# Coupled generation\*

Ben Dai<sup>1</sup>, Xiaotong Shen<sup>1</sup>, and Wing Wong<sup>2</sup>

#### Abstract

Instance generation creates representative examples to interpret a learning model, as in regression and classification. For example, representative sentences of a topic of interest describe the topic specifically for sentence categorization. In such a situation, a large number of unlabeled observations may be available in addition to labeled data, for example, many unclassified text corpora (unlabeled instances) are available with only a few classified sentences (labeled instances). In this article, we introduce a novel generative method, called a coupled generator, producing instances given a specific learning outcome, based on indirect and direct generators. The indirect generator uses the inverse principle to yield the corresponding inverse probability, enabling to generate instances by leveraging an unlabeled data. The direct generator learns the distribution of an instance given its learning outcome. Then, the coupled generator seeks the best one from the indirect and direct generators, which is designed to enjoy the benefits of both and deliver higher generation accuracy. For sentence generation given a topic, we develop an embedding-based regression/classification in conjuncture with an unconditional recurrent neural network for the indirect generator, whereas a conditional recurrent neural network is natural for the corresponding direct generator. Moreover, we derive finite-sample generation error bounds for the indirect and direct generators to reveal the generative aspects of both methods thus explaining the benefits of the coupled generator. Finally, we apply the proposed methods to a real benchmark of abstract classification and demonstrate that the coupled generator composes reasonably good sentences from a dictionary to describe a specific topic of interest.

Keywords: Classification, Natural language processing, Numerical embeddings, Semisupervised generation, Unstructured data.

#### 1 Introduction

Generating an essay or a text for given structured information is an important Artificial Intelligence (AI) problem, which automatically imitates a certain style of writing. Whereas solving this AI problem is rather challenging, we tackle its simpler version in this article,

 $<sup>^{\</sup>ast 1}$  School of Statistics, University of Minnesota, Minneapolis, MN 55455;  $^2$  Department of Statistics and Biomedical Data Science, Stanford University, CA 94305 helpful comments and suggestions. Research supported in part by NSF grants DMS-1712564, DMS-1721216, DMS-1952539, DMS-1952386, and NIH grants 1R01GM126002, R01HL105397, R01AG065636.

which we call instance (example) generation, that is, generation of representative instances given a specific outcome to describe and interpret the corresponding learning model, for instance, classification and regression.

The use of black-box predictive models such as deep neural networks has delivered a high empirical learning accuracy in many real-life applications [14, 15]. Yet, it is difficult to make a sense of such a learning model. From the generative perspective, instance generation can describe the relationship between an instance and an outcome retrospectively. Its applications include a topic description of sentence categorization, abstractive text summarization [12], and image captioning [25], where generated sentences render descriptive examples of topics, texts, and images. In such a situation, sentence generation allows us to compose a novel essay and image captioning when the structured information is supplied. For example, the UCI abstract categorization benchmark consists of sentences from abstracts of articles, which are labeled with one of five topic categories. The goal here is learning a sentence generation mechanism to compose a novel abstract given a specific topic, in which the generation performance is measured by the cross-entropy error based on a test sample.

In the literature, instance generation, despite its vast important applications in AI, remains largely unexplored, although some approaches have been suggested for sentence generation. For example, a computational linguistics approach represents words/phrases as trees to model linguistic dependencies [20], a learning approach uses a large text corpus to learn a sentence's structure without any access to linguistic annotation [5]. In [26], a sentence generating model is proposed to produce a document by sampling the latent topic of a sentence and then words of the sentence using a recurrent neural network (RNN). In [37, 17], image captioning links the image content to a language model through an interplay between a convolution neural network (CNN) and an RNN. Yet, there is a paucity of works on instance generation given structured information, and incorporating both labeled and unlabeled data.

 $<sup>^{1}</sup>$  ${
m https://archive.ics.uci.edu/ml/datasets/Sentence+Classification}$ 

One of the primary characteristics of topic-instance data is that the amount of unlabeled data may be significantly larger than that of labeled data. For example, in sentence generation, uncategorized sentences are about ten times more than categorized ones. This is in a parallel situation of semisupervised learning with a different focus on leveraging unlabeled data to enhance the predictive accuracy of supervised learning [42, 18], which is in contrast to our generation objective given a learning outcome.

Our main contribution lies in the development of a new semisupervised generation framework for producing instances given an outcome. On this ground, we propose three generative methods-indirect, direct, and coupled generators. The indirect generator uses the principle of inverse learning to estimate the conditional probability distribution of an outcome given an instance, enabling to leverage unlabeled data, if available. On the other hand, the direct generator estimates the corresponding conditional probability of an instance given an outcome in a supervised manner. Then, the coupled generator is designed to enjoy the benefits of both generations. The proposed generators are illustrated in sentence generation, where we generate a sentence through sequential next-word-prediction. Specifically, we develop regularized embedding-based regression/classification in conjuncture with an unconditional RNN for the indirect generator, whereas we use a conditional RNN for the direct generator.

To shed light on the generative performance of the three generators, we derive finite-sample generation error bounds for each method. Interestingly, the generation error of the indirect generator is governed by the complexity of the parameter space of the conditional densities of an outcome given an instance and that of marginal densities. Similarly, that of the direct generator is determined by the conditional densities of an instance given an outcome. As a result, the indirect and direct generators have their own advantages with respect to generation with the unlabeled data is large, and importantly the coupled generator enjoys the benefits of both in terms of generation accuracy. This, together with a real benchmark of sentence categorization, demonstrates the utility of the coupled generation for composing

reasonably good sentences to describe a specific topic. Numerically, the proposed method outperforms a separate RNN method and the indirect generator can leverage additional unlabeled data for further enhancing the performance.

This paper is organized as follows. Section 2 introduces the framework of coupled generation based on indirect and direct generations. Section 3 develops a theory of the generation performance of the proposed methods. Section 4 is devoted to the development of a novel sentence generative method given a topic of interest through sequential next-word prediction. Section 5 investigates the operating characteristics of the coupled generator and compares it with the direct and indirect generators as well as one competitor. The Appendix contains technical proofs.

### 2 Methods

Consider a generative model in which the goal is to generate an instance X given an outcome Y, where X and Y represent instance and response variables, which can be numerical or unstructured such as texts and documents that cannot be expressed in a predefined manner. In this article, we focus on instance generation under a generative model, based on the conditional distribution  $p_{X|Y}$  of X given an outcome of Y. As an example, in sentence generation [26], instance generation produces representative examples of X given a specific topic of Y, where X and Y represent a sentence and its associated topic.

For instance generation, a labeled training sample  $(\boldsymbol{x}^i, \boldsymbol{y}^i)_{i=1}^n$  is available as well as an instance-only sample  $(\boldsymbol{x}^j)_{j=1}^{\tilde{n}}$ , whose sample size  $\tilde{n}$  may greatly exceed or smaller than the sample size n. In our context, we leverage the unlabeled sample to enhance the generative accuracy of instance generation.

Indirect generator. An indirect generator produces instances using an estimate of  $p_{X|Y}$  through the inverse relation (1): an estimate of  $p_{Y|X}$  based on  $(x^i, y^i)_{i=1}^n$  in (2) and the

marginal density  $p_{\mathbf{x}}$  based on combined data  $(\mathbf{x}^i)_{i=1}^n$  and  $(\mathbf{x}^j)_{j=1}^{\tilde{n}}$  in  $\mathfrak{B}$ . That is,

Indirect: 
$$\widehat{p}_{X|Y}^{b}(x|y) = \frac{\widehat{p}_{Y|X}(y|x)\widehat{p}_{X}(x)}{\int_{x \in \mathcal{X}} \widehat{p}_{Y|X}(y|x)\widehat{p}_{X}(x)dx},$$
 (1)

$$\widehat{p}_{Y|X} = \underset{p_{Y|X} \in \mathcal{F}_b}{\operatorname{argmin}} - n^{-1} \sum_{i=1}^{n} \log \left( p_{Y|X}(\boldsymbol{y}^i | \boldsymbol{x}^i) \right) + \lambda_b J_b(p_{Y|X}), \tag{2}$$

$$\widehat{p}_{\mathbf{X}} = \underset{p_{\mathbf{X}} \in \mathcal{F}_m}{\operatorname{argmin}} - (n + \widetilde{n})^{-1} \left( \sum_{i=1}^{n} \log \left( p_{\mathbf{X}}(\mathbf{x}^i) \right) + \sum_{j=1}^{\widetilde{n}} \log \left( p_{\mathbf{X}}(\mathbf{x}^j) \right) \right) + \lambda_m J_m(p_{\mathbf{X}}), \tag{3}$$

where  $\hat{p}_{Y|X}$  and  $\hat{p}_{X}$  in  $\boxed{1}$  are regularized maximum likelihood estimates of  $p_{Y|X}$  and  $p_{X}$ ,  $J_{b}$  and  $J_{m}$  are regularizers, for example,  $L_{1}$ - or  $L_{2}$ -regularization in a neural network model,  $\lambda_{b} \geq 0$  and  $\lambda_{m} \geq 0$  are tuning parameters controlling the weights of regularization, and  $\mathcal{F}_{b}$  in  $\boxed{2}$  and  $\mathcal{F}_{m}$  in  $\boxed{3}$  are parameter spaces of  $p_{Y|X}$  and  $p_{X}$ , respectively. Note that  $\int_{x \in \mathcal{X}} \hat{p}_{X|Y}(x|y) \hat{p}_{X}(x) dx$  in  $\boxed{1}$  normalizes  $\hat{p}_{X|Y}^{b}$  to become a probability density, although normalization is unnecessary when only some aspects of the distribution such as the modes or percentiles are of concern, as opposed to the distribution itself. Importantly, the indirect generator leverages instance-only (unlabeled) data  $(x^{j})_{j=1}^{\tilde{n}}$ , but any potential bias in estimation of  $p_{X}$  based on  $(x^{j})_{j=1}^{\tilde{n}}$  could translate into that of  $p_{Y|X}$ .

