

# Accounting for Heterogeneity Across Multiple Imaging Sites Using Multi-Task Learning

Michelle T. Hromatka<sup>1</sup>, Wei Liu<sup>1</sup>, Jeffrey S. Anderson<sup>2</sup>, Brandon A. Zielinski<sup>3</sup>,  
Molly B. DuBray<sup>4</sup>, and P. Thomas Fletcher<sup>1</sup>

<sup>1</sup> School of Computing

<sup>2</sup> Division of Neuroradiology

<sup>3</sup> Division of Child Neurology

<sup>4</sup> Interdepartmental Program in Neuroscience  
University of Utah, Salt Lake City, UT 84112, USA

**Abstract.** Combining imaging data from multiple sites has the potential to increase statistical power in clinical studies. However, while pooling multiple data sources increases sample size, it also increases unwanted variance due to inconsistencies across sites, e.g., different scanners, protocols, and demographics. In this paper, we present an approach for combining multi-site imaging data in classification tasks that takes this heterogeneity into account. The idea is to treat the classification problem as a multi-task learning problem, where each imaging site is treated as a “task”. We employ a regularized support vector machine (SVM) that allows for differences in decision boundaries at individual sites, while at the same time leveraging the similarities in the decision boundaries across sites. We demonstrate the effectiveness of this approach in the classification of autism from multi-site functional magnetic resonance imaging (fMRI) from the Autism Brain Imaging Data Exchange (ABIDE). The proposed method achieves state-of-the-art accuracy and outperforms a comparable SVM classifier applied to pooled data as well as individual SVM classifiers applied per site.

## 1 Introduction

Recent years have seen a movement towards combining neuroimaging data collected across multiple sites. Such multi-site data has the potential to accelerate scientific discovery by increasing sample sizes, providing broader ranges of participant demographics, and making data publicly available. Different approaches include large, coordinated multi-site neuroimaging studies, such as the Alzheimer’s Disease Neuroimaging Initiative (ADNI) [?], as well as data sharing initiatives that combine multiple single-site studies, such as the Autism Brain Imaging Data Exchange (ABIDE) [3]. Analysis using these datasets, however, is not straightforward, due to differences across site scanners, protocols, populations, and diagnosis techniques. Treating a multi-site study as a single, homogeneous data set fails to account for this variability, which can be detrimental to the statistical power and counteract the gains made by increasing the sample size.

In addition, many sites do not have a large enough sample size to use learning algorithms on site-specific data. Meta-analysis, in which results across small studies are combined to extract patterns common in each, has previously been used to combine site-specific results, especially when the sample size is low. However, meta-analysis is not free of subjectivity of data variability, thus this method of combined analysis is also faulty [4].

Several groups have used the ABIDE dataset for classification with different, non-meta-analytic approaches. Nielsen, et al. combines the ABIDE dataset with a whole-brain approach, using a leave-one-out classifier to compute a classification score for each left-out subject based on age, gender and handedness. The correlations for each connection in turn were fit with a linear model, separating controls from ASDs, which was then adjusted by the difference between the subject’s site mean for that connection and the overall mean. This approach yielded a maximum overall accuracy of 60.0% despite finding significant positive correlation between the classification score and several of the phenotypic behavioral measures [9]. A different study used histogram of gradients and applied this to several multi-site imaging studies which was able to achieve 61.7% accuracy on the ABIDE dataset and 62.6% on the ADHD-200 dataset [5]. While [9] accounted for some of the site differences, both studies approached the differences across imaging sites as noise instead of extra data that can be leveraged when classifying an aggregate data set. This is the key principle in multi-task learning: heterogeneity across similar tasks can be accounted for while using a common mean to account for similarities between the different tasks. This approach will be applied within an SVM classifier on the ABIDE dataset.

## 2 Methods

Classification of multi-site imaging data can be thought of as a multi-task learning problem where each site,  $s$ , is treated as a separate task,  $t$ . For a dataset with  $S$  sites, the corresponding multi-task problem will have a set of  $T = S$  tasks and a function  $F(t)$  that describes the relationship between each task. Prior to using any sort of machine learning, however, nuisance factors (e.g. age) should be removed from the data to ideally isolate differences across all data that are exclusive to the disease, disorder being studied.

### 2.1 Multi-Task Learning

Evgeniou and Pontil (2004) introduces a method of multi-task learning based on kernel based methods typically used for single task learning. This method relies on minimizing regularization functions, such as that for SVM, to capture both overall similarity between tasks and individual task differences. The traditional minimization for a soft margin SVM is:

$$\frac{1}{2}w^2 + CF\left(\sum_{i=1}^t \xi_i\right) \quad (1)$$

where  $C$  is a constant and  $F(\mu)$  is a "monotonic convex function with  $F(0) = 0$  [2]. In the case of SVMs, the weight vector  $w$  is used to define the hyperplane,  $(w \cdot x + b)$ , which is the boundary between groups.

