

Supplementary Material for Learning Ordered Representations with Nested Dropout

A. Proofs for Section 3.2

Theorem 1. *Every optimal solution of the nested dropout problem is necessarily an optimal solution of the standard autoencoder problem.*

Proof. Let the nested dropout autoencoder be of latent dimension K . Recall that the nested dropout objective function in Equation (11) is a strictly positive mixture of the K different b -truncation problems. As described in Subsection 3.1, an optimal solution to each b -truncation must be of the form $\mathbf{X}_b^* = \mathbf{T}_b \boldsymbol{\Sigma}_{\downarrow b} \mathbf{R}^T$, $\boldsymbol{\Gamma}_b^* = \mathbf{Q}_{\downarrow b} \mathbf{T}_b^{-1}$ for some invertible transformation \mathbf{T}_b . We note that the PCA decomposition is a particular optimal solution for each b that is given for the choice $\mathbf{T}_b = \mathbb{I}_b$. As such, the PCA decomposition exactly minimizes every term in the nested dropout mixture, and therefore must be a global solution of the nested dropout problem. This means that every optimal solution of the nested dropout problem must exactly minimize every term in the nested dropout mixture. In particular, one of these terms corresponds to the K -truncation problem, which is in fact the original autoencoder problem. ■

Denote $\mathbf{T}_{\downarrow b} = \mathbf{J}_{K \rightarrow b} \mathbf{T} \mathbf{J}_{K \rightarrow b}^T$ as the b -th leading principal minor and its its bottom right corner as $t_b = T_{bb}$.

Lemma 1. *Let $\mathbf{T} \in \mathbb{R}^{K \times K}$ be commutative in its truncation and inversion. Then all the diagonal elements of \mathbf{T} are nonzero, and for each $b = 2, \dots, K$, either $\mathbf{A}_b = \mathbf{0}$ or $\mathbf{B}_b = \mathbf{0}$.*

Proof. We have $\det \mathbf{T}_{\downarrow b} = \det \mathbf{T}_{\downarrow b-1} \det(t_b - \mathbf{B}_b \mathbf{T}_{\downarrow b-1}^{-1} \mathbf{A}_b) \neq 0$ since $\mathbf{T}_{\downarrow b-1}$ is invertible. Since $\mathbf{T}_{\downarrow b-1}$ is also invertible, then $t_b - \mathbf{B}_b \mathbf{T}_{\downarrow b-1}^{-1} \mathbf{A}_b \neq 0$. As such, we write $\mathbf{T}_{\downarrow b}$ in terms of blocks $\mathbf{T}_{\downarrow b-1}$, \mathbf{A}_b , \mathbf{B}_b , t_b , and apply blockwise matrix inversion to find that $\mathbf{T}_{\downarrow b-1}^{-1} = \mathbf{T}_{\downarrow b-1}^{-1} + \mathbf{T}_{\downarrow b-1}^{-1} \mathbf{A}_b (t_b - \mathbf{B}_b \mathbf{T}_{\downarrow b-1}^{-1} \mathbf{A}_b)^{-1} \mathbf{B}_b \mathbf{T}_{\downarrow b-1}^{-1}$ which reduces to $\mathbf{A}_b \mathbf{B}_b = \mathbf{0}$. Now, assume by contradiction that $t_b = 0$. This means that either bottom row or the rightmost column of $\mathbf{T}_{\downarrow b}$ must be all zeros, which contradicts with the invertibility of $\mathbf{T}_{\downarrow b}$. ■

Theorem 2. *Every optimal solution of the nested dropout problem must be of the form*

$$\mathbf{X}^* = \mathbf{T} \boldsymbol{\Sigma} \mathbf{R}^T \quad (1)$$

$$\boldsymbol{\Gamma}^* = \mathbf{Q} \mathbf{T}^{-1}, \quad (2)$$

for some matrix $\mathbf{T} \in \mathbb{R}^{K \times K}$ that is commutative in its truncation and inversion.

