

A Derivation of the f -divergence Lower Bound

A detailed derivation of (6) is given here. The lower bound of the f -divergence is given as

$$\begin{aligned}
 & d_F(P^T(\mathbf{h}), Q^s(\mathbf{h})) \\
 &= \int P^T(\mathbf{h}) F\left(\frac{Q^s(\mathbf{h})}{P^T(\mathbf{h})}\right) d\mathbf{h} \\
 &\geq \int Q^s(\mathbf{h}) r^s(\mathbf{h}) d\mathbf{h} - \int P^T(\mathbf{h}) \left(\frac{r^s(\mathbf{h})^2}{2} + r^s(\mathbf{h})\right) d\mathbf{h} \\
 &= -\frac{1}{2} \int P^T(\mathbf{h}) r^s(\mathbf{h})^2 d\mathbf{h} + \int Q^s(\mathbf{h}) r^s(\mathbf{h}) d\mathbf{h} - \\
 &\quad \int P^T(\mathbf{h}) r_s(\mathbf{h}) d\mathbf{h} \\
 &= -\frac{1}{2} \int P^T(\mathbf{h}) (\boldsymbol{\alpha}^s \boldsymbol{\phi}^s(\mathbf{h}))^\top (\boldsymbol{\alpha}^s \boldsymbol{\phi}^s(\mathbf{h})) d\mathbf{h} \\
 &\quad + \int P^s(\mathbf{h}|\mathbf{y}) \boldsymbol{\gamma}^T (\boldsymbol{\alpha}^s)^\top \boldsymbol{\phi}^s(\mathbf{h}) d\mathbf{h} - 1
 \end{aligned}$$

The second step follows the lower bound of f -divergence. The last step is derived because that $r^s(\mathbf{h})$ is an estimation of $\frac{Q^s(\mathbf{h})}{P^T(\mathbf{h})}$. (7) can be derived by re-arranging the terms.

B Proof of Theorem 3.1

In this section, we give the proof sketch of Theorem 3.1. For completeness, we repeat the theorem here.

Theorem. *Assume that $P^s(\mathbf{h}|\mathbf{y}) = P^T(\mathbf{h}|\mathbf{y}) = P(\mathbf{h}|\mathbf{y})$, the variance in the feature space is finite, and the label proportions are all non-zero. When the number of training and testing samples goes to infinity, $\hat{\boldsymbol{\gamma}}^T$ is asymptotically consistent for $\boldsymbol{\gamma}^T$ if $(\mathbf{M}^s)^\top \mathbf{M}^s$ is invertible for all s .*

Proof. Considering first a single source domain, the quadratic form in Equation 3 would give that the estimator is

$$\hat{\boldsymbol{\gamma}} = ((\mathbf{M}^s)^\top (\mathbf{M}^s))^{-1} (\mathbf{M}^s)^\top \boldsymbol{\mu}^T. \quad (13)$$

From the assumption that we have finite variance and that conditional distributions are equivalent for all source and target domains, the central limit theorem gives that as the number of samples increases

$$(\boldsymbol{\mu}_l^s - \boldsymbol{\mu}_l^*) \sim \mathcal{N}(\mathbf{0}, \frac{1}{n_s \gamma_l^s} \boldsymbol{\Sigma}_l), \quad (14)$$

This is equivalent to noting that \mathbf{M}^s is asymptotically centered around \mathbf{M}^* , because \mathbf{M}^s is just the concatenation of these individual vectors. Likewise, we have that asymptotically

$$\boldsymbol{\epsilon} = (\boldsymbol{\mu}^T - \mathbf{M}^* \boldsymbol{\gamma}^T) \sim \mathcal{N}(\mathbf{0}, \sum_{l=1}^L \frac{1}{n_T \gamma_l^T} \boldsymbol{\Sigma}_l). \quad (15)$$

Using the equality from (15) in our estimator of (13), we have that

$$\hat{\boldsymbol{\gamma}} = ((\mathbf{M}^s)^\top (\mathbf{M}^s))^{-1} (\mathbf{M}^s)^\top (\mathbf{M}^* \boldsymbol{\gamma}^T + \boldsymbol{\epsilon}). \quad (16)$$

Note that asymptotically the errors go to 0 from (14) and (15) on all terms in (16), so the estimator has the asymptotic expectation of

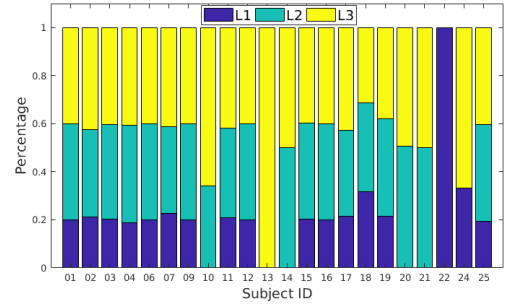
$$\begin{aligned}
 & \lim_{n_s \rightarrow \infty, n_T \rightarrow \infty} \mathbb{E}[\hat{\boldsymbol{\gamma}}^T] \\
 &= \lim_{n_s \rightarrow \infty, n_T \rightarrow \infty} \mathbb{E}[(\mathbf{M}^s)^\top (\mathbf{M}^s))^{-1} (\mathbf{M}^s)^\top (\mathbf{M}^* \boldsymbol{\gamma}^T + \boldsymbol{\epsilon})] \\
 &= ((\mathbf{M}^*)^\top (\mathbf{M}^*))^{-1} (\mathbf{M}^*)^\top \mathbf{M}^* \boldsymbol{\gamma}^T = \boldsymbol{\gamma}^T,
 \end{aligned}$$

if the inverse exists. Therefore, by the weak law of large numbers, $\lim_{n_s \rightarrow \infty, n_T \rightarrow \infty} \hat{\boldsymbol{\gamma}}^T = \boldsymbol{\gamma}^T$. \square

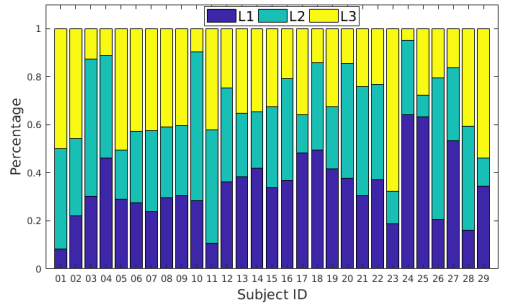
In two domains, $\lim_{n_s \rightarrow \infty, n_T \rightarrow \infty} \hat{\boldsymbol{\gamma}}^T = \boldsymbol{\gamma}^T$ shows the final result of Theorem 3.1.

C Time Series Data Label Proportion

The label proportion of the two EEG datasets used in experiment Section 5.2 are visualized in Figure 3.



(a) ASD Dataset Label Proportion



(b) LFP Dataset Label Proportion

Figure 3: (Left) Label constitution of the ASD dataset in each domain. ‘L1’, ‘L2’, ‘L3’ means before treatment, six months after treatment and twelve months after treatment, respectively. (Right) Label constitution of the LFP dataset. This dataset is for mice behavior classification using LFP signal. ‘L1’, ‘L2’, ‘L3’ mean home cage, open field and tail suspension, respectively. Please refer the the text for more detailed introduction.