

## A Appendix

### A.1 Experimental Settings

We provide the detailed model configurations and hyperparameter settings for DEEPCTRL in the following three experiments: double pendulum dynamics, sales forecasting, and cardiovascular classification.

#### A.1.1 Model Configurations

**Double Pendulum** The input and output states are 4 dimensional (two angular displacements and two angular velocities). The input state is fed into a shared layer whose configuration is [FC64,ReLU,FC16] where  $FC(n)$  denotes a fully-connected layer with  $n$  units. Then, output from the shared layer is fed into two encoders: Rule encoder [FC64,ReLU,FC64,ReLU,FC64] and Data encoder [FC64,ReLU,FC64,ReLU,FC64]. The combined representation from the two encoders is fed into a decision block [FC64,ReLU,FC4]. We use Adam optimization algorithm for training with learning rate 0.001.

**Sales Forecasting** The input dimension is 13, consisting of an item price as well as derivative features, and the output dimension is 1 (total weekly sales). Both encoders have same configuration [FC64,ReLU,FC64,ReLU,FC16] followed by a decision block [FC64,ReLU,FC1]. We use Adam optimization algorithm for training with learning rate 0.001.

**Cardiovascular Classification** The number of original input features is 11: AGE, HEIGHT, WEIGHT, GENDER, SYSTOLIC BLOOD PRESSURE, DIASTOLIC BLOOD PRESSURE, CHOLESTEROL, GLUCOSE, SMOKING, ALCOHOL INTAKE, PHYSICAL ACTIVITY. We expand categorical features to one-hot encoding (the input dimension is increased to 19) and the output dimension is 1 (Presence or absence of cardiovascular disease). Both encoders have same configuration [FC100,ReLU,FC16] followed by a decision block [FC1,Sigmoid]. We use Adam optimization algorithm for training with learning rate 0.001.

#### A.1.2 Data Splits

**Double Pendulum** The total length of simulated dynamics is 30,000 and it is split into training (18,000), validation (3,000), and testing (9,000).

**Sales Forecasting** Per the Kaggle task definition, first 273 weeks are used in a training set, the following 4 weeks are for a validation set, and the last 4 weeks are for a testing set. In each set, there are 22,543, 700, and 700 records from the preprocessed dataset.

**Cardiovascular Classification** The data partitioning is described in Section 4.3. For SOURCE partition, we have 70% training samples, 10% validation samples, and 20% testing samples to train and evaluate a model. Then, the trained model is applied to TARGET 1, TARGET 2, and TARGET 3, respectively.

Table 1: Dataset splits and USUAL vs. UNUSUAL partitioning.

DATA	SOURCE	TARGET 1	TARGET 2	TARGET 3
USUAL	6,007	20,000	6,000	4,000
UNUSUAL	14,018	6,009	6,009	6,009
RATIO	0.30	0.77	0.50	0.40

#### A.1.3 Training Settings

We commonly train DEEPCTRL for all tasks with a batch size of 32 on a single GPU (Nvidia T4 GPU) for 1000 epochs with early stopping where a validation error is not improved for 10 epochs. All results in the paper are mean values from 10 different random seeds.

### A.2 Perturbations

#### A.2.1 How to set $\delta x$ ?

In Section 3, we provide a method to integrate non-differentiable  $\mathcal{L}_{rule}$  via input perturbations. In this paper, we use the rules obtained via perturbation-based method for two tasks: sales forecasting (Section 4.2) and cardiovascular classification (Section 4.3). We summarize how the perturbations are generated and used to define

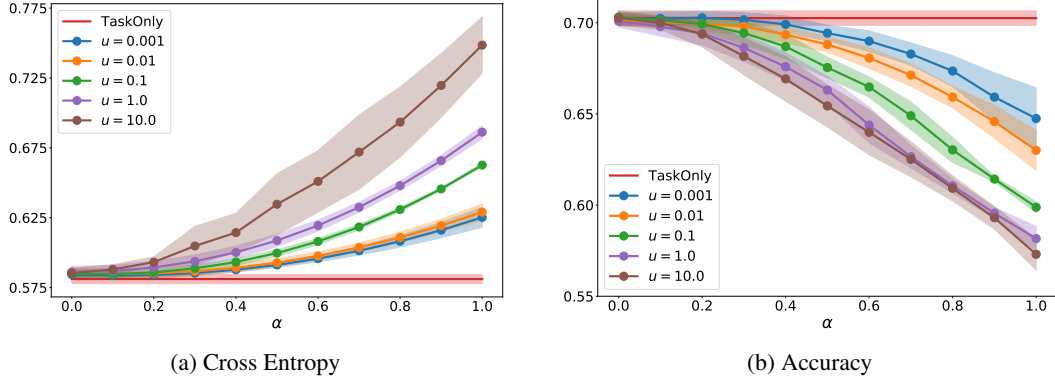


Figure 8: (Left) Cross entropy and (Right) Accuracy vs. rule strength for various upper bounds of perturbation scale.