Proof. Consider an optimal solution $\mathbf{X}^*, \boldsymbol{\Gamma}^*$ of the nested dropout problem. For each b -truncation, as established in the proof of Theorem 1, it must hold that

$$\mathbf{X}_b^* = \mathbf{T}_b \mathbf{J}_{K \rightarrow b} \boldsymbol{\Sigma} \mathbf{R}^T \quad (3)$$

$$\boldsymbol{\Gamma}_b^* = \mathbf{Q} \mathbf{J}_{K \rightarrow b}^T \mathbf{T}_b^{-1}. \quad (4)$$

However, it must also be true that $\mathbf{X}_b = \mathbf{X}_{\downarrow b}$, $\boldsymbol{\Gamma}_b = \boldsymbol{\Gamma}_{\downarrow b}$ by the definition of the nested dropout objective in Equation (11). The first equation thus gives that $\mathbf{T}_b \mathbf{J}_{K \rightarrow b} = \mathbf{J}_{K \rightarrow b} \mathbf{T}_K$, and therefore $\mathbf{T}_b = \mathbf{J}_{K \rightarrow b} \mathbf{T}_K \mathbf{J}_{K \rightarrow b}^T = \mathbf{T}_{\downarrow b}$. This establishes the fact that the optimal solution for each b -truncation problem simply draws the b -th leading principal minor from the same “global” matrix $\mathbf{T} := \mathbf{T}_K$. The second equation implies that for every b , it holds that $\mathbf{J}_{K \rightarrow b} \mathbf{T}^{-1} \mathbf{J}_{K \rightarrow b}^T = (\mathbf{J}_{K \rightarrow b} \mathbf{T} \mathbf{J}_{K \rightarrow b}^T)^{-1}$ and as such \mathbf{T} is commutative in its truncation and inversion. ■

Theorem 3. *Under the orthonormality constraint $\boldsymbol{\Gamma}^T \boldsymbol{\Gamma} = \mathbb{I}_K$, there exists a unique optimal solution for the nested dropout problem, and this solution is exactly the set of the K top eigenvectors of the covariance of \mathbf{Y} , ordered by eigenvalue magnitude. Namely, $\mathbf{X}^* = \boldsymbol{\Sigma} \mathbf{R}^T$, $\boldsymbol{\Gamma}^* = \mathbf{Q}$.*

Proof. The orthonormality constraint implies $(\mathbf{T}^{-1} \mathbf{Q})^T \mathbf{Q} \mathbf{T}^{-1} = \mathbb{I}_K$ which gives $\mathbf{T}^T = \mathbf{T}^{-1}$. Hence every row and every column must have unit norm. We also have that for every $b = 1, \dots, K$

$$\mathbf{T}_{\downarrow b}^T = (\mathbf{J}_{K \rightarrow b} \mathbf{T} \mathbf{J}_{K \rightarrow b}^T)^T \quad (5)$$

$$= \mathbf{J}_{K \rightarrow b} \mathbf{T}^T \mathbf{J}_{K \rightarrow b}^T \quad (6)$$

$$= \mathbf{J}_{K \rightarrow b} \mathbf{T}^{-1} \mathbf{J}_{K \rightarrow b}^T \quad (7)$$

$$= (\mathbf{J}_{K \rightarrow b} \mathbf{T} \mathbf{J}_{K \rightarrow b}^T)^{-1} \quad (8)$$

$$= \mathbf{T}_{\downarrow b}^{-1} \quad (9)$$

110	where in the last equation we applied Lemma 1 to Theorem	165
111	2. As such, every leading principal minor is also orthonor-	166
112	mal. For the sake of contradiction, assume there exist some	167
113	$m, n, m \neq n$ such that $T_{mn} \neq 0$. Without loss of general-	168
114	ity assume $m < n$. Then $\sum_{p=1}^{n-1} T_{mp}^2 < 1$, but this violates	169
115	the orthonormality of T_{n-1} . Thus it must be that the di-	170
116	agonal elements of T are all identically 1, and therefore	171
117	$T = \mathbb{I}_K$. The result follows. ■	172
118		173
119		174
120		175
121		176
122		177
123		178
124		179
125		180
126		181
127		182
128		183
129		184
130		185
131		186
132		187
133		188
134		189
135		190
136		191
137		192
138		193
139		194
140		195
141		196
142		197
143		198
144		199
145		200
146		201
147		202
148		203
149		204
150		205
151		206
152		207
153		208
154		209
155		210
156		211
157		212
158		213
159		214
160		215
161		216
162		217
163		218
164		219