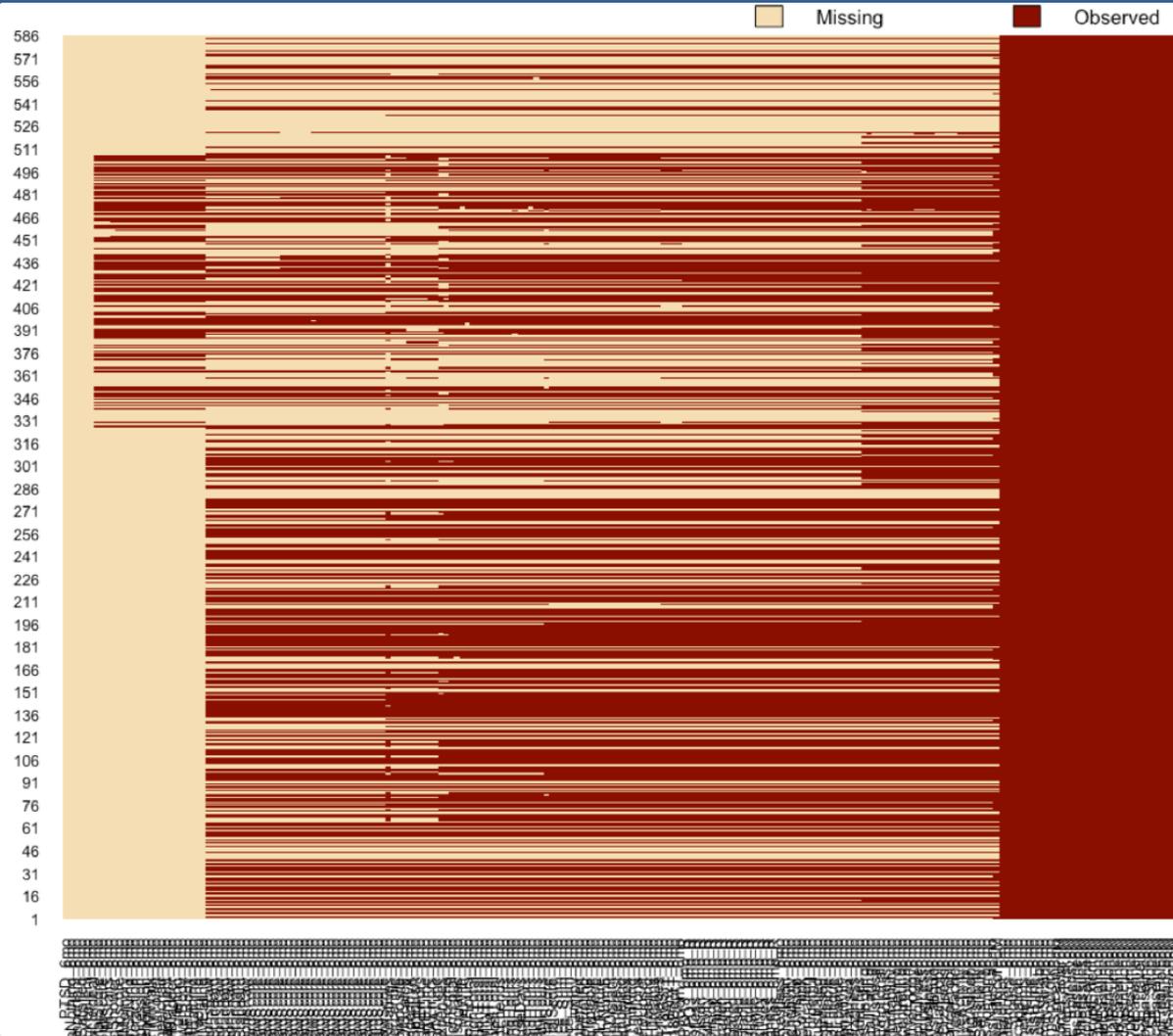
Data Structure



Nelson LD, Ranson J, Ferguson AR, Giacino J, Okonkwo DO, Valadka A, Manley G, McCrea M, (2017). Validating Multidimensional Outcome Assessment Using the TBI Common Data Elements: An Analysis of the TRACK-TBI Pilot Sample." J Neurotrauma. doi: 10.1089/neu.2017.5139.

Tbi.6m all r x c

- Data cleaning
 - Delete rows and columns with most/all values missing
 - Delete variables with constant values
 - Convert all categorical variables to numerical variables
 - Convert factors level into ordinal numerical variables if level relationship can be understood
 - Convert factors level to binary codes (with increased dimensions) if level relationship cannot be understood