the corresponding rule-based constraint for each task. Note that both tasks have a similar non-differentiable form of rule-constraint that the input  $\mathbf{x}$  and output  $\mathbf{y}$  have a negative or positive correlation.

**Perturbations in Sales Forecasting** There are a number of weekly-sales-records per an item at a particular store. For each record (week  $t$ ), we have input features including the price of an item  $\mathbf{x}_t$  and the target sales  $\mathbf{y}_t$ . The correlation coefficient between  $\mathbf{x}$  and  $\mathbf{y}$  is:

$$R = \frac{\sum_{t=1}^T (\mathbf{x}_t - \bar{\mathbf{x}})(\mathbf{y}_t - \bar{\mathbf{y}})}{\sqrt{\sum_{t=1}^T (\mathbf{x}_t - \bar{\mathbf{x}})^2} \sqrt{\sum_{t=1}^T (\mathbf{y}_t - \bar{\mathbf{y}})^2}}, \quad (4)$$

where  $\bar{\mathbf{x}}$  and  $\bar{\mathbf{y}}$  are the sample mean of  $\mathbf{x}_t$  and  $\mathbf{y}_t$ , respectively.

The rule-based constraint we want to impose is *price and sales should have a negative correlation coefficient*, i.e.  $R < 0$ . While the constraint is based on the exact definition of the correlation coefficient over  $T$  samples (Eq. 4), we use a constraint based on individual sample instead:

$$\frac{\Delta \mathbf{y}}{\Delta \mathbf{x}} = \frac{\mathbf{y}_p - \mathbf{y}}{\mathbf{x}_p - \mathbf{x}} < 0, \quad (5)$$

where  $\Delta \mathbf{x}$  is a price-difference and  $\Delta \mathbf{y}$  is a sales-difference, respectively, and  $\mathbf{x}_p$  is a perturbed price where  $\mathbf{x}_p = \mathbf{x} + \delta \mathbf{x}$  and  $\mathbf{y}_p$  is an output from the perturbed price. Note that Eq. 5 is identical to  $\mathbf{y}_p < \mathbf{y}$  once  $\delta \mathbf{x} > 0$ . There are two reasons not to use Eq. 4 directly. First, it is costly to compute the coefficient at every iteration. Second, it is possible to control  $\delta \mathbf{x}$  in Eq. 5. We set  $\delta \mathbf{x} = \gamma |\mathbf{x}|$ , where  $\gamma \sim \mathcal{U}[0, u]$  such that  $\gamma$  is a perturbation scale parameter and  $u$  is an upper bound of  $\gamma$ . Thus, the magnitude of  $\delta \mathbf{x}$  is bounded by  $[0, u|\mathbf{x}|]$ .

**Perturbations in Cardiovascular Classification** Similarly, we impose a positive correlation between blood pressure  $\mathbf{x}$  and a risk of the cardiovascular disease  $\mathbf{y}$ :

$$\frac{\Delta \mathbf{y}}{\Delta \mathbf{x}} = \frac{\mathbf{y}_p - \mathbf{y}}{\mathbf{x}_p - \mathbf{x}} > 0. \quad (6)$$

Eq. 6 is analogous to  $\mathbf{y}_p > \mathbf{y}$  as long as the perturbation  $\delta \mathbf{x} > 0$ .

**Upper Bound  $u$ :** We set the upper bound  $u$  as 0.1 for both tasks via the analysis described in Section A.2.2.

### A.2.2 Impact of Perturbation Scale

In previous section, we define  $\delta \mathbf{x}$  as a random scalar that is upper-bounded by  $u|\mathbf{x}|$ . We experiment on the cardiovascular classification task with different upper bounds,  $u$  to see how the scale of perturbation affects the model’s behavior.