**Direct generator.** A direct generator uses  $p_{X|Y}$  to generate instances, estimated by minimizing the negative regularized conditional likelihood of X given Y based on  $(x^i, y^i)_{i=1}^n$ :

Direct: 
$$\widehat{p}_{\mathbf{X}|\mathbf{Y}}^{\mathbf{f}}(\mathbf{x}|\mathbf{y}) = \widehat{p}_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}),$$

$$\widehat{p}_{\mathbf{X}|\mathbf{Y}} = \underset{p_{\mathbf{X}|\mathbf{Y}} \in \mathcal{F}_{\mathbf{f}}}{\operatorname{argmin}} - n^{-1} \sum_{i=1}^{n} \log \left( p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}^{i}|\mathbf{y}^{i}) \right) + \lambda_{\mathbf{f}} J_{\mathbf{f}}(p_{\mathbf{X}|\mathbf{Y}}), \tag{4}$$

where  $\mathcal{F}_f$  is a parameter space of  $p_{X|Y}$ ,  $J_f$  is a regularizer, and  $\lambda_f \geq 0$  is a tuning parameter controlling the weight of regularization.

It appears that  $\P$  can be extended to leverage additional unlabeled data through the conditional likelihood of  $p_{X|Y}$  and a mixture relation  $\int_{\mathbf{y}} p_{X|Y}(\mathbf{x}|\mathbf{y}) p_{Y}(\mathbf{y}) d\mathbf{y} = p_{X}(\mathbf{x})$ . Unfor-

tunately, however, the mixture approach may suffer from an asymptotic bias when additional unlabeled data is included, thus degrading the estimation performance of  $p_{X|Y}$  [S. 9]. 39]. This is because the aforementioned mixture relation may not hold when  $\mathcal{F}_f$  is misspecified, and moreover its impact could be minimal even it holds, especially when the support of Y is large. As suggested by the theorem in Section 4 [39], the supervised and semisupervised maximum likelihood estimates may converge to different values, and thus more unlabeled data produces a larger estimation bias as measured by the Kullback-Leibler divergence, when the model is misspecified in that  $p_X^0$  does not belong to the parameter space  $\mathcal{F}_f = \{p_X(x) = \int_y p_{X|Y}(x|y)p_Y(y)dy; p_{X|Y} \in \mathcal{F}_f\}$  or the mixture relation is not satisfied. Furthermore, as demonstrated by Figures 1 and 2 of [3] and Figure 4.1 of [7], empirical studies indicate that an EM algorithm based on both labeled and unlabeled data tends to degrade performance solely based on the labeled data when the size of labeled data exceeds 30 in SecStr dataset. As a result,  $p_{X|Y}$  estimated from labeled data renders a better performance than that on labeled and unlabeled data.

In summary, how to leverage unlabeled data to enhance the generation performance remains an open question, which depends on model assumptions that may not be verifiable in practice. It is worth mentioning that (4) is a general formulation without assuming any specific assumption on how  $p_{\mathbf{x}}$  is related to  $\mathcal{F}_{\mathbf{f}}$ . However, if such an assumption becomes available in practice, (4) can be generalized based on it to incorporate unlabeled data for improvement. At present, we shall not pursue this aspect as the indirect method can benefit from additional unlabeled data, as suggested by Theorem (1) in Section 3.

Coupled generator. The level of difficulty of estimating  $\hat{p}_{X|Y}^f$  and that of  $\hat{p}_{X|Y}^b$  may differ, particularly when  $p_X$  can be well-estimated from both instance-only and unlabeled data. Depending on situations, the former may be more difficult than the latter, and vice versa. Some theoretical results for this aspect are illustrated in Section 4.5. Then we propose coupled generation by choosing, between the two, the one maximizing a predictive

log-likelihood, or minimizing a negative log-likelihood, such as (23) in the sentence generation example. In particular, a coupled generator is defined as,

$$\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{c} = \begin{cases}
\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{f} & \text{if } \widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{f} \text{ has a higher log-likelihood value on a validation set than } \widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{b}, \\
\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{b} & \text{otherwise.}
\end{cases}$$
(5)

The probability density  $\widehat{p}_{X|Y}^c$  gives the whole spectrum of values of X given Y. First, we may generate representative instances using the mode of  $p_{X|Y}$  to give one representation or sampling-based on  $p_{X|Y}$  for multiple representations. Second, discriminative features X with respect to Y can be extracted by comparing  $\widehat{p}_{X|Y}^c$  at different Y-values retrospectively. For example, in classification with  $Y = \pm 1$ , a comparison of  $\widehat{p}_{X|Y=1}^c$  and  $\widehat{p}_{X|Y=-1}^c$  leads to discriminative features. This aspect will be further investigated elsewhere.

Coupled learning has its distinct characteristics although it appears remotely related to semisupervised variational auto-encoders [18] and inverse autoregressive flows [19]. In particular, [19] uses a generative model  $p_{X|Y}$  and  $p_X$  to enhance a discriminative model  $p_{Y|X}$  regarding the marginal distribution as a mixture of conditional distributions, whereas the proposed indirect generator integrates the unlabeled data to separately estimate the marginal distribution. Furthermore, [18] estimates the marginal density of X  $p_X$  via a chain of latent factors and inevitable transformations of autoregressive neural networks and connects blocks by invertible relations. Yet, the proposed method links two conditional densities by Bayes' law. Finally, the theoretical justification of [19] and [18] remains unknown.

# 3 Theory

This section develops a learning theory to investigate the generation errors of direct, indirect, and coupled generators. In particular, we derive finite-sample generation error bounds for estimators  $\hat{p}_{X|Y}^b$ ,  $\hat{p}_{X|Y}^f$ , and  $\hat{p}_{X|Y}^c$  of (1), (4) and (5).

The generation error for generating X given Y is defined as the expected Hellingerdistance between two conditional densities  $p_{X|Y}$  and  $q_{X|Y}$  with respect to Y:

$$d(p_{\boldsymbol{X}|\boldsymbol{Y}},q_{\boldsymbol{X}|\boldsymbol{Y}}) = (\mathbb{E}_{\boldsymbol{Y}}h^2(p_{\boldsymbol{X}|\boldsymbol{Y}},q_{\boldsymbol{X}|\boldsymbol{Y}}))^{1/2} \equiv (\mathbb{E}_{\boldsymbol{Y}}\int (p_{\boldsymbol{X}|\boldsymbol{Y}}^{1/2} - q_{\boldsymbol{X}|\boldsymbol{Y}}^{1/2})^2 d\mu)^{1/2},$$

where  $\mu$  is the Lebesgue measure on  $\boldsymbol{x}$ , and  $\mathbb{E}_{\boldsymbol{Y}}$  is the expectation with respect to  $\boldsymbol{Y}$ .

Three parameter spaces  $\mathcal{F}_b$ ,  $\mathcal{F}_m$ , and  $\mathcal{F}_f$  are defined for estimating  $p_{Y|X}$  in (2),  $p_X$  in (3), and  $p_{X|Y}$  in (4), each of which is allowed to depend on the corresponding sample size. Then their regularized parameter spaces are given as follows:  $\mathcal{F}_{b,k} = \{p_{Y|X} \in \mathcal{F}_b : J(p_{Y|X}) \leq k\}$  for (2),  $\mathcal{F}_{m,k} = \{p_X \in \mathcal{F}_m : J(p_X) \leq k\}$  for (3), and  $\mathcal{F}_{f,k} = \{p_{X|Y} \in \mathcal{F}_f : J(p_{X|Y}) \leq k\}$  for (4). On this ground, we define the metric entropy to measure their complexities to be used for our theory.

The *u*-bracketing metric entropy  $H(u, \mathcal{F})$  of space  $\mathcal{F}$  with respect to a distance D is defined as the logarithm of the cardinality of the *u*-bracketing of  $\mathcal{F}$  of the smallest size. A *u*-bracketing of  $\mathcal{F}$  is a finite set (of pairs of functions)  $\{(p_j^L, p_j^U), j = 1, \dots, N\}$  such that for any  $p \in \mathcal{F}$ , there is a j such that  $p_j^L \leq p \leq p_j^U$  with  $d(p_j^L, p_j^U) \leq u$ ;  $j = 1, \dots, N$ . Note that  $d^2(p_{Y|X}, q_{Y|X}) = \mathbb{E}_{\mathbf{X}} \left(h^2(p_{Y|X}, q_{Y|X})\right)$ ,  $h^2(p_X, q_X)$ , and  $\mathbb{E}_{\mathbf{Y}} h^2(p_{X|Y}, q_{X|Y})$ , respectively for  $\mathcal{F}_{b,k}$ ,  $\mathcal{F}_{m,k}$ , and  $\mathcal{F}_{f,k}$ .

To quantify the degree of approximation of the true density  $p_{Y|X}^0$  by  $\mathcal{F}_b$ , we introduce a distance  $\rho_b(p_{Y|X}^0, p_{Y|X}) = \mathbb{E}_{\mathbf{X}} \mathbb{E}_{\mathbf{Y}|\mathbf{X}} g_{\alpha}(p_{Y|X}^0/p_{Y|X})$ , where  $g_{\alpha}(x) = \alpha^{-1}(x^{\alpha} - 1)$  for  $\alpha \in (0, 1)$ . As suggested in Section 4 of [38], this distance is stronger than the corresponding Hellinger distance. Similarly,  $\rho_m(p_X^0, p_X) = \mathbb{E}_{\mathbf{X}} g_{\alpha}(p_X^0/p_X)$  and  $\rho_f(p_{X|Y}^0, p_{X|Y}) = \mathbb{E}_{\mathbf{X}} \mathbb{E}_{\mathbf{Y}|X} g_{\alpha}(p_{X|Y}^0/p_{X|Y})$  are defined for approximating the true densities  $p_X^0$  and  $p_{X|Y}^0$  by  $\mathcal{F}_m$  and  $\mathcal{F}_f$ , respectively.