Original GOSE Categorical Variables

	GOSE_Reponder6M	GOSE_SimpleCommand6M	GOSE_NeedAssistance6M	GOSE_NeedFreqHelp6M	GOSE_AssistanceBaseline6M
	: 1	: 1	: 1	:230	:230
Patient alone	:231	No : 0	No :229	2 : 1	No : 5
Patient plus relative	: 3	Yes:234	Yes: 5	Yes: 4	Yes: 0
RelativeΣ≤_/friendΣ≤_/caretaker:	0				
	GOSE_Shop6M	GOSE_ShopBaseline6M	GOSE_Travel6M	GOSE_TravelBaseline6M	GOSE_WorkResume6M
	: 2	: 19	: 2	: 18	: 2
No :	3	No : 0	No : 4	No : 0	No : 60
Yes:	230	Yes:216	Yes:229	Yes:217	Yes:173
	GOSE_WorkRestriction6M	GOSE_WorkBaseline6M	GOSE_SocialResume6M	GOSE_SocialRestriction6M	
	:180	:175	: 1	:166	
Reduced work capacity	: 38	No : 16	No : 70	Participate a bit less: 22	
Sheltered workΣ≤_/unable to work:	17	Yes: 44	Yes:164	Participate much less : 29	
				Unable to participate : 18	
	GOSE_SocialBaseline6M	GOSE_DisruptRelation6M	GOSE_DisruptExtent6M	GOSE_DisruptBaseline6M	GOSE_OtherIssues6M
	:165	: 2	:155	:157	: 1
No :	6	No :151	Constant : 19	No : 63	No : 83
Yes: 64	Yes: 82	Frequent : 22	Occasional: 39	Yes: 15	Yes:151
	GOSE_OtherIssuesBaseline6M	GOSE_Epilepsy6M	GOSE_EpilepsyRisk6M	GOSE_OutcomeFactor6M	
	: 20	: 1	: 1	: 1	
No :	204	No :224	No :197	Both : 81	
Yes: 11	Yes: 10	Yes: 37	Extracranial: 10	11	
			TBI :143	10	
				01	



Matrix Completion

- to impute missing values

- Reconstruct a low rank matrix from a subset of its entries
 - The low rank assumption is applicable here due to the redundancy among variables
- Let M be the original incomplete matrix and X be the estimated low rank complete matrix

- Define $(M^E)_{ij} = \begin{cases} M_{ij} & \text{if } M_{ij} \text{ is known} \\ 0 & \text{otherwise} \end{cases}$ and $(X^E)_{ij} = \begin{cases} X_{ij} & \text{if } M_{ij} \text{ is known} \\ 0 & \text{otherwise} \end{cases}$

- Then, $\min \|M^E - X^E\|_F^2$ s.t. $\text{rank}(X) \leq r$

- Finally,

$$\hat{M}_{ij} = \begin{cases} M_{ij} & \text{if } M_{ij} \text{ is known} \\ X_{ij} & \text{otherwise} \end{cases}$$

Method of Non-negative Matrix Factorization (NMF) for Similarity Data

– for data clustering

- Similarity matrix $S_{N \times N}$

$$S_{ij} = e^{-\frac{\|x_i - x_j\|^2}{\sigma^2}}$$

- Probability matrix $P = S / 1_N^t S 1_N$
- Factorization $P = WH^t$, where

$$p_{ij} = P(x_i, x_j) = \sum_{v=1}^K w_{iv} h_{jv}$$

$$w_{iv} = P(x_i, c_v)$$

$$h_{jv} = P(x_j | c_v)$$

- The objective is to

$$\min_{W, H} C(P || WH^t) := - \sum_{i,j} [p_{ij} \log \sum_v (w_{iv} h_{jv})]$$

subject to $0 \leq w_{iv}, h_{jv} \leq 1, \forall i, j, v, \sum_{i,v} w_{iv} = 1, \sum_j h_{jv} = 1$.

- The label of x_i is

$$l_{x_i} = \arg \max_v w_{iv}$$

Model Selection

- to determine K

- The cross entropy keeps on decreasing with increased K
- Stability measure

$$S = \mathbb{E} \left[\min_{\pi \in \Theta_K} \frac{1}{n} \sum_{i=1}^n \mathbf{1} \{ \pi(h_{\mathbf{X}, \alpha_{\mathbf{X}}(\mathbf{X})}(\mathbf{X}_i')) \neq \alpha_{\mathbf{X}'}(\mathbf{X}_i') \} \right]$$

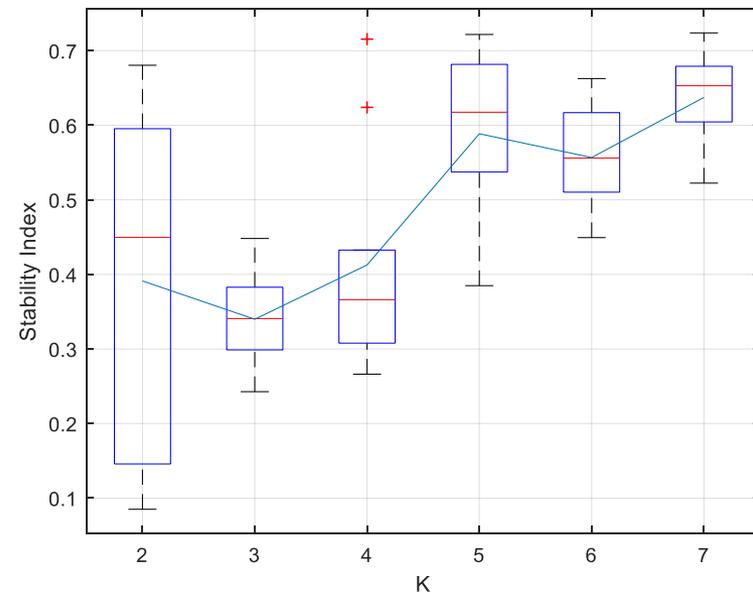
where $\alpha_{\mathbf{X}'}(\mathbf{X}_i')$ is the set of cluster labels for data \mathbf{X}_i' using model order K and clustering method α , and $h_{\mathbf{X}, \alpha_{\mathbf{X}}(\mathbf{X})}(\mathbf{X}_i')$ is the set of predicted labels of data \mathbf{X}_i' using \mathbf{X} , $\alpha_{\mathbf{X}}(\mathbf{X})$ as the training data. Here, \mathbf{X} and \mathbf{X}_i' are from the same distribution. Θ_k is the set of all permutations of $\{1, \dots, K\}$. On average, the smaller S , the more stable the model is.

