
Supplementary Material: Representational Similarity Learning

1. Extending methods to matrices that are not Positive Semi-Definite

For any symmetric S we have an eigendecomposition UDU^T , where the eigenvalues may be negative. We can think of W as having the form BDB^T for some B and using the D from S .

So we consider the following optimizations instead.

$$\min_B \|S - XBDBX^T\|_F^2,$$

and

$$\min_B \|U - XB\|_F^2.$$

Then given a B construct $W = BDB^T$. The rest remains the same.

2. Clustering properties of GrOWL with absolute error loss function

Proof of Theorem 3.1

Proof. The proof is divided into two steps. First, we show $\|\hat{\beta}_j\| = \|\hat{\beta}_k\|$ and then we further show that the rows are equal. We proceed by contradiction. Assume $\|\hat{\beta}_j\| \neq \|\hat{\beta}_k\|$ and, without loss of generality, suppose $\|\hat{\beta}_j\| > \|\hat{\beta}_k\|$. We see that there exists a modification of the solution with a smaller GrOWL norm and same data-fitting term, and thus smaller overall objective value which contradicts our assumption that \hat{B} is the minimizer of $L(B) + G(B)$.

Consider the modification, $V = \hat{B}$ except $\hat{v}_j = \hat{\beta}_j - \varepsilon$ and $\hat{v}_k = \hat{\beta}_k + \varepsilon$ where $\varepsilon = \delta \hat{\beta}_j$ and δ is chosen such that $\|\varepsilon\| \in (0, \frac{\|\hat{\beta}_j\| - \|\hat{\beta}_k\|}{2}]$

Let $L(B) = \|Y - XB\|_1 = \|Y' - x_{.j}\hat{\beta}_j - x_{.k}\hat{\beta}_k\|_1$ where Y' is the residual term given by $Y' = Y - \sum_{i \neq j,k} x_{.i}\hat{\beta}_i$. Since $x_{.j} = x_{.k}$, L is invariant under this transformation, i.e., $L(V) = L(\hat{B})$. Same is true for $L(B) = \|Y - XB\|_F^2$.

Observe that the GrOWL norm of B is equal to the OWL norm of the vector of euclidean norms of rows of B . Since $\|v_k\| = \|\beta_k + \varepsilon\| \leq \|\beta_k\| + \|\varepsilon\|$, this transformation is equivalent to that defined in Lemma 3.1 and we have

$$G(\hat{B}) - G(V) \geq \Delta \|\varepsilon\|$$

This leads to a contradiction to our assumption that \hat{B} is the minimizer of $L(B) + G(B)$ and completes the proof that $\|\hat{\beta}_j\| = \|\hat{\beta}_k\|$. Now, let $\hat{\beta}_j + \hat{\beta}_k = z$, then the minimizer satisfies

$$\min_{\hat{\beta}_j, \hat{\beta}_k} w_j \|\hat{\beta}_j\| + w_k \|\hat{\beta}_k\|$$

$$\text{such that } \hat{\beta}_j + \hat{\beta}_k = z \text{ and } \|\hat{\beta}_j\| = \|\hat{\beta}_k\|$$

It is easy to see that the solution to this optimization is $\hat{\beta}_j = \hat{\beta}_k = z/2$ \square

Proof of Theorem 3.2

Proof. The proof is similar to the identical columns theorem. By contradiction and without loss of generality, suppose $\|\hat{\beta}_j\| > \|\hat{\beta}_k\|$. We show that there exists a transformation of \hat{B} such that the increase in the data fitting term is smaller than the decrease in the GrOWL norm.

Consider the modification, V , as defined in the proof of Theorem 3.1. By triangle inequality, the difference in loss function L that results from this modification satisfies

$$L(V) - L(\hat{B}) \leq \|x_{.j} - x_{.k}\|_1 \|\varepsilon\|_1$$

Invoking Lemma 3.1 as in the previous theorem and $\|\varepsilon\|_1 \leq \sqrt{r} \|\varepsilon\|$, we get

$$\begin{aligned} L(V) + G(V) - (L(\hat{B}) + G(\hat{B})) \\ \leq \sqrt{r} (\|x_{.j} - x_{.k}\|_1 - \frac{\Delta}{\sqrt{r}}) \|\varepsilon\| < 0 \end{aligned}$$

This contradicts our assumption that \hat{B} is the minimizer of $L(B) + G(B)$ and completes the proof for absolute loss. The proof with squared Frobenius loss can easily be extended using the inequality derived in Appendix B. \square

Proof of Theorem 3.3

Proof. The proof is similar to the identical columns theorem. By contradiction, suppose $\|\hat{\beta}_j - \hat{\beta}_k\| \geq \frac{8\phi \|\hat{\beta}_k\|}{4\phi^2 + 1} \geq \frac{2\|\hat{\beta}_k\|}{\phi}$. We show that there exists a transformation of \hat{B} such that the increase in the data fitting term is smaller than the decrease in the GrOWL norm.