Let  $p_{Y|X}^* \in \mathcal{F}_b$  and  $p_X^* \in \mathcal{F}_m$  be two approximating points of  $p_{Y|X}^0$  and  $p_X^0$  in that  $\rho_b(p_{Y|X}^0, p_{Y|X}^*) \leq \gamma_b$  and  $\rho_m(p_X^0, p_X^*) \leq \gamma_m$  for some sequences  $\gamma_b \geq 0$  and  $\gamma_m \geq 0$ . Of course,  $\gamma_b = 0$  when  $p_{Y|X}^0 \in \mathcal{F}_b$  and  $\gamma_m = 0$  when  $p_X^0 \in \mathcal{F}_m$ .

**Theorem 1** (Indirect generator). Suppose there exist some positive constants  $c_1$ - $c_6$ , such that, for any  $\epsilon_b > 0$  and  $\lambda_b \geq 0$ ,

$$\sup_{k>1} \int_{2^{-8}L_k}^{2^{1/2}L_k^{1/2}} H^{1/2}(u/c_3, \mathcal{F}_{b,k}) du/L_k \le c_2 n^{1/2}, \quad L_k = c_1 \epsilon_b^2 + \lambda_b(k-1), \tag{6}$$

and, for any  $\epsilon_m > 0$  and  $\lambda_m \geq 0$ ,

$$\sup_{k\geq 1} \int_{2^{-8}L_k}^{2^{1/2}L_k^{1/2}} H^{1/2}(u/c_6, \mathcal{F}_{m,k}) du/L_k \leq c_5(n+\tilde{n})^{1/2}, \quad L_k = c_4 \epsilon_m^2 + \lambda_m(k-1), \quad (7)$$

then

$$P(d(\hat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{b}, p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) \geq 2(\eta_{b} + \eta_{m})) \leq 8 \exp(-c_{7}n\eta_{b}^{2}) + 8 \exp(-c_{8}(n + \tilde{n})\eta_{m}^{2}),$$

$$\eta_{b} = \max(\epsilon_{b}, \gamma_{b}^{1/2}), \quad \eta_{m} = \max(\epsilon_{m}, \gamma_{m}^{1/2}),$$

$$(8)$$

provided that  $\lambda_b \max(J_b(p_{\mathbf{Y}|\mathbf{X}}^*), J_b(p_{\mathbf{Y}|\mathbf{X}}^0), 1) \leq c_9 \eta_b^2$  and  $\lambda_m \max(J_m(p_{\mathbf{X}}^*), J_m(p_{\mathbf{X}}^0), 1) \leq c_9 \eta_m^2$ , and  $c_7$ - $c_9$  are some positive constants. Consequently,  $d(\widehat{p}_{\mathbf{X}|\mathbf{Y}}^b, p_{\mathbf{X}|\mathbf{Y}}^0) = O_p(\eta_b + \eta_m)$  as  $n, \tilde{n} \to \infty$  under under  $p_{\mathbf{X},\mathbf{Y}}^0$ .

Theorem  $\square$  indicates that the generation error of the indirect generator is governed by the estimation errors  $\epsilon_b$  and  $\epsilon_m$  from  $\square$  and  $\square$ , and the approximation errors  $\gamma_b$  and  $\gamma_m$ , where  $\epsilon_b$  and  $\epsilon_m$  can be obtained by solving the entropy integral equations in  $\square$  and  $\square$ . Moreover, optimal rates of convergence can be obtained through tuning of  $\lambda_b$  and  $\lambda_m$ . Note that  $\eta_m$  could be tuned as a smaller order of  $\eta_b$  when the size of unlabeled data greatly exceeds the size of labeled data. Then the generalization error of the indirect method is governed primarily by  $\eta_b$ . In other words, the indirect generator's performance is mainly determined by the estimation of  $p_{Y|X}$ .

For direct generator, let  $p_{X|Y}^* \in \mathcal{F}_f$  be an approximation of  $p_{X|Y}^0$  in that  $\rho_f(p_{X|Y}^0, p_{X|Y}^*) \leq \gamma_f$ 

for some  $\gamma_f \geq 0$ .

**Theorem 2** (Direct generator). Suppose there exist some positive constants  $c_{10}$ - $c_{12}$ , such that, for any  $\epsilon_f > 0$ , and  $\lambda_f \geq 0$ ,

$$\sup_{k\geq 1} \int_{2^{-8}L_k}^{2^{1/2}L_k^{1/2}} H^{1/2}(u/c_{12}, \mathcal{F}_{f,k}) du/L_k \leq c_{11}n^{1/2}, \quad L_k = c_{10}\epsilon_f^2 + \lambda_f(k-1), \tag{9}$$

then

$$P(d(\widehat{p}_{\mathbf{X}|\mathbf{Y}}^f, p_{\mathbf{X}|\mathbf{Y}}^0) \ge \eta_f) \le 8 \exp(-c_{13}n\eta_f^2), \quad \eta_f = \max(\epsilon_f, \gamma_f^{1/2}), \tag{10}$$

provided that  $\lambda_f \max(J_f(p_{\boldsymbol{X}|\boldsymbol{Y}}^*), J_f(p_{\boldsymbol{X}|\boldsymbol{Y}}^0), 1) \leq c_9 \eta_f^2$ , and  $c_{13} > 0$  is a constant. Consequently,  $d(\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^f, p_{\boldsymbol{X}|\boldsymbol{Y}}^0) = O_p(\eta_f)$  as  $n \to \infty$  under  $p_{\boldsymbol{X},\boldsymbol{Y}}^0$ .

In contrast to the indirect generation, the generation error  $\eta_f$  of the direct generation could be much larger or smaller than that of the indirect generation  $\eta_b$  depending on the complexities of  $\mathcal{F}_f$ ,  $\mathcal{F}_b$ , and the corresponding approximation errors  $\gamma_f$  and  $\gamma_b$ , when  $p_{\boldsymbol{X}}$  can be sufficiently well estimated. This suggests that either may outperform the other depending on the model assumptions.

Note that  $\gamma_f$ ,  $\gamma_b$ , and  $\gamma_m$  are the approximation errors of the approximation capabilities of function spaces  $\mathcal{F}_f$ ,  $\mathcal{F}_b$  and  $\mathcal{F}_m$  [35, 40]. In particular, when the function space is defined by a ReLU deep neural network to approximate a function in a Sobolev space, the approximation error is available and related to the scale of the neural network [40].

Theorem 3 says that the coupled generator performs no worse than the indirect and direct generators when 5 is used to select based on an independent cross-validation sample of size N.

**Theorem 3** (Coupled generation). Under  $p_{\mathbf{X},\mathbf{Y}}^0$ , as  $N \to \infty$ , the coupled generator defined in (5) satisfies  $K(p_{\mathbf{X}|\mathbf{Y}}^0, \widehat{p}_{\mathbf{X}|\mathbf{Y}}^c) \le \min \left(K(p_{\mathbf{X}|\mathbf{Y}}^0, \widehat{p}_{\mathbf{X}|\mathbf{Y}}^b), K(p_{\mathbf{X}|\mathbf{Y}}^0, \widehat{p}_{\mathbf{X}|\mathbf{Y}}^f)\right)$ , where  $K(p_{\mathbf{X}|\mathbf{Y}}, q_{\mathbf{X}|\mathbf{Y}})$  is the Kullback-Leibler divergence between  $p_{\mathbf{X}|\mathbf{Y}}$  and  $q_{\mathbf{X}|\mathbf{Y}}$ .

Remarks: In Theorem [3], if  $K(p_{\boldsymbol{X}|\boldsymbol{Y}}^0, p_{\boldsymbol{X}|\boldsymbol{Y}}) \leq c_{14}^2 d^2(p_{\boldsymbol{X}|\boldsymbol{Y}}^0, p_{\boldsymbol{X}|\boldsymbol{Y}})$  for some constant  $c_{14} > 0$ , then  $d(p_{\boldsymbol{X}|\boldsymbol{Y}}^0, \widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^c) \leq c_{15} \min \left(d(p_{\boldsymbol{X}|\boldsymbol{Y}}^0, \widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^b), d(p_{\boldsymbol{X}|\boldsymbol{Y}}^0, \widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^c)\right)$ , which occurs when the likelihood ratio is bounded.

## 4 Sentence generation given a topic

This section derives generative methods of sentence generation, which integrates the likelihood methods developed previously with language models to compose a sentence. As a result, a new sentence can be generated, which may not appear in training data; see Table 3 for an example.

A complete sentence is represented by a word vector  $X_{1:T} = (X_1, \dots, X_T)'$ , where  $X_t$  is the t-th word, T is a sentence-specific length, and ' denotes the transpose of a vector. For convenience, we write  $X_1 =$  "START" and  $X_{T+1} =$  "END" as the null words of the first and last words of a sentence, respectively. For example,  $X_1 =$  "START",  $X_2 =$  "Football",  $X_3 =$  "is",  $X_4 =$  "a",  $X_5 =$  "popular",  $X_6 =$  "sport", and  $X_7 =$  "END". Together with  $X_{1:T}$ , its associated topic category  $Y = (Y_1, \dots, Y_K)'$  is available, where  $Y_j \in \{0, 1\}$  or  $Y_j \in \mathbb{R}$ . Finally, we construct a dictionary  $\mathcal{D} = (w_1, \dots, w_{|\mathcal{D}|})'$  to contain all composing words, that is  $X_t \in \mathcal{D}$ ;  $t = 1, \dots, T$ , with  $|\mathcal{D}|$  denoting  $\mathcal{D}$ 's size.