TRACK-TBI Pilot Outcome Variables - 6 months: Summary Scores

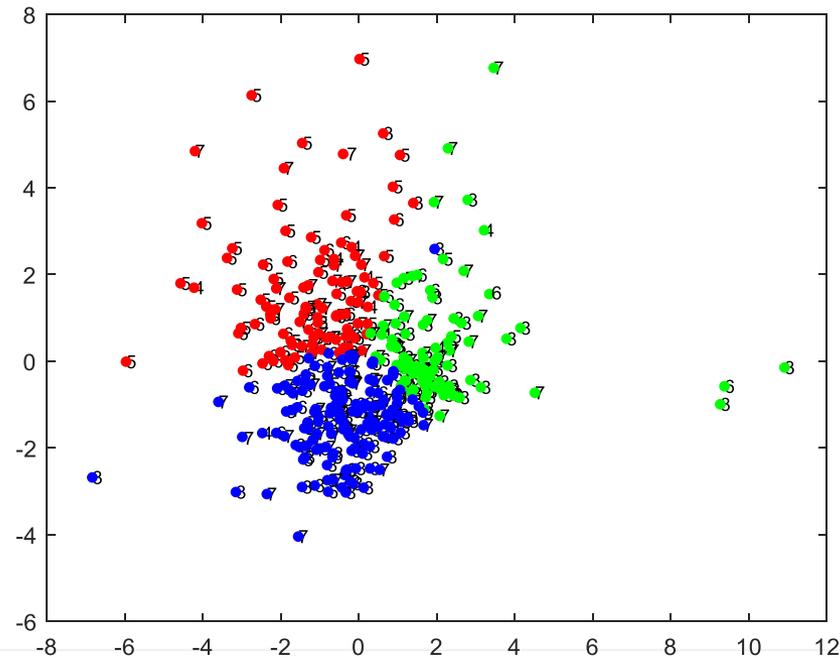
CDE Outcome Domain	Name of Measure	Variable Name	Coding
Global Outcome	Glasgow Outcome Scale Extended (GOSE)	GOSE_OverallScore	(1=death, 2=vegetative state, 3/4=lower/upper severe, 5/6=lower/ upper moderate, 7/8=lower/upper good recovery)
Post-concussive/TBI-related symptoms	Rivermead Post Concussion Symptoms Questionnaire (RPQ)	RPQ_Total	0 - 64 (higher = more injury-related symptoms)
Psychiatric and Psychological Status	Brief Symptom Inventory 18-item (BSI-18)*	BSI18DeprT	T scores with population M = 50, SD = 10; higher=more symptoms
		BSI18AnxT	T scores with population M = 50, SD = 10; higher=more symptoms
		BSI18SomT	T scores with population M = 50, SD = 10; higher=more symptoms
	PTSD Checklist--Civilian (PCL-C)	PCLTotalScore	17-84 (higher = more posttraumatic stress disorder symptoms)
Generic Quality of Life	Satisfaction With Life Scale (SWLS)	SWLSTotalScore	1-35 (higher=more satisfied with life)
Neuropsychological Impairment	Trail Making Test	TMTPartATime	# seconds (higher = slower psychomotor speed)
		TMTPartBTime	# seconds (higher = slower speed/mental flexibility/set-shifting)
	Wechsler Adult Intelligence Scale (WAIS-IV) Processing Speed Index (PSI)	WAIS_PSI_Composite	Standard score (population M = 100, SD = 15); higher = better/faster processing speed
	California Verbal Learning Test - Second Edition (CVLT-II)	CVLTTrial1To5RawScore	Higher = better memory
		CVLTTrialBRawScore	Higher = better memory
		CVLTShortDelayFreeRecall	Higher = better memory
		CVLTLongDelayFreeRecall	Higher = better memory
		CVLTTotalRecognitionDiscriminability	Higher = better memory
		CVLTTotalIntrusionsRaw	Higher = more intrusion errors (worse performance)
CVLTTotalRepetitionsRaw		Higher = more repetitive responses (worse performance)	

Analysis of Summary Scores: model selection

- Stability assessment boxplot showing the distribution of 10 repeats of stability value on y-axis for each model across different model on x-axis.