For multi-task learning, the relationship between  $T$  tasks must be described, which Evgeniou and Pontil approach using the hierarchical Bayes method. This assumes that each task function comes from a class of probability distributions. The relationship is defined as:

$$w_t = w_0 + v_t, \quad (2)$$

where  $w_0$  is the mean of the data and each task  $t$  has its own weight vector,  $v_t$ . Multi-task learning allows for simultaneous learning of the mean of all tasks,  $w_0$ , and each task weight vector,  $v_t$ , so the minimization function then becomes:

$$C \sum_{t=1}^T \sum_{i=1}^m \xi_{it} + \frac{\lambda_1}{T} \sum_{t=1}^T \|v_t\|^2 + \lambda_2 \|w_0\|^2, \quad (3)$$

where  $\lambda_1, \lambda_2$  are "positive regularization parameters" and  $C$  is still a constant. For high similarity between tasks, the  $v_t$  will be small in relation to  $w_0$ ; this relationship is described by the hyperparameters  $\lambda_1, \lambda_2$  that must be chosen by the user.

The dual of equation 3 can be found by defining a set of functions  $f_t(x) = w_t \cdot x$  which can be simplified to  $F(x, t) = f_t(x)$ . This can be described by a kernel function  $\phi((x, t))$  which allows us to relate the dual of a multi-task learning problem to the dual of Equation 1.

$$\max_{\alpha_{it}} \left\{ \sum_{i=1}^m \sum_{t=1}^T \alpha_{it} - \sum_{i=1}^m \sum_{s=1}^T \sum_{j=1}^m \sum_{t=1}^T \alpha_{is} y_{is} \alpha_{jt} y_{jt} \phi((x, t)) \right\} \quad (4)$$

where

$$\phi((x, t)) = \left( \frac{x}{\sqrt{\mu}}, \underbrace{0, \dots, 0}_{t-1}, x, \underbrace{0, \dots, 0}_{T-t} \right), \quad \text{for } \mu = \frac{T\lambda_2}{\lambda_1}. \quad (5)$$

As you can see in Equation 4, this is the same dual problem as for a single task-SVM, with the data transformed by  $\phi((x, t))$  into the multi-task kernel space.

## 2.2 Feature Selection

Data extraction in imaging studies typically leads to very high dimensional data spaces. For f-MRI, a typical choice for data is the pairwise correlation between  $n$  predefined regions of the brain. This yields a dataspace of  $\frac{n(n+1)}{2}$  dimensionality, which, even for a relatively small number of regions, can be computationally expensive. The multi-task learning above further increases dimensionality with the number of tasks. For a feature space  $x \in \mathbb{R}^d$ , the multi-task kernel  $\phi((x, t))$  would yield a feature space of  $d(t+1)$  dimensions. Feature selection can and

should be employed to remove redundancy and increase relevancy of the data while reducing computation time [6].

One approach is to use the values of the weight vector  $w$  to choose the most relevant features. Recall that the decision boundary in an SVM is defined by  $w \cdot x$ , so the highest magnitudes in the weight vector denote the features that best define the decision boundary between groups. This approach requires choosing a user defined  $n$  size subset of features which introduces another parameter and layer of complexity but overall reduces computation time by dramatically reducing the dimensionality of the feature space.

### 3 Evaluation

#### 3.1 Data

The Autism Brain Imaging Data Exchange (ABIDE) database is an online consortium of resting-state functional-MRI data from 17 international sites, resulting in brain imaging data for 539 individuals with ASD and 573 typically developing (TD) controls [3]. All ASD subjects were diagnosed by either the Autism Diagnosis Observation Schedule-General (ADOS-G) or the Autism Diagnostic Interview-Revised tests and removed from the study if other co-morbid disorders were present [7] [8] [3]. Further inclusion details can be found at (put url? citation to... website? or the abide paper?)

