

1 We thank all the reviewers for their valuable suggestions and feedback. We kindly appeal
 2 to all the reviewers that they will reconsider and improve their scores, because as shown
 3 again in the new results provided in this rebuttal our method achieved the best performance
 4 in comparison to all the benchmarks across the real-world datasets. We believe our work
 5 presents a valuable and general regularization method for supervised learning models.

6 **[[Reviewer 1]] ■ Strength of empirical results:** We would like to point out that the benefit
 7 of our method is not only (the magnitude of) the improved performance; we also show that
 8 our method is consistently providing the best regularization (across every dataset - both
 9 synthetic and real-data). For example, in the encircled portion of Figure A in this response,
 10 our method consistently achieved the best MSE in comparison to all the benchmarks across
 11 the real-world datasets. This consistency is not seen in the other benchmark regularizers,
 12 which exhibit higher variance in their average rank across all the datasets. This notion is
 13 also conveyed in the synthetic experiments in Figure 3 and for real-data (classification and
 14 regression) in Figure 5 in the original manuscript. ■ **Experiments with Mixup:** We appreciate the recommendation
 15 for additional benchmarks. We have conducted the recommended experiments for Mixup and Manifold Mixup and
 16 show that CASTLE still outperforms the other benchmarks as shown in Table A (please compare with CASTLE in
 17 Table 3 in the original manuscript) and Figure A for the real datasets.

18 **[[Reviewer 2]] ■ Notation clarification:** You are assuming correctly - $V(W)$ is the ℓ_1 norm in our methodology. We
 19 will clarify this in the revised manuscript. ■ **Generalization bounds:** The primary goal of our method is improving
 20 out-of-sample prediction performance which we use a generalization bound as justification. We did not prove the
 21 consistency of using a reconstruction loss and a norm-based regularizer in DAG learning which has already been proven
 22 in [49] and [50], respectively. ■ **CNN Limitations:** We agree with you and will clarify the limitations of our method as
 23 we did for CNNs in the Broader Impact statement. ■ **Additional results:** We reinforce the superiority of our proposed
 24 method by providing additional results in Table A and Figure A in the response to Reviewer 1.

25 **[[Reviewer 3]] ■ Experimental hyperparameters:** As mentioned in lines
 26 255-256, we performed a grid-search over a wide range of hyperparameters.
 27 We believe that this is fairly done for each benchmark, as we conducted the
 28 same grid-search for our model that we did for the other benchmark methods
 29 (lines 257-258) and applied early stopping for each. The performance gain
 30 from our regularizer is not due to improper hyperparameter tuning. ■ **L2**
 31 **missing:** Table 2 contains L2 regularization. ■ **Definition of W :** We are
 32 not defining W twice. The adjacency matrix can be represented by a matrix
 33 containing negative values - see NOTEARS [39] and Non-parametric DAGs
 34 [40]. Because of this, we can embed the adjacency matrix in the input layers,
 35 W 's, of the proposed neural network (Section 3.3). ■ **Prediction using**
 36 **causal parents:** Each feature is constructed using every other feature based
 37 on the DAG structure embedded in the neural network input layers. When
 38 a DAG is learned, the parent features (non-zero weights W) are obligated to construct each child (see Figure 2 and
 39 lines 166-168). ■ **Variable u :** In Def. 1, each variable u_i is specific for a feature X_i in the DAG, they are not hidden
 40 confounders, but the random noise to generate the feature X_i . ■ **Causal neighbors:** Consider the case where a variable
 41 is just noise, and therefore does not have any causally adjacent nodes (neighbors). Reconstruction methods, such as
 42 SAE, naively (and inefficiently) learns to reconstruct noise variables that have no causal implications on the target
 43 variable. Through DAG learning, our method does not reconstruct these variables as the input weight matrices get
 44 forced to zero (see Figures 3, 6, and 7). ■ **Sibling variables:** We mean the function generating the sibling variables
 45 may share the some similarities. We will elaborate this point with concrete examples in the revised submission. ■
 46 **X vs X :** We define the random variables as X and the corresponding data matrix as X . This is standard notation in
 47 machine learning. ■ **Regularization terms:** The description of the regularization terms is given in lines 162-168. We
 48 will describe them in more detail in the revised paper. ■ **Layer size:** We provided the layer size in Section 4 on line
 49 264. ■ **Acyclicity constraint:** Starting on line 165, we introduce and describe Theorem 1 from [39]. We describe
 50 Theorem 1 by saying, "the graph given by W is a DAG if and only if $R_W = 0$." ■ **Writing quality:** Although the
 51 other reviewers have positively acknowledged our exposition, we will work to improve the writing quality.

52 **[[Reviewer 4]] ■ Scalability:** For a typical dataset with hundreds of features or less, the computational training time
 53 does not differ significantly between the regularizers. For example, on simulated data, an experimental run with
 54 200 features, 2000 samples, and 200 epochs had an average training time of $\sim 55s$ and $\sim 64s$ for SAE and CASTLE,
 55 respectively, on an Intel i7-6850K CPU at 3.60GHz. We will incorporate a computational complexity analysis of our
 56 method as well as a demonstration of the computational trade-offs between our method and improved performance in
 57 the revised manuscript. ■ **Notations:** We will correct the suggested typos in the revision.

Table A: Additional Experiments (real data).

DATA	MIXUP	MIXUP-MAN
REGRESSION (MSE)		
BH	0.134 ± 0.019	0.130 ± 0.023
WQ	0.717 ± 0.030	0.712 ± 0.028
FB	0.329 ± 0.184	0.387 ± 0.316
BC	0.339 ± 0.025	0.325 ± 0.039
SP	0.208 ± 0.040	0.282 ± 0.023
CM	0.333 ± 0.025	0.386 ± 0.031
CLASSIFICATION (AUROC)		
CC	0.763 ± 0.008	0.772 ± 0.007
PD	0.814 ± 0.018	0.808 ± 0.009
BC	0.719 ± 0.020	0.728 ± 0.013
LV	0.586 ± 0.041	0.571 ± 0.025
SH	0.904 ± 0.015	0.916 ± 0.011
RP	0.798 ± 0.016	0.804 ± 0.018

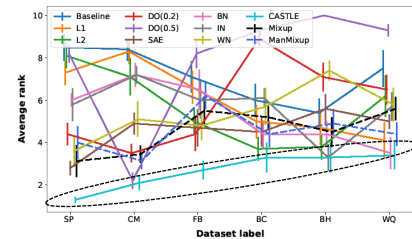


Figure A: Comparison in terms of average rank (in terms of MSE - lower is better). CASTLE has the best and most stable performance across all datasets.