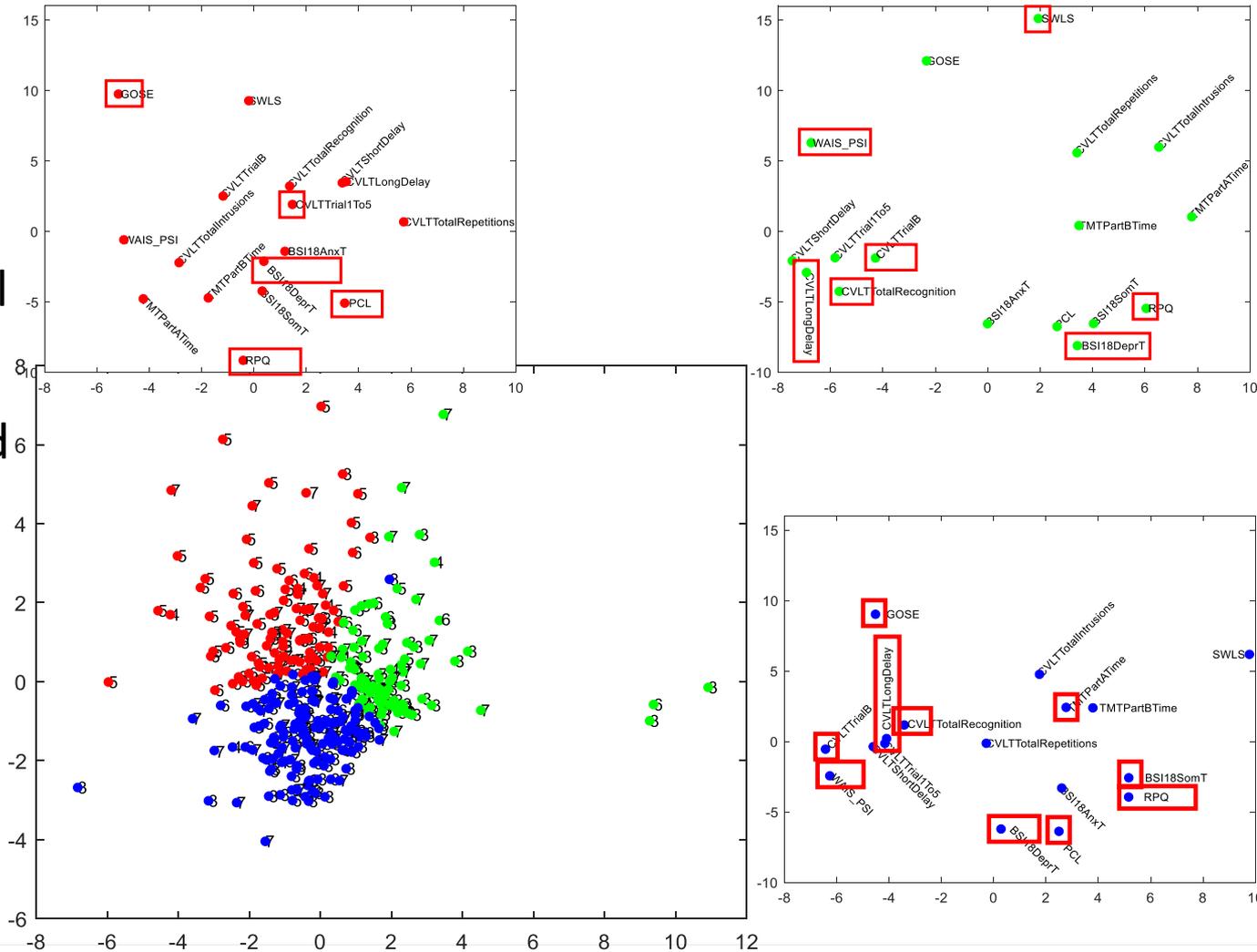


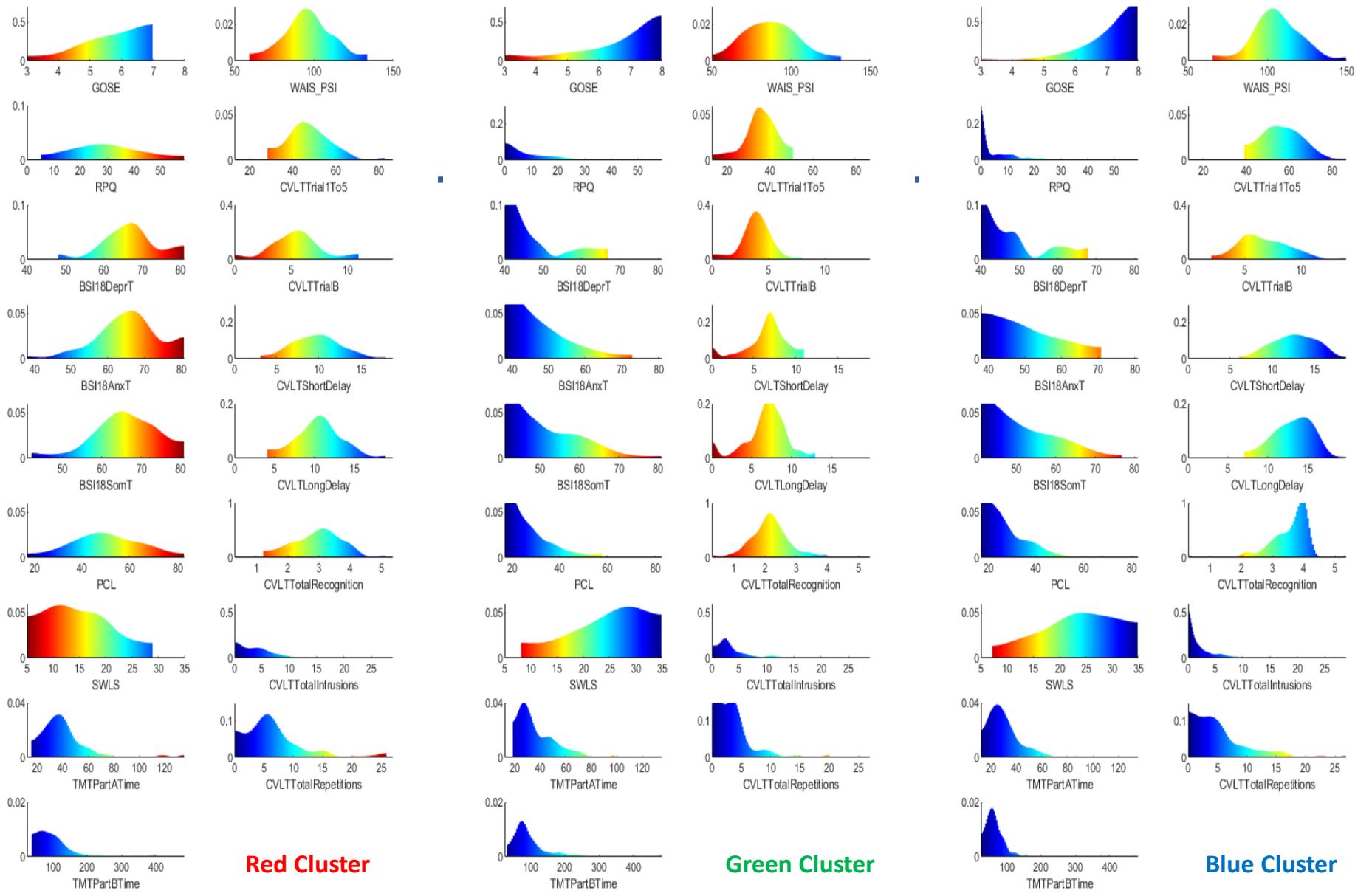
Analysis of Summary Scores: clustering for K=3