For simplicity, we consider the case of a fixed T, where sentences of different lengths can be processed with a fixed length, as illustrated in Table  $\mathbb{I}$ . Sentence generation given a topic Y generates a sentence  $X_{1:T+1}$  using the conditional probability  $P(X_{1:T+1} = x_{1:T+1}|Y = y)$ . However, estimation of this probability at the sentence level is infeasible. Therefore, we decompose it at the word level by the probability chain rule:

$$\log (p(\boldsymbol{X}_{1:T+1} = \boldsymbol{x}_{1:T+1} | \boldsymbol{Y} = \boldsymbol{y})) = \sum_{t=1}^{T} \log (p(X_{t+1} = x_{t+1} | \boldsymbol{X}_{1:t} = \boldsymbol{x}_{1:t}, \boldsymbol{Y} = \boldsymbol{y})).$$
(11)

This decomposition (11) permits sequential generation of a sentence through next-word-prediction given existing words by learning  $p(X_{t+1} = x_{t+1} | \boldsymbol{x}_{1:t}, \boldsymbol{y})$  from data;  $t = 1, \dots, T$ .

Yet, estimation of  $p(X_{t+1} = x_{t+1}|\mathbf{x}_{1:t}, \mathbf{y})$  remains challenging for unstructured  $\mathbf{X}_{1:t}$  because of a lack of observations in any conditioning event of  $\mathbf{X}_{1:t}$  given  $\mathbf{Y}$  even with large training data. Furthermore, it is difficult to utilize unlabeled data to estimate  $p(X_{t+1} = x_{t+1}|\mathbf{x}_{1:t},\mathbf{y})$ .

#### 4.1 Indirect generator

In this context, we derive a version of (2) and that of (3) through (11) to estimate the inverse probabilities. Specifically,  $p(x_{t+1}|\mathbf{x}_{1:t},\mathbf{y})$  can be written as

$$p(x_{t+1}|\mathbf{x}_{1:t}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x}_{1:t+1})p(x_{t+1}|\mathbf{x}_{1:t})}{\sum_{x_{t+1} \in \mathcal{D}} p(\mathbf{y}|\mathbf{x}_{1:t}, x_{t+1})p(x_{t+1}|\mathbf{x}_{1:t})};$$
(12)

for  $t = 1, \dots, T$ . Then, we estimate the inverse probability  $p(\boldsymbol{y}|\boldsymbol{x}_{1:t+1})$  based on labeled data  $(\boldsymbol{x}_{1:t}^i, \boldsymbol{y}^i)_{i=1}^n$  and estimate  $p(x_{t+1}|\boldsymbol{x}_{1:t})$  based on  $(\boldsymbol{x}_{1:t}^i)_{i=1}^n$  for  $t = 1, \dots, T^i$ , and unlabeled data  $(\boldsymbol{x}_{1:t}^j)_{j=1}^{\tilde{n}}$  for  $t = 1, \dots, T^j$ .

Estimation of  $p(\boldsymbol{y}|\boldsymbol{x}_{1:t})$  may proceed with unstructured predictors  $\boldsymbol{x}_{1:t}$ . To proceed, we map a sentence  $\boldsymbol{x}_{1:t}$  to a numerical vector  $\mathcal{E}(\boldsymbol{x}_{1:t}) \in \mathbb{R}^p$ , known as a numerical embedding of size p via a pre-trained embedding model such as Doc2Vec [23, 24] and BERT [II]. If a pre-trained embedding model is sufficient in that  $p(\boldsymbol{y}|\boldsymbol{X}_{1:t}=\boldsymbol{x}_{1:t})=p(\boldsymbol{y}|\mathcal{E}(\boldsymbol{X}_{1:t})=\mathcal{E}(\boldsymbol{x}_{1:t}))$  [II], the numerical embedding  $\mathcal{E}(\boldsymbol{x}_{1:t})$  captures word-to-word relations expressed in terms of co-occurrence of words, which may raise the level of predictability of unstructured predictors  $\boldsymbol{X}_{1:t}$ . Next, we model  $p(\boldsymbol{y}|\boldsymbol{x}_{1:t})$  through  $p(\boldsymbol{y}|\mathcal{E}(\boldsymbol{x}_{1:t}))$  when  $\boldsymbol{Y} \in \{0,1\}^K$  is categorical or

 $\boldsymbol{Y} \in \mathbb{R}^{K}$  is continuous with an embedded label  $\boldsymbol{Y}$ :

$$p(\boldsymbol{y}|\boldsymbol{x}_{1:t}) = \begin{cases} \boldsymbol{y}' \boldsymbol{\sigma} \left( f(\mathcal{E}(\boldsymbol{x}_{1:t})) \right), & \text{if } \boldsymbol{y} \in \{0, 1\}^K, \\ (2\pi)^{-K/2} \exp\left( -\frac{1}{2} \|\boldsymbol{y} - f(\mathcal{E}(\boldsymbol{x}_{1:t}))\|_2^2 \right), & \text{if } \boldsymbol{y} \in \mathbb{R}^K, \end{cases}$$
(13)

where K is the dimension of  $\mathbf{Y}$ ,  $\boldsymbol{\sigma}(\cdot)$  is the softmax function  $[\![1\!]]$ , and f is a nonparametric classification or regression function forest  $[\![3\!]]$  or linear function  $f(\mathcal{E}(\boldsymbol{x}_{1:t})) = \boldsymbol{\theta}_b \mathcal{E}(\boldsymbol{x}_{1:t})$  with  $\boldsymbol{\theta}_b \in \mathbb{R}^{K \times p}$ . For illustration, we use a linear representation  $f(\mathcal{E}(\boldsymbol{x}_{1:t})) = \boldsymbol{\theta}_b \mathcal{E}(\boldsymbol{x}_{1:t})$  in  $[\![1\!]]$  sequentially. Now the cost function  $L_b(\boldsymbol{\theta}_b)$  in  $[\![2\!]]$  becomes

$$L_b(\boldsymbol{\theta}_b) = -\frac{1}{n} \sum_{i=1}^n (T^i)^{-1} \sum_{t=1}^{T^i} \log \left( p(\boldsymbol{y}^i | \mathcal{E}(\boldsymbol{x}_{1:t}^i)) + \lambda_b J_b(f), \right)$$
(14)

where  $\lambda_b \geq 0$  is a tuning parameter and  $J_b(f) \geq 0$  is a regularizer, for example,  $J_b(f) = \|\boldsymbol{\theta}_b\|_F^2$  if  $f(\mathcal{E}(\boldsymbol{x}_{1:t})) = \boldsymbol{\theta}_b \mathcal{E}(\boldsymbol{x}_{1:t})$ , where  $\|\cdot\|_F$  is the Frobenius-norm of a matrix.

On the other hand, the next word probability is estimated by a RNN in a form of

$$p(x_{t+1}|\mathbf{x}_{1:t}) = \mathbf{o}_{[x_{t+1}]}(x_t, \mathbf{h}_t; \mathbf{\theta}_m), \text{ with } \mathbf{h}_t = \mathbf{h}(x_t, \mathbf{h}_{t-1}; \mathbf{\theta}_m), \mathbf{h}_0 = \mathbf{0},$$
 (15)

where  $[x_{t+1}] = \{j : w_j = x_{t+1}\}$ ,  $\boldsymbol{o}_j(x_t, \boldsymbol{h}_t, \boldsymbol{\theta}_m)$  is the probability of occurrence of the j-th word in  $\mathcal{D}$ , and  $\boldsymbol{h}(x_t, \boldsymbol{h}_{t-1}, \boldsymbol{\theta}_m)$  is a hidden state function, such as a long short-term memory unit (LSTM) [16], a bidirectional unit [32], a gated recurrent unit (GRU) [6], and GPT2 [30],  $\boldsymbol{\theta}_m$  is the parameter of a specific RNN model, for example,  $\boldsymbol{\theta}_m = (\boldsymbol{W}_m^o, \boldsymbol{W}_m^x, \boldsymbol{W}_m^h)$  in a basic RNN,

$$\boldsymbol{o}(x_t, \boldsymbol{h}_t, \boldsymbol{\theta}_m) = \boldsymbol{\sigma}(\boldsymbol{W}_m^o \boldsymbol{h}_t), \ \boldsymbol{h}_t = \phi(\boldsymbol{W}_m^x \mathbf{1}_{[x_t]} + \boldsymbol{W}_m^h \boldsymbol{h}_{t-1}), \text{ and } \boldsymbol{h}_0 = \boldsymbol{0},$$
 (16)

where  $\sigma(\cdot)$ , as defined before, is the softmax and  $\phi$  is an activation function such as the ReLU function  $\coprod$ ,  $\mathbf{W}_m^o \in \mathbb{R}^{r_m \times |\mathcal{D}|}$ ,  $\mathbf{W}_m^x \in \mathbb{R}^{|\mathcal{D}| \times r_m}$ , and  $\mathbf{W}_m^h \in \mathbb{R}^{r_m \times r_m}$ , and  $r_m$  is the number of latent factors of the RNN. See Figure  $\coprod$  for a display of the architecture of a basic RNN.