Fig 8 shows how the cross entropy and accuracy are changed as the rule strength is changed. When the upper bound  $u$  is increased, the perturbation scale  $\gamma$  can be also increased, and thus, it leads to generate larger perturbations  $\delta \mathbf{x}$ . As  $u$  increases, the performance of a classifier is degraded when rule strength is non zero ( $\alpha > 0$ ). The curves imply that larger  $u$  leads more degraded performance when the rule strength is higher. When  $\delta \mathbf{x}$  increases, the perturbed output  $\mathbf{y}_p$  is more different to  $\mathbf{y}$  as  $\mathbf{y}_p$  is a function of  $\mathbf{x} + \delta \mathbf{x}$ . According to Eq. 5 and 6, the rule-based objective  $\mathcal{L}_{rule}$  is a function of  $\mathbf{y}_p - \mathbf{y}$ , and thus, the larger  $\delta \mathbf{x}$  eventually causes larger  $\mathcal{L}_{rule}$ . In other words, as  $\mathcal{L}_{rule}$  increases, DEEPCTRL is more driven by the rule-based objective when  $\alpha$  is non zero and thus, the task-based performance is degraded. As discussed, it is desired to have

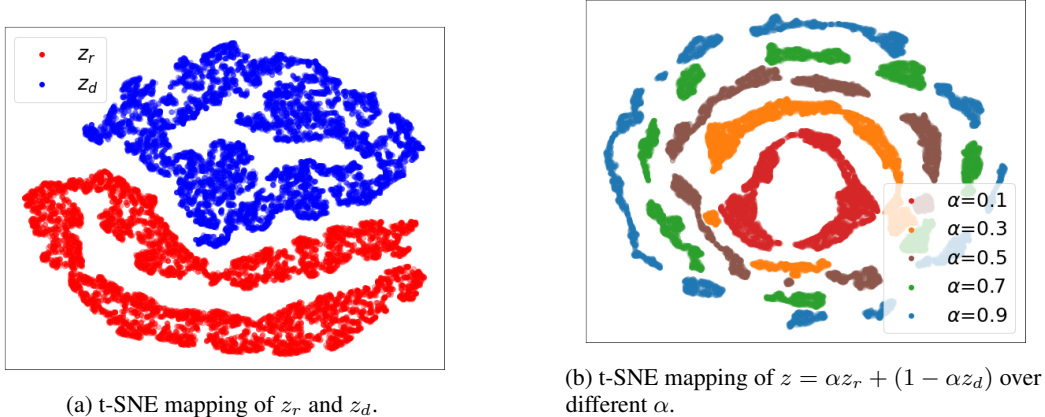


Figure 9: (Left) t-SNE visualization of  $z_r$  and  $z_d$  from rule encoder  $\phi_r$  and data encoder  $\phi_d$ , respectively, from the double pendulum task.  $z_r$  and  $z_d$  from 9,000 test samples do not have overlapped representations and it implies that two representations are distinct. (Right) t-SNE visualization of  $z = \alpha z_r + (1 - \alpha) z_d$  over different  $\alpha$ . It shows that task-/rule-specific representations are placed in inner/outer space, respectively. Note that the proposed representations gradually interpolate two extremes rather than being abruptly crossed.

Table 2: Training time (in seconds) for Sales forecasting (Retail) and Cardiovascular classification (Healthcare).

DATA	TASKONLY (TOTAL)	DEEPCTRL (TOTAL)	TASKONLY (PER EPOCH)	DEEPCTRL (PER EPOCH)
RETAIL	362.7 $\pm$ 62.3	328.2 $\pm$ 91.0	2.31 $\pm$ 0.04	2.32 $\pm$ 0.01
HEALTHCARE	154.7 $\pm$ 11.0	168.1 $\pm$ 23.1	3.33 $\pm$ 0.24	3.45 $\pm$ 0.17

distinct model’s behavior when  $\alpha = 0$  and  $\alpha = 1$  and thus, too small perturbation less incorporate rule-based representations. However, if too large perturbations are considered, the model is dominated by the rule mostly and the performance can be worse when  $\alpha$  is close to 0 (the brown curve is slightly higher than others when  $\alpha \rightarrow 0$  in Fig. 8a).