Analysis of Summary Scores: clustering for K=3

- To understand and characterize each cluster of patients, a binomial logistic regression model is built based on the assigned cluster labels and the step-wise feature selection procedure is conducted to select the contributing features in distinguishing these clusters.

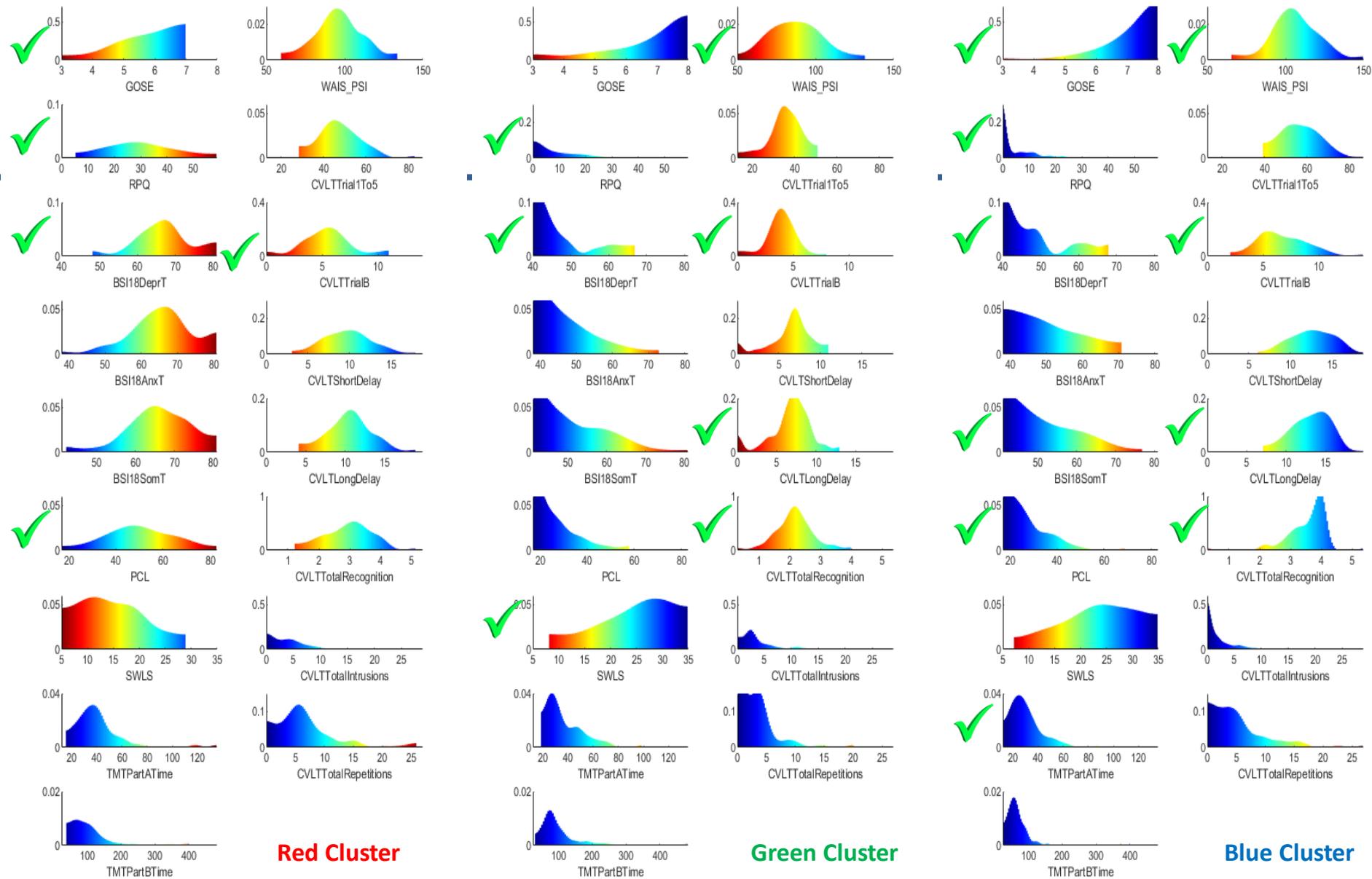




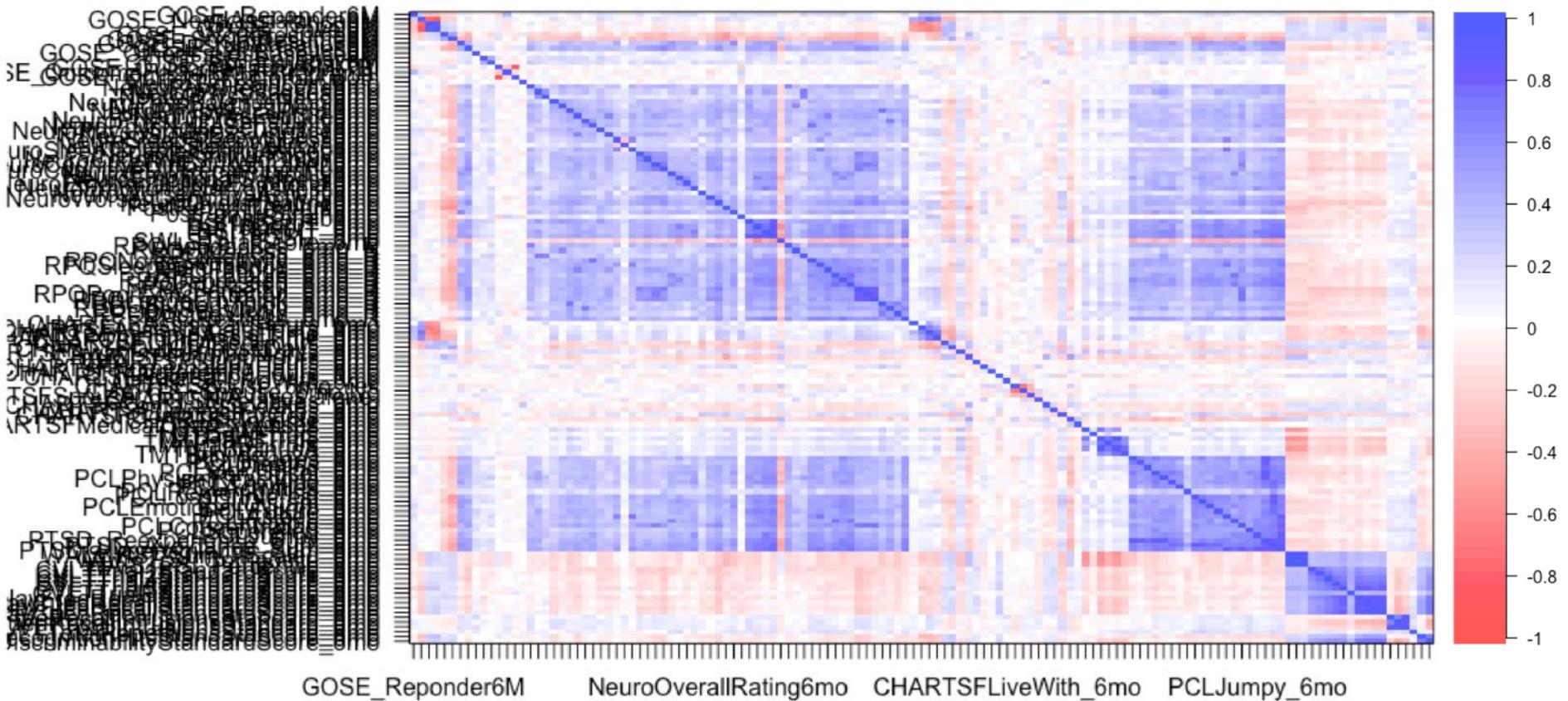
Red Cluster

Green Cluster

Blue Cluster



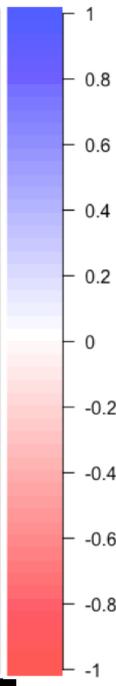
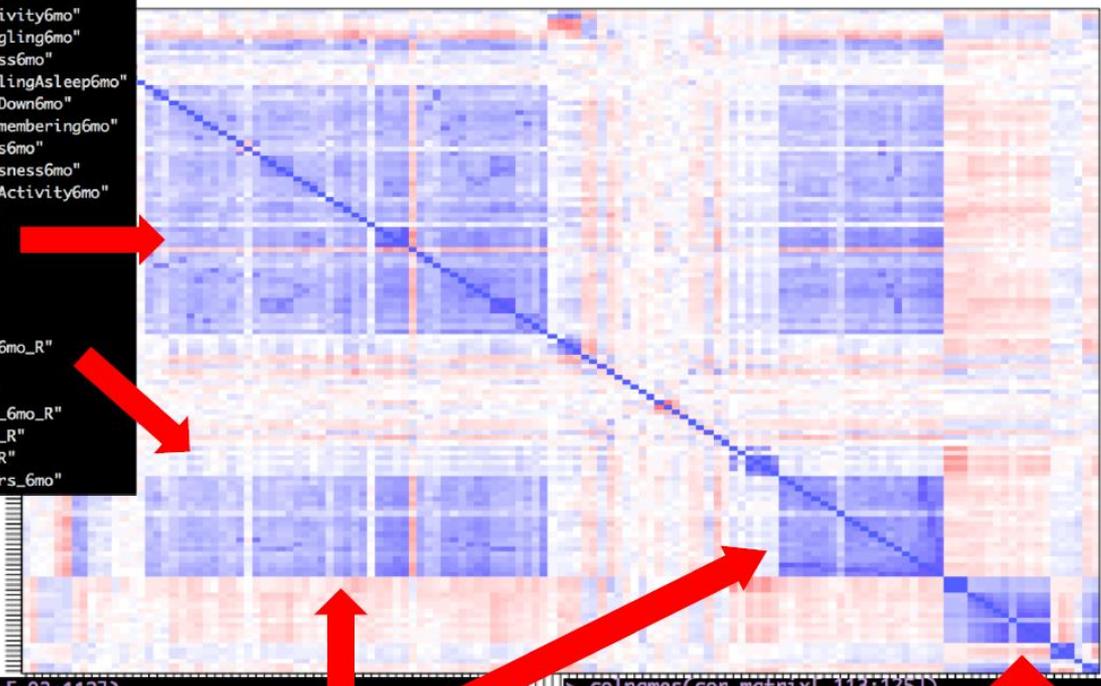
Tbi.6m (132 variables)– correlation plot