On the ground of (15), the cost function  $L_m(\boldsymbol{\theta}_m)$  in (3) becomes

$$L_{m}(\boldsymbol{\theta}_{m}) = -(n+\tilde{n})^{-1} \sum_{i=1}^{n} (T^{i})^{-1} \sum_{t=1}^{T^{i}} \log \left( \boldsymbol{o}_{[x_{t+1}^{i}]}(x_{t}^{i}, \boldsymbol{h}_{t}^{i}, \boldsymbol{\theta}_{m}) \right)$$
$$-(n+\tilde{n})^{-1} \sum_{j=1}^{\tilde{n}} (T^{j})^{-1} \sum_{t=1}^{T^{j}} \log \left( \boldsymbol{o}_{[x_{t+1}^{j}]}(x_{t}^{j}, \boldsymbol{h}_{t}^{j}, \boldsymbol{\theta}_{m}) \right) + \lambda_{m} J_{m}(\boldsymbol{\theta}_{m}), \qquad (17)$$

where  $\lambda_m \geq 0$  is a tuning parameter and  $J_m(\boldsymbol{\theta}_m)$  is a regularizer regularizing the weights matrix and the activation layer [22].

Minimizing (14) and (17) yields estimators  $\widehat{\boldsymbol{\theta}}_b$  and  $\widehat{\boldsymbol{\theta}}_m$ , respectively. Then the conditional probabilities are estimated as  $\widehat{p}(\boldsymbol{y}|\boldsymbol{x}_{1:t+1};\widehat{\boldsymbol{\theta}}_b) = \sigma(\widehat{\boldsymbol{\theta}}_b \mathcal{E}(\boldsymbol{x}_{1:t+1}))$  and  $\widehat{p}(x_{t+1}|\boldsymbol{x}_{1:t};\widehat{\boldsymbol{\theta}}_m) = \boldsymbol{o}_{[x_{t+1}]}(x_t,\boldsymbol{h}_{t-1},\widehat{\boldsymbol{\theta}}_m)$ . Plugging these estimates into (12), we obtain the estimated probability, and the process is summarized as,

$$\widehat{p}^{b}(X_{t+1} = x | \boldsymbol{x}_{1:t}, \boldsymbol{y}) = \frac{\widehat{p}(\boldsymbol{y} | \boldsymbol{x}_{1:t}, x; \widehat{\boldsymbol{\theta}}_{b}) \widehat{p}(X_{t+1} = x | \boldsymbol{x}_{1:t}; \widehat{\boldsymbol{\theta}}_{m})}{\sum_{x \in \mathcal{D}} \widehat{p}(\boldsymbol{y} | \boldsymbol{x}_{1:t}, x; \widehat{\boldsymbol{\theta}}_{b}) \widehat{p}(X_{t+1} = x | \boldsymbol{x}_{1:t}; \widehat{\boldsymbol{\theta}}_{m})}$$

$$\widehat{\boldsymbol{\theta}}_{b} = \underset{\boldsymbol{\theta}_{b}}{\operatorname{argmin}} L_{b}(\boldsymbol{\theta}_{b}), \qquad \widehat{\boldsymbol{\theta}}_{m} = \underset{\boldsymbol{\theta}_{m}}{\operatorname{argmin}} L_{m}(\boldsymbol{\theta}_{m}).$$

$$(18)$$

Then, a sentence is sequentially generated as follows:

$$\widehat{x}_{t+1} = \operatorname*{argmax}_{x \in \mathcal{D}} \widehat{p}^b(X_{t+1} = x | \boldsymbol{X}_{1:t} = \boldsymbol{x}_{1:t}, \boldsymbol{Y} = \boldsymbol{y}); \quad t = 1, \dots, \widehat{T}.$$
(19)

This generation process begins with  $\hat{x}_1 = \text{"START"}$  or pre-specified  $t_0$ -words  $\hat{x}_{1:t_0}$  and proceeds until  $\hat{x}_{\widehat{T}} = \text{"END"}$  is reached, where  $\widehat{T}$  is an index at termination. It is worth mentioning that the denominator in (18) normalizes the probability but may not need to be

computed when a maximizer of (18) is desired in (19).

#### 4.2 Direct generator

The direct generation is inspired by a conditional RNN (C-RNN; [37, 17]) by estimating

$$p(x_{t+1}|\boldsymbol{x}_{1:t},\boldsymbol{y}) = \boldsymbol{o}_{[x_{t+1}]}(x_t,\boldsymbol{h}_{t-1},\boldsymbol{y},\boldsymbol{\theta}_f), \text{ with } \boldsymbol{h}_t = \boldsymbol{h}(x_t,\boldsymbol{h}_{t-1},\boldsymbol{\theta}_f), \boldsymbol{h}_0 = \boldsymbol{h}_0(\boldsymbol{y},\boldsymbol{\theta}_f),$$
(20)

where  $\boldsymbol{\theta}_{\mathrm{f}}$  represents the parameters of a RNN, and  $\boldsymbol{h}_{0}$  is built on the label information as opposed to  $\boldsymbol{h}_{0} = \boldsymbol{0}$  in (16). As in (16), the direct generator requires additional parameters  $\boldsymbol{W}_{\mathrm{f}}^{y} \in \mathbb{R}^{r_{\mathrm{f}} \times K}$  for  $\boldsymbol{\theta}_{\mathrm{f}} = (\boldsymbol{W}_{\mathrm{f}}^{o}, \boldsymbol{W}_{\mathrm{f}}^{x}, \boldsymbol{W}_{\mathrm{f}}^{h}, \boldsymbol{W}_{\mathrm{f}}^{y})$  to model the effect from  $\boldsymbol{y}$  as follows:

$$o(x_t, h_t, y, \theta_f) = \sigma(W_f^o h_t), h_t = \phi(W_f^x \mathbf{1}_{[x_t]} + W_f^h h_{t-1}), \text{ and } h_0 = \phi(W_f^y y),$$
 (21)

where  $\mathbf{W}_{\mathrm{f}}^{o} \in \mathbb{R}^{r_{\mathrm{f}} \times |\mathcal{D}|}$ ,  $\mathbf{W}_{\mathrm{f}}^{x} \in \mathbb{R}^{|\mathcal{D}| \times r_{\mathrm{f}}}$ ,  $\mathbf{W}_{\mathrm{f}}^{h} \in \mathbb{R}^{r_{m} \times r_{\mathrm{f}}}$ , and  $r_{\mathrm{f}}$  is the number of latent factors of the RNN. On this ground, the cost function in (4) becomes

$$L_{\mathrm{f}}(\boldsymbol{\theta}_{\mathrm{f}}) = -n^{-1} \sum_{i=1}^{n} (T^{i})^{-1} \sum_{t=1}^{T_{i}} \log \left( \boldsymbol{o}_{[x_{t+1}^{i}]}(x_{t}^{i}, \boldsymbol{h}_{t-1}^{i}, \boldsymbol{y}^{i}, \boldsymbol{\theta}_{\mathrm{f}}) \right) + \lambda_{\mathrm{f}} J_{\mathrm{f}}(\boldsymbol{\theta}_{\mathrm{f}}), \tag{22}$$

where  $\lambda_{\rm f} \geq 0$  is a tuning parameter and  $J_{\rm f}(\boldsymbol{\theta}_{\rm f})$  is a nonnegative regularizer. Minimizing (22) in  $\boldsymbol{\theta}_{\rm f}$  yields an estimate  $\hat{\boldsymbol{\theta}}_{\rm f}$ , thus the estimated probability  $\hat{p}^{\rm f}(x|\boldsymbol{x}_{1:t},\boldsymbol{y}) = \boldsymbol{o}_{[x]}(x_t,\boldsymbol{h}_{t-1},\boldsymbol{y},\hat{\boldsymbol{\theta}}_{\rm f})$ , from (20). Then, sentence generation proceeds as in (19).

Worth of note is that the direct and indirect generators can be respectively implemented using different RNN models, for example, GPT2 for the direct RNN (20) while LSTM for the indirect RNN in (15). Moreover, different model architectures of RNNs may yield different empirical results. This aspect is illustrated in Section 5.

#### 4.3 Coupled generator

Given the estimated probabilities  $\hat{p}^f(x_{t+1}|\boldsymbol{x}_{1:t},\boldsymbol{y})$  and  $\hat{p}^b(x_{t+1}|\boldsymbol{x}_{1:t},\boldsymbol{y})$ . The coupled generator chooses one between  $\hat{p}^f$  and  $\hat{p}^b$  on a validation set by minimizing an empirical version of the log-likelihood loss,

$$\operatorname{Ent}(\widehat{p}) = -T^{-1} \sum_{t=1}^{T} \log \widehat{p}(X_{t+1} = x_{t+1} | \boldsymbol{X}_{1:t} = \boldsymbol{x}_{1:t}, \boldsymbol{Y} = \boldsymbol{y}).$$
 (23)

#### 4.4 Large-scale computation

This section develops a computational scheme for the indirect generator in (14)-(17), and the direct generator in (22) can be treated by a standard RNN implementation as in [36, 29]. In particular, when stochastic backpropagation is used through the time gradient method, the computation complexity is of order of the number of parameters per time step [27].

In what follows, we apply gradient descent [41] or stochastic gradient descent [31] to solve [14]. For (17), we apply a classical back-propagation algorithm. In each case, we use analytic a gradient expression for updates.