**Preprocessing** All data was preprocessed using the Functional Connectomes-1000 preprocessing scripts [1]. This includes:

1. MRI Deoblique, reorient, skull strip
2. f-MRI Reorient, motion correct, skull strip, smooth
3. registration
4. Segmentation - csf, white matter
5. extracting global signal, from csf and wm
6. extract time series, Z-transform correlations
7. spatial smoothing, register to atlas
8. some sort of regression

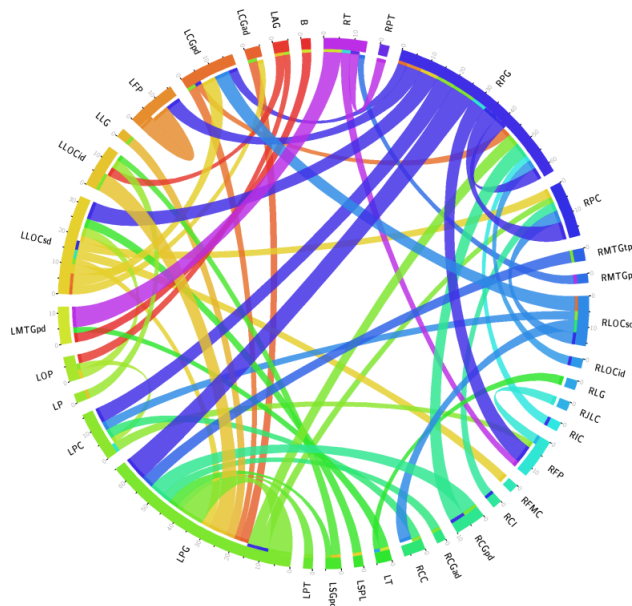
Twelve subjects were removed because of failure during the preprocessing. Two Oregon subjects were missing the resting fMRI file and 10 UCLA subjects were missing the anatomical scan file which is required in step 1 of the preprocessing pipeline above. This resulted in 1100 subjects for analysis, 530 ASD and 570 TD controls.

**Data Extraction** From each subject's postprocessed image, the time series for each of 264 regions is extracted based on Power's regions of interest [10]. These 264 regions are spread out among the cerebral cortex, subcortical structures and cerebellum, where each region is a sphere of 5mm in radius and regions are separated by a minimum distance of 10mm so as to avoid detection of a shared

signal. The Fisher transformed Pearson correlation coefficient is then found between each region and the other 263 regions, resulting in a 34,716 dimensional feature space for each subject. After feature selection, this number is reduced to 312 features per subject. \*\*\*add in featsel, multiplication/interdependence of features??

### 3.2 Results

## 4 Discussion and Conclusion



**Fig. 1.** The pairwise connections, represented as ribbons connecting each region, selected by the hypothesis test described in 2.2. The width of each ribbon is determined by the weights from the  $w_0$  vector in the SVM.

## References

1. Bharat B Biswal, Maarten Mennes, Xi-Nian Zuo, Suril Gohel, Clare Kelly, Steve M Smith, Christian F Beckmann, Jonathan S Adelstein, Randy L Buckner, Stan Colcombe, et al. Toward discovery science of human brain function. *Proceedings of the National Academy of Sciences*, 107(10):4734–4739, 2010.
2. Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.

3. A Di Martino, CG Yan, Q Li, E Denio, FX Castellanos, K Alaerts, JS Anderson, M Assaf, SY Bookheimer, M Dapretto, et al. The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain architecture in autism. *Molecular psychiatry*, 2013.
4. Hans J Eysenck. Meta-analysis and its problems. *BMJ: British Medical Journal*, 309(6957):789, 1994.
5. Sina Ghiassian, Russell Greiner, Ping Jin, and Matthew RG Brown. Learning to classify psychiatric disorders based on fmri images: Autism vs healthy and adhd vs healthy.
6. Isabelle Guyon and André Elisseeff. An introduction to variable and feature selection. *The Journal of Machine Learning Research*, 3:1157–1182, 2003.
7. Catherine Lord, Susan Risi, Linda Lambrecht, Edwin H Cook Jr, Bennett L Leventhal, Pamela C DiLavore, Andrew Pickles, and Michael Rutter. The autism diagnostic observation schedulegeneric: A standard measure of social and communication deficits associated with the spectrum of autism. *Journal of autism and developmental disorders*, 30(3):205–223, 2000.
8. Catherine Lord, Michael Rutter, and Ann Le Couteur. Autism diagnostic interview-revised: a revised version of a diagnostic interview for caregivers of individuals with possible pervasive developmental disorders. *Journal of autism and developmental disorders*, 24(5):659–685, 1994.
9. Jared A Nielsen, Brandon A Zielinski, P Thomas Fletcher, Andrew L Alexander, Nicholas Lange, Erin D Bigler, Janet E Lainhart, and Jeffrey S Anderson. Multisite functional connectivity mri classification of autism: Abide results. *Frontiers in human neuroscience*, 7, 2013.
10. Jonathan D Power, Alexander L Cohen, Steven M Nelson, Gagan S Wig, Kelly Anne Barnes, Jessica A Church, Alecia C Vogel, Timothy O Laumann, Fran M Miezin, Bradley L Schlaggar, et al. Functional network organization of the human brain. *Neuron*, 72(4):665–678, 2011.