Correlation plot

```
[1] "NeuroPhysHeadache6mo"
[3] "NeuroPhysVomiting6mo"
[5] "NeuroPhysDizziness6mo"
[7] "NeuroPhysFatigue6mo"
[9] "NeuroPhysNoiseSensitivity6mo"
[11] "NeuroSleepDrowsiness6mo"
[13] "NeuroSleepSleepingMore6mo"
[15] "NeuroCognitiveFoggy6mo"
[17] "NeuroCognitiveDiffConcentrating6mo"
[19] "NeuroEmotionalIrritability6mo"
[21] "NeuroEmotionalMoreEmotional6mo"
[23] "NeuroWorsenPhysActivity6mo"
[25] "NeuroOverallRating6mo"
[27] "PostFamilyStrain6mo"
[29] "BSI18SomT_6mo"
[31] "BSI18AnxT_6mo"
[33] "SWLSTotalScore_6mo"
[35] "RPQDizziness_6mo_R"
[37] "RPQNoiseSensitivity_6mo_R"
[39] "RPQFatigue_6mo_R"
[41] "RPQDepressed_6mo_R"
[43] "RPQForgetful_6mo_R"
[45] "RPQLongerToThink_6mo_R"
[47] "RPQLightSensitivity_6mo_R"
[49] "RPQRestless_6mo_R"
```

```
"NeuroPhysNausea6mo"
"NeuroPhysBalanceProbl6mo"
"NeuroPhysVisualProbl6mo"
"NeuroPhysLightSensitivity6mo"
"NeuroPhysNumbnesTingling6mo"
"NeuroSleepSleepingLess6mo"
"NeuroSleepTroubleFallingAsleep6mo"
"NeuroCognitiveSlowedDown6mo"
"NeuroCognitiveDiffRemembering6mo"
"NeuroEmotionalSadness6mo"
"NeuroEmotionalNervousness6mo"
"NeuroWorsenCognitiveActivity6mo"
"PostReturnToWork6mo"
"PostRehab6mo"
"BSI18DeprT_6mo"
"BSI18GSIT_6mo"
"RPQHeadaches_6mo_R"
"RPQNausea_6mo_R"
"RPQSleepDisturbance_6mo_R"
"RPQIrritable_6mo_R"
"RPQFrustrated_6mo_R"
"RPQPoorConcentration_6mo_R"
"RPQBlurredVision_6mo_R"
"RPQDoubleVision_6mo_R"
"CHARTSFAssistPaidHours_6mo"
```



```
> colnames(cor.matrix[,93:112])
[1] "PCLImages_6mo"
[3] "PCLFeeling_6mo"
[5] "PCLPhysicalReactions_6mo"
[7] "PCLActivities_6mo"
[9] "PCLLossOfInterest_6mo"
[11] "PCLEmotionallyNumb_6mo"
[13] "PCLAsleep_6mo"
[15] "PCLConcentrating_6mo"
[17] "PCLJumpy_6mo"
[19] "PTSD_Avoidance_Sum_6mo"
"PCLDreams_6mo"
"PCLVeryUpset_6mo"
"PCLThinking_6mo"
"PCLRemembering_6mo"
"PCLDistant_6mo"
"PCLFuture_6mo"
"PCLIrritable_6mo"
"PCLSuperAlert_6mo"
"PTSD_Reexperience_Sum_6mo"
"PTSD_Hypervigilance_Sum_6mo"
```

```
> colnames(cor.matrix[,113:125])
[1] "WAIS_PSI_SumOfScaled_6mo"
[2] "WAIS_PSI_Composite_6mo"
[3] "WAIS_PSI_Percentile_6mo"
[4] "CVLTTrial1StandardScore_6mo"
[5] "CVLTTrial2StandardScore_6mo"
[6] "CVLTTrial3StandardScore_6mo"
[7] "CVLTTrial4StandardScore_6mo"
[8] "CVLTTrial5StandardScore_6mo"
[9] "CVLTTrial8StandardScore_6mo"
[10] "CVLTShortDelayFreeRecallStandardScore_6mo"
[11] "CVLTShortDelayCuedRecallStandardScore_6mo"
```

Data Fusion

- Recall, the objective is to

$$\min_{W,H} C(P||WH^t) := - \sum_{i,j} [p_{ij} \log \sum_v (w_{iv} h_{jv})]$$

subject to $0 \leq w_{iv}, h_{jv} \leq 1, \forall i, j, v, \sum_{i,v} w_{iv} = 1, \sum_j h_{jv} = 1$.