Consider the modification, \mathbf{V} , as defined in the proof of Theorem 3.1 with $\varepsilon = \frac{\hat{\beta}_j - \hat{\beta}_k}{2}$. By triangle inequality, the difference in loss function L that results from this modification satisfies

$$L(\mathbf{V}) - L(\hat{\mathbf{B}}) \leq \|\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}\|_1 \|\varepsilon\|_1$$

We now bound the decrease in the GrOWL norm. Note by parallelogram law,

$$\begin{aligned} & \|\hat{\beta}_j + \hat{\beta}_k\|^2 \\ &= 2\|\hat{\beta}_j\|^2 + 2\|\hat{\beta}_k\|^2 - \|\hat{\beta}_j - \hat{\beta}_k\|^2 \\ &\leq 2\|\hat{\beta}_j\|^2 + 2\|\hat{\beta}_k\|^2 + \left(\frac{1}{4\phi^2} - \frac{1}{4\phi^2} - 1\right) \|\hat{\beta}_j - \hat{\beta}_k\|^2 \\ &\leq 4\|\hat{\beta}_j\|^2 + \left(\frac{\|\hat{\beta}_j - \hat{\beta}_k\|}{2\phi}\right)^2 - \frac{1 + 4\phi^2}{4\phi^2} \|\hat{\beta}_j - \hat{\beta}_k\|^2 \\ &\leq 4\|\hat{\beta}_j\|^2 + \left(\frac{\|\hat{\beta}_j - \hat{\beta}_k\|}{2\phi}\right)^2 - 2\frac{\|\hat{\beta}_j\| \|\hat{\beta}_j - \hat{\beta}_k\|}{\phi} \\ &\leq \left(\|\hat{\beta}_j\| + \|\hat{\beta}_k\| - \frac{\|\hat{\beta}_j - \hat{\beta}_k\|}{2\phi}\right)^2 \end{aligned}$$

Thus, we have

$$\begin{aligned} G(\hat{\mathbf{B}}) - G(\mathbf{V}) &\geq \Delta \left(\|\hat{\beta}_j\| + \|\hat{\beta}_k\| - \|\hat{\beta}_j + \hat{\beta}_k\| \right) \\ &\geq \frac{\Delta \|\hat{\beta}_j - \hat{\beta}_k\|}{2\phi} = \frac{\Delta \|\varepsilon\|}{\phi} \end{aligned}$$

Combining this with $\|\varepsilon\|_1 \leq \sqrt{r} \|\varepsilon\|$, we get

$$\begin{aligned} L(\mathbf{V}) + G(\mathbf{V}) - (L(\hat{\mathbf{B}}) + G(\hat{\mathbf{B}})) \\ \leq \left(\sqrt{r} \|\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}\|_1 - \frac{\Delta}{\phi} \right) \|\varepsilon\| < 0 \end{aligned}$$

This contradicts our assumption that $\hat{\mathbf{B}}$ is the minimizer of $L(\mathbf{B}) + G(\mathbf{B})$ and completes the proof for absolute loss. The proof with squared Frobenius loss can easily be extended using the inequality derived in Appendix B. \square

3. Clustering properties of GrOWL with squared Frobenius loss function

In this section, we consider the optimization

$$\min_{\mathbf{X}} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2 + G(\mathbf{B}) \quad (1)$$

Here we derive an upper bound on the increase in the squared loss term after applying the transformation, \mathbf{V} . We assume that the columns of the matrix, \mathbf{X} , are normalized

to a common norm, *i.e.*, ($\|\mathbf{x}_{\cdot i}\| = c$ for $i = 1, \dots, p$). Define $L(\mathbf{X}) = \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2 = \|\mathbf{Y}' - \mathbf{x}_{\cdot j}\beta_j - \mathbf{x}_{\cdot k}\beta_k\|_F^2$ where \mathbf{Y}' is again the residual term.

Lemma 1. *Let $\hat{\mathbf{B}} \in \mathbb{R}^{p \times r}$ and if \mathbf{V} is as defined in the respective theorems, then we have*

$$L(\mathbf{V}) - L(\hat{\mathbf{B}}) \leq \|\varepsilon\| \|\mathbf{Y}'\|_F \|\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}\|$$

Proof.

$$\begin{aligned} L(\mathbf{V}) - L(\hat{\mathbf{B}}) &= \frac{1}{2} \|\mathbf{Y}' - \mathbf{x}_{\cdot j}(\hat{\beta}_j - \varepsilon) - \mathbf{x}_{\cdot k}(\hat{\beta}_k + \varepsilon)\|_F^2 \\ &\quad - \frac{1}{2} \|\mathbf{Y}' - \mathbf{x}_{\cdot j}\hat{\beta}_j - \mathbf{x}_{\cdot k}\hat{\beta}_k\|_F^2 \end{aligned}$$