### A.3 Visualization of $z_r$ , $z_d$ and $z$

In this section, we analyze the learned representations to demonstrate the rule vs. data disentanglement capability of DEEPCTRL. We first visualize what the rule encoder  $\phi_r$  and data encoder  $\phi_d$  learn to support that the two encoders actually handle distinct representations rather than similar or overlapped representations. Then, we show how  $z = \alpha z_r + (1 - \alpha) z_d$  is changed over different  $\alpha$  to show meaningful combined representations with varying rule strength. Fig. 9 demonstrates t-SNE mapping of  $z_r$ ,  $z_d$ , and  $z$ , respectively. The learned representations of  $z_r$  and  $z_d$  are observed to be separated, while their linear combinations are clustered together with an orientation of the clusters highly dependent on the rule strength.

### A.4 Computational Complexity

Compared to TASKONLY training, the proposed method (DEEPCTRL) does not cause any additional computations that are proportional to the sample size. However, if perturbations are included, teaching rules via perturbation-based method cause an additional computation. For each batch, it is required to compute both the main forward pass ( $x$  to  $y$ ) and the perturbation-based forward pass ( $x_p$  to  $y_p$ ). Thus, the time complexity of DEEPCTRL is  $2C_{forward} + C_{backward}$  and that of TASKONLY is  $C_{forward} + C_{backward}$  where  $C_{forward}$  is the time complexity of forward propagation and  $C_{backward}$  of backward propagation. As  $C_{backward}$  takes the majority of training time, the extra computations from the perturbation-based method is not significant. Furthermore, since the computational complexity is independent on the size of samples, DEEPCTRL is still scalable. Table 2 shows the training time from TASKONLY and DEEPCTRL over 10 repeats. Note that the running time per an epoch from DEEPCTRL is similar to that of TASKONLY and it proves that the computational complexity of DEEPCTRL is not significantly increased.

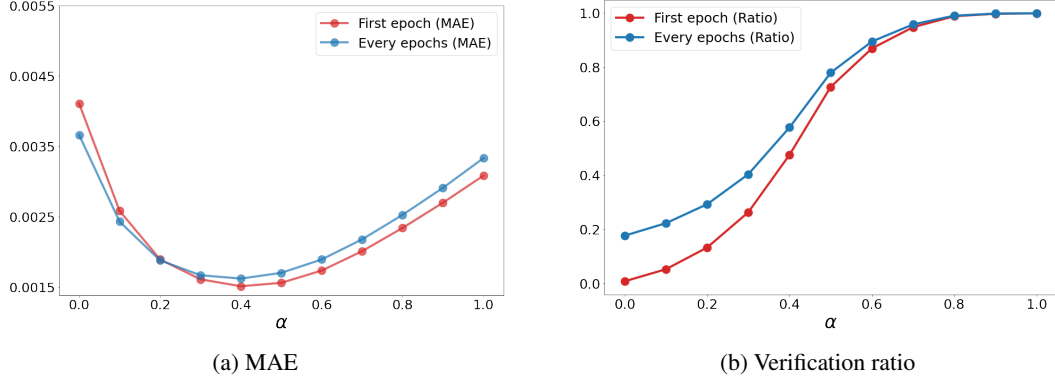


Figure 10: Results on a double pendulum task (different  $\rho$ ).

### A.5 Adaptive loss combination:

To balance the contributions from  $\mathcal{L}_{task}$  and  $\mathcal{L}_{rule}$ , we propose to use the scale parameter  $\rho = \mathcal{L}_{rule,0}/\mathcal{L}_{task,0}$ . One potential issue with this proposal is that at the beginning of training, the model is far from convergence, and the initial estimates of loss values may not be representative. One straightforward idea to address could be adapting  $\rho$  as the ratio after every epoch. In Fig. 10, we compare this approach to the proposed  $\rho = \mathcal{L}_{rule,0}/\mathcal{L}_{task,0}$  and we indeed show that the results are quite similar, and indeed  $\rho = \mathcal{L}_{rule,0}/\mathcal{L}_{task,0}$  yields better decoupling of  $\alpha = 0$  cases (as evident from 0 verification ratio).

We attribute this to the fact that the scale mismatching mostly happens at early phase of training, and  $\mathcal{L}_{rule}$  and  $\mathcal{L}_{task}$  become rapidly very small and the parameters are not significantly changed. This property is dependent on the type of rule, type of task, and dataset so that it is necessary to choose a proper method beforehand. Having a fixed coefficient is particularly beneficial for stable training because the model has a fixed target, rather than a varying one.