- Now $P = \sum_{l=1}^L \alpha_l P_l, \sum_{l=1}^L \alpha_l = 1, \alpha_l \geq 0, \forall l$.

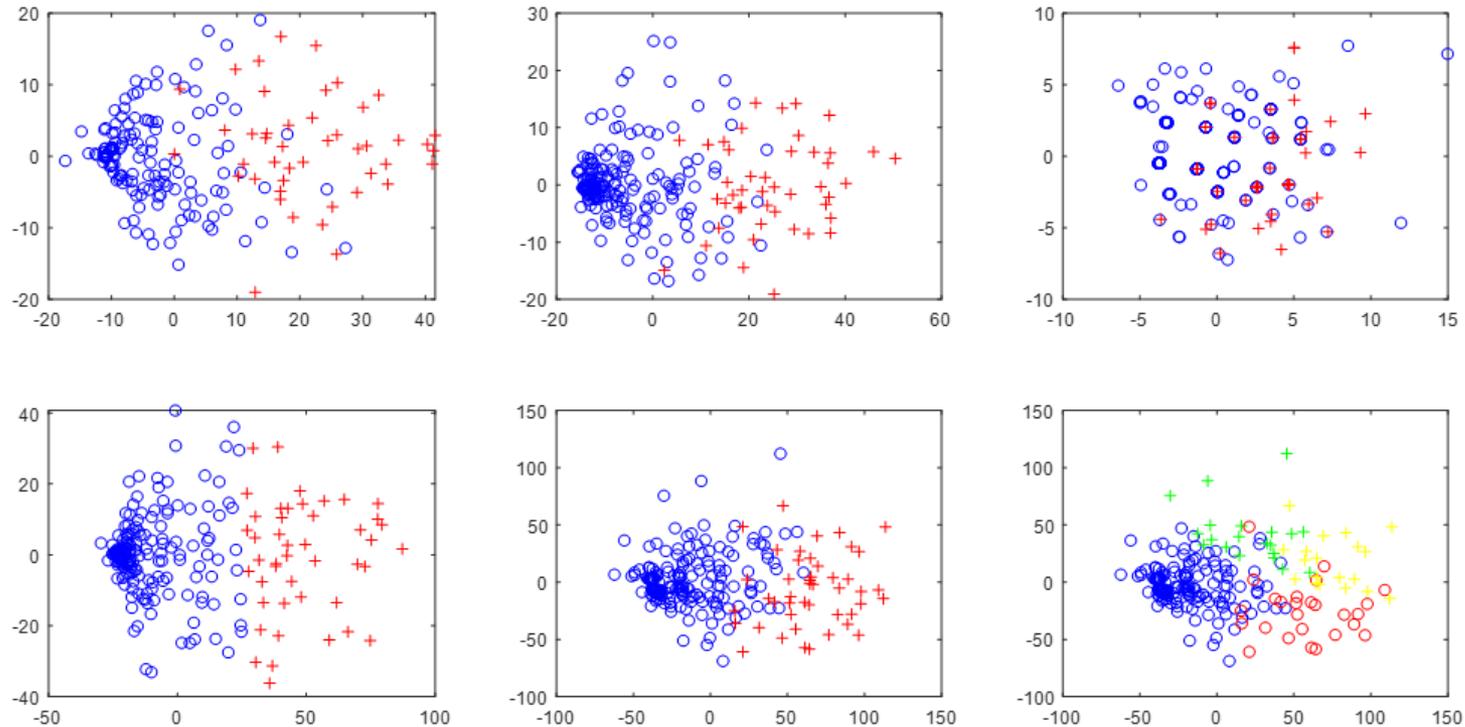
$$\min_{\alpha, W, H} C(P||WH^t) - \gamma En(\alpha) := E_{\alpha} [C(P_l||WH^t)] - \gamma En(\alpha)$$

- Normalized mutual information (NMI) criterion for selecting data for fusion

$$NMI(label_1, label_2) = \frac{2I(label_1, label_2)}{En(label_1) + En(label_2)}$$

	RPQ & GOSE	RPQ & PCL	GOSE & PCL
NMI	0.1036	0.2918	0.0842

Clustering Fusion



- Data from RPQ and PCL can be fused because of their high correlation/NMI. Two sets of labels are generated, one based on both RPQ and PCL (label_1), and the other based on GOSE (label_2).
- (top left) label_1 plotted against RPQ data using MDS. (top middle) label_1 plotted against PCL data using MDS. (top right) label_1 plotted against GOSE data using MDS.
- (bottom left) label_1 plotted against RPQ+PCL data using MDS. (bottom middle) label_1 plotted against RPQ+PCL+GOSE data using MDS. (bottom right) label_1 & label_2 plotted against RPQ+PCL+GOSE data using MDS.

Summary

- Complete integration of the diverse data for TBI diagnosis and patient stratification remains an unmet challenge.
- NMF is applied to patient similarity data for clustering which identifies meaningful phenotypes based on 6-month outcome summary scores
- Step-wise feature selection is used to characterize each cluster
- Future work
 - The proposed framework will be applied to a larger dataset to test if the conclusions hold
 - The data fusion approach will be applied to core variables to characterize the trend of disease development
 - Determine predictors of outcome at 6 months by including pre-injury, acute injury and post injury variables in prediction models.



Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.