Gradient for indirect generation. The gradient expression for  $\theta_m$  in (17) is given in [29], while that for  $\theta_b$  in (14) is computed as

$$\frac{\partial L_b}{\partial \boldsymbol{\theta}_{b,k}} = \begin{cases}
\lambda_b \boldsymbol{\theta}_{b,k} - n^{-1} \sum_{i=1}^n (T^i)^{-1} \sum_{t=1}^{T^i} \left( y_k^i - \sigma_k(\boldsymbol{\theta}_b \mathcal{E}(\boldsymbol{x}_{1:t}^i)) \mathcal{E}(\boldsymbol{x}_{1:t}^i), \quad \boldsymbol{y}^i \in \{0,1\}^K, \\
\lambda_b \boldsymbol{\theta}_{b,k} - n^{-1} \sum_{i=1}^n (T^i)^{-1} \sum_{t=1}^{T^i} \left( y_k^i - \boldsymbol{\theta}_{b,k}^\prime \mathcal{E}(\boldsymbol{x}_{1:t}^i) \right) \mathcal{E}(\boldsymbol{x}_{1:t}^i), \quad \boldsymbol{y}^i \in \mathbb{R}^K,
\end{cases} (24)$$

where  $\boldsymbol{\theta}_{b,k}$  denotes the kth column of  $\boldsymbol{\theta}_b$ .

The detail of gradient descent for the indirect generator is summarized as follows.

#### Algorithm 1 (Gradient descent for indirect generator):

Step 1 (Initialization): Specify a RNN architecture, randomly initialize  $\theta_b$  and  $\theta_m$ , a step size for gradient descent, tuning parameters  $\lambda_b$  and  $\lambda_m$ , and the number of maximal

training iterations.

**Step 2** (Regression): Updating  $\theta_b$  in (14) based on the gradient in (24).

Step 3 (RNN Gradient update): Solving (17) by updating the indirect RNN parameters  $\theta_m$  via gradient descent based on back-propagation in [29].

Step 4 (Termination): Iterate Step 3 until a stopping criterion is met.

Algorithm 1 can be updated by a stochastic gradient scheme [2]. Lemma [1] describes computational properties of Algorithm 1.

**Lemma 1.** If the cost functions  $L_b$  in (14) and  $L_m$  in (17) are continuously twice differential, and the probability measure of random initialization is absolutely continuous with respect to the Lebesgue measure. Then,  $\hat{\theta}_b$  is a global minimizer of (14), while  $\hat{\theta}_m$  is a local minimizer of (17) almost surely, provided that the step size in Algorithm 1 is sufficiently small.

#### 4.5 Theory for sentence generation

This section generalizes the theoretical result of Section 3 to the problem of next-word-prediction.

Now we use  $p_{X|Y}$ ,  $p_{Y|X}$ , and  $p_X$  to respectively represent  $\{p_{X_{t+1}|X_{1:t},Y}\}_{t=1}^T$ ,  $\{p_{Y|X_{1:t}}\}_{t=1}^T$ , and  $\{p_{X_{t+1}|X_{1:t}}\}_{t=1}^T$ . The expected square Hellinger-distance for next-word-prediction is

$$d(\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}, \boldsymbol{q}_{\boldsymbol{X}|\boldsymbol{Y}}) = \left(\bar{\mathbb{E}}h^2(p_{X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y}}, q_{X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y}})\right)^{\frac{1}{2}},$$
(25)

where  $\bar{\mathbb{E}}(\cdot) = T^{-1} \sum_{t=1}^{T} \mathbb{E}_{\boldsymbol{X}_{1:t}, \boldsymbol{Y}}(\cdot)$ .

The metric entropy of  $\mathcal{F}_{b,k}$  is defined by a distance  $\kappa^2(\boldsymbol{p}_{\boldsymbol{Y}|\boldsymbol{X}},\boldsymbol{q}_{\boldsymbol{Y}|\boldsymbol{X}}) = \bar{\mathbb{E}}h^2(p_{\boldsymbol{Y}|\boldsymbol{X}_{1:t}},q_{\boldsymbol{Y}|\boldsymbol{X}_{1:t}})$ . Similarly,  $\kappa^2(\boldsymbol{p}_{\boldsymbol{X}},\boldsymbol{q}_{\boldsymbol{X}}) = \bar{\mathbb{E}}h^2(p_{X_{t+1}|\boldsymbol{X}_{1:t}},q_{X_{t+1}|\boldsymbol{X}_{1:t}})$ , and  $d^2(\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}},\boldsymbol{q}_{\boldsymbol{X}|\boldsymbol{Y}})$  are used for  $\mathcal{F}_{m,k}$  and  $\mathcal{F}_{f,k}$ , respectively.

The approximation error for  $\boldsymbol{p}_{Y|X}^0$  is  $\rho_b(\boldsymbol{p}_{Y|X}^0, \boldsymbol{p}_{Y|X}) = \bar{\mathbb{E}}g_{\alpha}(\frac{p^0(Y|X_{1:t+1})}{p(Y|X_{1:t+1})})$ . Similarly, the approximation errors  $\rho_m(\boldsymbol{p}_X^0, \boldsymbol{p}_X) = \bar{\mathbb{E}}g_{\alpha}(\frac{p^0(X_{t+1}|X_{1:t})}{p(X_{t+1}|X_{1:t})})$  and  $\rho_f(\boldsymbol{p}_{X|Y}^0, \boldsymbol{p}_{X|Y}) = \bar{\mathbb{E}}g_{\alpha}(\frac{p^0(X_{t+1}|X_{1:t},Y)}{p(X_{t+1}|X_{1:t},Y)})$ 

are used for  $p_X^0$  and  $p_{X|Y}^0$ .

Corollary 1 (Sequential generation). All the results in Theorems  $\boxed{1}$  and  $\boxed{2}$  continue to hold with the distance  $d(\cdot, \cdot)$  defined in  $\boxed{25}$ .

Next we provide a theoretical example to illustrate Corollary 1

Theoretical example. Suppose that the RNN in (15) is a basic recurrent network with  $\boldsymbol{\theta}_m = (\boldsymbol{W}_m^o, \boldsymbol{W}_m^x, \boldsymbol{W}_m^h)$ , that is,  $\boldsymbol{o}(x_t, \boldsymbol{h}_{t-1}, \boldsymbol{\theta}_m) = \boldsymbol{\sigma}(\boldsymbol{W}_m^o \boldsymbol{h}_{t-1})$ ,  $\boldsymbol{h}_t = \phi(\boldsymbol{W}_m^x \mathbf{1}_{[x_t]} + \boldsymbol{W}_m^h \boldsymbol{h}_{t-1})$ , and  $\boldsymbol{h}_0 = \mathbf{0}_{r_m}$ , where  $\boldsymbol{W}_m^o \in \mathbb{R}^{r_m \times |\mathcal{D}|}$ ,  $\boldsymbol{W}_m^x \in \mathbb{R}^{|\mathcal{D}| \times r_m}$ , and  $\boldsymbol{W}_m^h \in \mathbb{R}^{r_m \times r_m}$ ,  $r_m$  is the number of latent factors of the RNN, and  $\phi(\boldsymbol{z})$  is an activation function, such as the sigmoid function  $\phi(\boldsymbol{z}) = 1/(1 + \exp(-\boldsymbol{z}))$ , the tanh function  $\phi(\boldsymbol{z}) = \tanh(\boldsymbol{z})$ , and the Rectified linear unit (ReLU)  $\phi(\boldsymbol{z}) = \boldsymbol{z}_+$ . For illustration, we focus on the sigmoid function.

The RNN in (20) is that  $o(x_t, h_{t-1}, \theta_f) = \sigma(W_f^o h_{t-1})$ ,  $h_t = \phi(W_f^x \mathbf{1}_{[x_t]} + W_f^h h_{t-1})$ , and  $h_0 = \phi(W_f^y y)$ . The network parameters are  $\theta_f = (W_f^o, W_f^x, W_f^h, W_f^y)$ , where  $W_f^o \in \mathbb{R}^{r_f \times |\mathcal{D}|}$ ,  $W_f^x \in \mathbb{R}^{|\mathcal{D}| \times r_f}$ ,  $W_f^h \in \mathbb{R}^{r_f \times r_f}$ , and  $W_f^y \in \mathbb{R}^{r_f \times K}$ , and  $r_f$  is the number of latent factors of the RNN in the direct generation.

Corollary 2 gives the generation errors of the direct and indirect generators.

Corollary 2 (Theoretical example). For the estimated next-word probabilities  $\widehat{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}}^{\mathrm{f}}$  by the direct generator in (22), we have that  $d(\widehat{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}}^{\mathrm{f}}, \boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) = O_{p}(\eta_{\mathrm{f}})$ , where

$$\eta_{\rm f} = \max\left\{ \left(\frac{\Lambda_{\rm f}}{n}\log\left(\frac{n\max(r_{\rm f},2c_{15})2^TT^{-1/2}}{\Lambda_{\rm f}}\right)\right)^{\frac{1}{2}}, \gamma_{\rm f}^{\frac{1}{2}}\right\},\,$$

 $\Lambda_{\rm f} = r_{\rm f}(2|\mathcal{D}| + r_{\rm f} + K), \ \lambda_{\rm f} = c_{17}\eta_{\rm f}^2, \ and \ c_{15} > 0 \ and \ c_{16} > 0 \ are \ constants \ with \mathbb{E}_{\boldsymbol{Y}} \|\boldsymbol{Y}\|_2^2 \leq c_{15}.$ Similarly, the estimated next-word probabilities  $\hat{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}}^b$  by the indirect generator in (14) and

(17) satisfies:  $d(\widehat{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}}^b, \boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^0) = O_p(\eta_b + \eta_m)$ , where

$$\eta_b = \max\left\{ \left( \frac{\Lambda_b}{n} \log \left( \frac{c_{16}n}{\Lambda_b} \right) \right)^{\frac{1}{2}}, \gamma_b^{\frac{1}{2}} \right\},$$

$$\eta_m = \max\left\{ \left( \frac{\Lambda_m}{n+\tilde{n}} \log \left( \frac{r_m(n+\tilde{n})2^T T^{-1/2}}{\Lambda_m} \right) \right)^{\frac{1}{2}}, \gamma_m^{\frac{1}{2}} \right\},$$