Expanding the Frobenius norm terms, canceling the common $\frac{1}{2} \|\mathbf{Y}'\|_F^2$ terms and using the common norm of columns ($\|\mathbf{x}_{\cdot i}\| = c$ for $i = 1, \dots, p$) we get

$$\begin{aligned} L(\mathbf{V}) - L(\hat{\mathbf{B}}) &= \frac{c^2}{2} \text{tr}((\hat{\beta}_j - \varepsilon)(\hat{\beta}_j - \varepsilon)^T + (\hat{\beta}_k + \varepsilon)(\hat{\beta}_k + \varepsilon)^T \\ &\quad - \hat{\beta}_j \hat{\beta}_j^T - \hat{\beta}_k \hat{\beta}_k^T) + \text{tr}(\mathbf{Y}'^T (\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}) \varepsilon) \\ &\quad + \text{tr}((\hat{\beta}_j - \varepsilon) \mathbf{x}_{\cdot j}^T \mathbf{x}_{\cdot k} (\hat{\beta}_k + \varepsilon)^T - \hat{\beta}_j \mathbf{x}_{\cdot j}^T \mathbf{x}_{\cdot k} \hat{\beta}_k^T) \end{aligned}$$

Expanding terms and making further cancellations gives

$$\begin{aligned} L(\mathbf{V}) - L(\hat{\mathbf{B}}) &= \text{tr}(\mathbf{Y}'^T (\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}) \varepsilon) - (c^2 - \mathbf{x}_{\cdot j}^T \mathbf{x}_{\cdot k}) \text{tr}((\hat{\beta}_j - \hat{\beta}_k - \varepsilon) \varepsilon^T) \\ &\leq \text{tr}(\mathbf{Y}'^T (\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}) \varepsilon) \\ &\quad - (c^2 - \mathbf{x}_{\cdot j}^T \mathbf{x}_{\cdot k}) \|\varepsilon\| (\|\hat{\beta}_j\| - \|\hat{\beta}_k\| - \|\varepsilon\|) \\ &\leq \text{tr}(\mathbf{Y}'^T (\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}) \varepsilon^T) \\ &\leq \|\mathbf{Y}'\|_F \|(\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}) \varepsilon\|_F \\ &= \|\varepsilon\| \|\mathbf{Y}'\|_F \|\mathbf{x}_{\cdot j} - \mathbf{x}_{\cdot k}\| \end{aligned}$$

where the first inequality follows from simplification and Cauchy-Schwarz inequality. The second inequality follows from $c^2 > \mathbf{x}_{\cdot j}^T \mathbf{x}_{\cdot k}$ and $\|\hat{\beta}_j\|_2 - \|\hat{\beta}_k\|_2 - \|\varepsilon\| > 0$ (by assumption). The third inequality follows, again, by Cauchy-Schwarz inequality. \square

Using this Lemma one can easily extend the clustering properties of GrOWL to the optimization in (1).

4. Proximal algorithms for GrOWL

Proof. Outline: the proof proceeds by finding a lower bound for the objective function in (5) and then we show that the proposed solution achieves this lower bound.

First, note that the following is true for any \mathbf{B} and \mathbf{V} ,

$$\begin{aligned}\|\mathbf{B} - \mathbf{V}\|_F^2 &= \sum_{i=1}^p \|\beta_{i\cdot} - \mathbf{v}_{i\cdot}\|^2 \\ &\geq \sum_{i=1}^p (\|\beta_{i\cdot}\| - \|\mathbf{v}_{i\cdot}\|)^2 = \|\tilde{\beta} - \tilde{\mathbf{v}}\|^2\end{aligned}$$

where the inequality follows from reverse triangle inequality.

Combining this with $G(\mathbf{B}) = \Omega_{\mathbf{w}}(\tilde{\beta})$, we have a lower bound on the objective function in (5). For all $\mathbf{B} \in \mathbb{R}^{p \times r}$

$$\frac{1}{2}\|\mathbf{B} - \mathbf{V}\|_F^2 + G(\mathbf{B}) \geq \frac{1}{2}\|\text{prox}_{\Omega_{\mathbf{w}}}(\tilde{\mathbf{v}}) - \tilde{\mathbf{v}}\|^2 + \Omega_{\mathbf{w}}(\text{prox}_{\Omega_{\mathbf{w}}}(\tilde{\mathbf{v}}))$$

Finally, we show that $\mathbf{B} = \hat{\mathbf{V}}$ achieves this lower bound,

$$\begin{aligned}&\frac{1}{2}\|\hat{\mathbf{V}} - \mathbf{V}\|_F^2 + G(\hat{\mathbf{V}}) \\ &= \frac{1}{2}\sum_{i=1}^p \left\| \left(\text{prox}_{\Omega_{\mathbf{w}}}(\tilde{\mathbf{v}}) \right)_i \frac{\mathbf{v}_{i\cdot}}{\|\mathbf{v}_{i\cdot}\|} - \mathbf{v}_{i\cdot} \right\|_2^2 + \Omega_{\mathbf{w}}(\text{prox}_{\Omega_{\mathbf{w}}}(\tilde{\mathbf{v}})) \\ &= \frac{1}{2}\|\text{prox}_{\Omega_{\mathbf{w}}}(\tilde{\mathbf{v}}) - \tilde{\mathbf{v}}\|_2^2 + \Omega_{\mathbf{w}}(\text{prox}_{\Omega_{\mathbf{w}}}(\tilde{\mathbf{v}}))\end{aligned}$$

□