 $\Lambda_b = Kp, \ \Lambda_m = r_m(2|\mathcal{D}| + r_m), \ \lambda_b = c_{18}\eta_b^2, \ \lambda_m = c_{18}\eta_m^2, \ and \ c_{16} > 0 \ and \ c_{18} > 0 \ are constants with \ \bar{\mathbb{E}} \|\mathcal{E}(\boldsymbol{X}_{1:t})\|_2^2 \leq c_{16}.$ 

Corollary 2 says that the generation error of the indirect generator in (1) becomes  $d(\widehat{p}_{X|Y}^b, p_{X|Y}^0) = O_p(\frac{\Lambda_b}{n}\log(\frac{n}{\Lambda_b}))^{\frac{1}{2}}$  when  $(\tilde{n}+n)/n = O(\frac{\Lambda_m \log(r_m n/\Lambda_m)}{\Lambda_b \log(c_1 n/\Lambda_b)})$ , when  $\gamma_b = \gamma_m = 0$ . In fact, the generation error is primarily dominated by its estimation error of  $p(Y|X_{1:t})$ , because  $p(X_{t+1}|X_{1:t})$  can be well estimated with the help of large unlabeled data with  $\tilde{n} \gg n$ . In this situation, the indirect method outperforms the direct method, particularly when  $\Lambda_b < \Lambda_f$ , suggesting that the estimation complexity of the indirect method is less than that of the direct method. Interestingly, the generation error of the direct generator agrees with that of the maximum likelihood estimates under the Hellinger-distance [34, 38]. With respect to tuning, a large value of  $\Lambda_b$ ,  $\Lambda_m$ , and  $\Lambda_f$  increases the complexity of the corresponding functional space for probability estimation, thus reduces the generation errors. Consequently, the generation errors of the direct and indirect generators indeed depend on the model complexity of parameter spaces  $\mathcal{F}_b$  and  $\mathcal{F}_f$ .

To illustrate the synergy of indirect and direct generators' respective strengths, we consider two situations. First,  $d(\hat{p}_{X|Y}^f, p_{X|Y}^0) = o_p(1)$  but  $d(\hat{p}_{X|Y}^b, p_{X|Y}^0)$  is bounded away from zero if an unlabeled sample  $(X_{1:T^j}^j)_{j=1}^{\tilde{n}}$  follows a different marginal distribution of the labeled sample  $(X_{1:T^i}^i)_{i=1}^n$ . Second,  $d(\hat{p}_{X|Y}^b, p_{X|Y}^0) = o_p(1)$  but  $d(\hat{p}_{X|Y}^f, \hat{p}_{X|Y}^0)$  is bounded away from zero in the presence of a new word in labeled but unlabeled samples. However, in both situations,  $d(\hat{p}_{X|Y}^c, p_{X|Y}^0) = o_p(1)$ , when Kullback-Leibler divergence is equivalent to the Hellinger-distance. In other words, only the coupled generator has a generation error

tending to zero in both situations.

### 5 Benchmark

This section examines the performance of the coupled, indirect, and direct generators in one benchmark example, and compares with a baseline method "Separate RNN", which fits RNNs for each topic as in [36]. The benchmark concerns sentence categorization based on a text corpus in the UCI machine learning repository. This corpus contains 1,039 labeled sentences collected from abstracts and introductions of 30 articles, in which five topic categories are AIM (a specific aim of the present paper), OWN (description of own work presented in the present paper), CONTRAST (comparison statements with other works, including strengths and weaknesses), BASIS (statements of agreement with other works or continuation of other works), and MISC (generally accepted scientific background or description of other works). These labels originate from three scientific domains: computational biology (PLOS), the machine learning repository on arXiv (ARXIV), and the psychology journal judgment and decision making (JDM). For example, a typical sentence such as "The instantaneous loss bound of SYMBOL implies only convergence in probability." is labeled as "MISC" according to scientific topic classification. In addition to the aforementioned labeled sentences, this corpus contains 34,481 unlabeled sentences from 300 articles in PLOS, ARXIV, and JDM.

Before proceeding, we pre-process the text corpus to filter out redundant each sentence's component so that numerical embeddings are applied for the indirect generator. First, we replace all numerical values, symbolic values, and citations by "NUMBER", "SYMBOL", and "CITATION", respectively, and remove all standalone punctuation marks except commas, periods, and semicolons. For unlabeled sentences, we remove words appearing less than 20 times in the corpus, which leads to an unlabeled corpus of 8,286 sentences. On this ground,

<sup>&</sup>lt;sup>2</sup>https://archive.ics.uci.edu/ml/datasets/Sentence+Classification

a dictionary is constructed, consisting of 5,369 words extracted from labeled and unlabeled sentences.

For training, we generate word strings for next-word-prediction based on the maximum length of all sentences in the dataset. Thus, all the previous tokens in the sentence are contributing to predict the next word. Specifically, we create next-word-prediction sequences consisting of consecutive previous words and fill with the null word "NULL" to pad all word strings as the same length. An example of such a next-word-prediction sequence is displayed in Table [I]. In this fashion, we gather enough training sentences as the null words do not impact our learning process. Now, 28,180 labeled next-word-prediction sequences are generated from the original 1,039 labeled sentences, together with 174,355 unlabeled sequences from the original 8,286 unlabeled sentences.

The generation performance is measured by two commonly used metrics, namely, the next-word entropy loss and the Bi-Lingual Evaluation Understudy (BLEU) loss [28] over a test sample, approximating the predicted Kullback-Leibler divergence and Jaccard distance [13], respectively. Given sentences  $(\hat{x}_{1:\hat{T}^i}^i)_{i=1}^{n_{\text{test}}}$  generated from  $\hat{p}$  and its reference sentence  $(x_{1:T^i}^i)_{i=1}^{n_{\text{test}}}$  given a topic y, the entropy loss is defined in as the empirical version of [23], while the BLEU<sub>l</sub>-loss  $(l=1,\cdots,4)$  can be written as

$$\mathrm{BLEU}_{l}\text{-loss}(\widehat{\boldsymbol{p}}) = 1 - n_{\mathrm{test}}^{-1} \sum_{i=1}^{n_{\mathrm{test}}} \exp\big(\min(1 - \frac{\widehat{T}^{i}}{T^{i}}, 0)\big) \frac{|\mathrm{gram}_{l}(\boldsymbol{x}_{1:T^{i}}^{i}) \cap \mathrm{gram}_{l}(\widehat{\boldsymbol{x}}_{1:\widehat{T}^{i}}^{i})|}{|\mathrm{gram}_{l}(\widehat{\boldsymbol{x}}_{1:\widehat{T}^{i}}^{i})|},$$

where  $n_{\text{test}}$  is the number of sentences in the testing set,  $|\cdot|$  denotes a set's size and  $\text{gram}_l(\cdot)$  is the l-gram set for a sentence. For a sentence "the cat in the hat", its 1-gram set is { "the", "cat", "in", "the", "hat"}, the 2-gram set is { "the cat", "cat in", "in the", "the hat"}, and the 3-gram set is { "the cat in", "cat in the", "in the hat" }. The BLEU $_l$ -loss can be computed using the NLTK library in Python. Whereas the entropy loss measures the occurrence probability of the reference sentences, the BLEU $_l$ -loss focuses on exact matching between l consecutive words of two sentence. Moreover, we also consider the SF-BLEU $_l$ -loss

to evaluate the diversity of a generated sentences [43], defined as

$$\text{SF-BLEU}_{l}\text{-loss}(\widehat{\boldsymbol{p}}) = 1 - n_{\text{test}}^{-1} \sum_{i=1}^{n_{\text{test}}} \exp\big(\min(1 - \frac{\widehat{T}^{i}}{\widehat{T}_{\min}^{-i}}, 0)\big) \frac{\max_{j \neq i} |\operatorname{gram}_{l}(\widehat{\boldsymbol{x}}_{1:\widehat{T}^{i}}^{i}) \cap \operatorname{gram}_{l}(\widehat{\boldsymbol{x}}_{1:\widehat{T}^{j}}^{j})|}{|\operatorname{gram}_{l}(\widehat{\boldsymbol{x}}_{1:\widehat{T}^{i}}^{i})|},$$

where  $\hat{T}_{\min}^{-i} = \operatorname{argmin}_{\hat{T}^j; j \neq i} |\hat{T}^j - \hat{T}^i|$ , and a high SF-BLEU<sub>l</sub>-loss score means more diverse.

For training, validation, and testing, we randomly split all the labeled articles into three sets with a partition ratio of 60%, 20%, and 20%, respectively. Moreover, for a sentence  $\boldsymbol{x}_{1:T}$  and its associated topic  $\boldsymbol{y}$  in a testing set, five starting words  $\boldsymbol{x}_{1:5}$  as opposed to the null word are given to predict the rest of a sentence.

Consider two situations of the semantic label: (1)  $\mathbf{Y} \in \{0,1\}^K$  is categorical and is coded as a 0-1 vector using the one-hot encoding from the topic category; (2)  $\mathbf{Y} \in \mathbb{R}^K$  is continuous with each topic as a p = 128-dimensional vector based on Doc2Vec. In (2), each topic is represented by the averaged embedding of all the sentences in this topic category in training data.

In the case of  $\mathbf{Y} \in \{0,1\}^K$ , the indirect generator involves [14] and [17]. For [14], we perform regularized multinomial logistic regression using the Python library  $\mathbf{sklearn}^3$  on the embedded next-word prediction sequence training samples  $(\mathcal{E}(\mathbf{x}_{1:t}^i), \mathbf{y}^i)$ , where  $\mathcal{E}(\mathbf{x}_{1:t}^i)$  is the numerical embedding of  $\mathrm{Doc2Ved}^3$  of size p=128 and the optimal  $\lambda_b$  is obtained by minimizing the entropy loss based on validation data over a set of grids  $\{.0001, .001, .01, .1, 1, 10, 100\}$ . For [17], the indirect RNN is trained based on both labeled and unlabeled next-word prediction sequences in training data. The indirect RNN model in [17] is structured in four layers, including an embedding layer consisting of 5,369 nodes with each node corresponding one word in the dictionary  $\mathcal{D}$ , an LSTM layer of 128 latent factors, a dense layer with output dimension 5,369. Note that the tuning parameter in [17] is fixed as  $\lambda_m = .0001$  through in the embedding layer to regularize words in the absence in a training set. Similarly, the direct

https://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.Ridge.html https://radimrehurek.com/gensim/models/doc2vec.html

generator trains the RNN model in (22), which has the same configuration as the indirect RNN expect that the input dimension is  $|\mathcal{D}| + K = 5,374$  in its embedding layer. Moreover, Separate RNN has the same structure as the indirect RNN given each topic.

As discussed in Section 4.2, different RNN model architectures may yield different empirical performance. Toward this end, we compare the LSTM architecture with GPT2 architecture for the direct RNNs. In particular, we consider the base GPT2 with 12 layers and 117M parameters [30] for the direct method, denoted as direct-GPT2. One key difference between LSTM and GPT2 lies in its masked self-attention layer, which masks future tokens and passes the attention information through tokens that are positioned at the left of the current position.

In the case of continuous  $Y \in \mathbb{R}^p$  after numerical embedding, the indirect generation proceeds as in the categorical case except that linear regression as opposed to multinomial logistic regression in (14) is performed using sklearn on the labeled embedded next-word prediction sequences in training data ( $\mathcal{E}(\boldsymbol{x}_{1:t}^i), \boldsymbol{y}^i$ ), where each  $\boldsymbol{y}^i$  is a 128-dimensional embedding vector.

All RNN models are trained using Keras with the batch and epoch sizes 200 and 100, and optimizer as Adam, and the over-fitting is prevented by early termination 4 of patience as 20. Moreover, the coupled generator is tuned as in 5.

Table 2 about here

As indicated in Table 2, when only labeled data is available, the coupled generator delivery higher accuracy than direct and indirect generators, which suggests the advantage of the proposed method. When combining with unlabeled data, the coupled generator outperforms the direct generator and separate RNN for both categorical and continuous labels, which selects the indirect generator in this case. With respect to the entropy loss, the amounts of

<sup>5</sup>https://keras.io/

improvement of the indirect generator over the separate RNN method and direct generator are 20.3% and 14.5% for the categorical case and 29.1% and 16.1% for the continuous case. With respect to BLEU<sub>1</sub>-BLEU<sub>4</sub> losses, a similar situation occurs, with the amounts of improvement vary with the best improvement around 15.6%. Concerning unlabeled data, a comparison between the indirect generator with and without unlabeled data suggests that unlabeled data indeed help to improve the performance of the indirect generation over 14.5%. Interestingly, in terms of the entropy loss, the direct generator based on fine-tuned GPT2 outperforms the direct generator and indirect generator based on LSTM without unlabeled data, while the coupled generator achieves the best performance between them. However, they perform similarly in terms of BLEU<sub>l</sub> scores. In view of the SF-BLEU<sub>l</sub> scores, sentences generated by the direct and indirect generators have a high degree of diversity. Moreover, the semantic label Y after sentence embeddings Doc2Vec performs slightly worse than its categorical counterpart for the indirect and direct generations, indicating that semantic relations or linguistics dependencies, as captured by the sentence embeddings, may not have an impact given that there are only five categories. Finally, as suggested in Table 3, an abstract generated based on the five categories is reasonable except for three grammatical errors that are correctable by a grammar checker<sup>6</sup>.

# Supplementary Materials

The supplementary materials provide Python codes used in real data application.

### Acknowledgments

The authors thank the editors, the associate editor, and two anonymous referees for helpful comments and suggestions.

http://www.grammarcheckforsentence.com/

## **Appendix**

**Proof of Lemma** 1. Note that  $L_b(\boldsymbol{\theta}_b)$  in (14) is convex in  $\boldsymbol{\theta}_b$  and  $L_b(\boldsymbol{\theta}_b)$  and  $L_m(\boldsymbol{\theta}_m)$  in (17) are continuously twice-differential. Then the result follows from Theorem 4 of [21]. This completes the proof.

Proof of Theorem 1. Note that  $\hat{p}_{Y}^{b}(y) = \int \hat{p}_{Y|X}^{b}(y|x)\hat{p}_{X}^{b}(x)dx$  and

$$d^2(\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^b, p_{\boldsymbol{X}|\boldsymbol{Y}}^0) = \int \left( \left( \frac{p_{\boldsymbol{Y}}^0(\boldsymbol{y}) \widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^b(\boldsymbol{y}|\boldsymbol{x}) \widehat{p}_{\boldsymbol{X}}^b(\boldsymbol{x})}{\widehat{p}_{\boldsymbol{Y}}^b(\boldsymbol{y})} \right)^{1/2} - \left( p_{\boldsymbol{Y}}^0(\boldsymbol{y}) p_{\boldsymbol{X}|\boldsymbol{Y}}^0(\boldsymbol{x}|\boldsymbol{y}) \right)^{1/2} \right)^2 d\boldsymbol{x} d\boldsymbol{y}.$$

Furthermore,  $\int \widehat{p}_{Y|X}^b(y|x)dy = 1$ . It follows from the triangular inequality that

$$d(\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{b}, p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) \leq \left(\int \left(\left(\frac{p_{\boldsymbol{Y}}^{0}(\boldsymbol{y})\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{b}(\boldsymbol{y}|\boldsymbol{x})\widehat{p}_{\boldsymbol{X}}^{b}(\boldsymbol{x})}{\widehat{p}_{\boldsymbol{Y}}^{b}(\boldsymbol{y})}\right)^{1/2} - \left(\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{b}(\boldsymbol{y}|\boldsymbol{x})\widehat{p}_{\boldsymbol{X}}^{b}(\boldsymbol{x})\right)^{1/2}\right)^{2} d\boldsymbol{x} d\boldsymbol{y}\right)^{1/2} + \left(\int \left(\left(\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{b}(\boldsymbol{y}|\boldsymbol{x})\widehat{p}_{\boldsymbol{X}}^{b}(\boldsymbol{x})\right)^{1/2} - \left(\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{b}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2}\right)^{2} d\boldsymbol{x} d\boldsymbol{y}\right)^{1/2} + \left(\int \left(\left(\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{b}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2} - \left(p_{\boldsymbol{Y}|\boldsymbol{X}}^{0}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2}\right)^{2} d\boldsymbol{x} d\boldsymbol{y}\right)^{1/2} + \left(\int \left(\left(\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{b}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2} - \left(p_{\boldsymbol{Y}|\boldsymbol{X}}^{0}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2}\right)^{2} d\boldsymbol{x} d\boldsymbol{y}\right)^{1/2} + \left(\int \left(p_{\boldsymbol{Y}|\boldsymbol{X}}^{0}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2} + \left(\int \left(p_{\boldsymbol{Y}|\boldsymbol{X}}^{0}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2}\right)^{2} d\boldsymbol{x} d\boldsymbol{y}\right)^{1/2} + \left(\int \left(p_{\boldsymbol{Y}|\boldsymbol{X}}^{0}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2} d\boldsymbol{x} d\boldsymbol{y}\right)^{1/2} d\boldsymbol{x} d\boldsymbol{y}\right)^{1/2} + \left(\int \left(p_{\boldsymbol{X}|\boldsymbol{X}}^{0}(\boldsymbol{y}|\boldsymbol{x})p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{1/2} d\boldsymbol{x} d\boldsymbol{y}\right)^{1/2} d\boldsymbol{x} d\boldsymbol{y}$$

Note that  $\widehat{p}_{\boldsymbol{X},\boldsymbol{Y}}^b(\boldsymbol{x},\boldsymbol{y}) = \widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^b(\boldsymbol{y}|\boldsymbol{x})\widehat{p}_{\boldsymbol{X}}^b(\boldsymbol{x})$ . By the triangle inequality,

$$h(p_{\boldsymbol{Y}}^{0}, \widehat{p}_{\boldsymbol{Y}}^{b}) = \left(\int \left(\left(\int p_{\boldsymbol{X},\boldsymbol{Y}}^{0}(\boldsymbol{x},\boldsymbol{y})d\boldsymbol{x}\right)^{\frac{1}{2}} - \left(\int \widehat{p}_{\boldsymbol{X},\boldsymbol{Y}}^{b}(\boldsymbol{x},\boldsymbol{y})d\boldsymbol{x}\right)^{\frac{1}{2}}\right)^{2}d\boldsymbol{y}\right)^{\frac{1}{2}}$$

$$\leq \left(\int \left(\left(p_{\boldsymbol{X},\boldsymbol{Y}}^{0}(\boldsymbol{x},\boldsymbol{y})\right)^{\frac{1}{2}} - \left(\widehat{p}_{\boldsymbol{X},\boldsymbol{Y}}^{b}(\boldsymbol{x},\boldsymbol{y})\right)^{\frac{1}{2}}\right)^{2}d\boldsymbol{x}d\boldsymbol{y}\right)^{\frac{1}{2}}$$

$$\leq \left(\int \left(\left(p_{\boldsymbol{X},\boldsymbol{Y}}^{0}(\boldsymbol{x},\boldsymbol{y})\right)^{\frac{1}{2}} - \left(p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{b}(\boldsymbol{y}|\boldsymbol{x})\right)^{\frac{1}{2}}\right)^{2}d\boldsymbol{x}d\boldsymbol{y}\right)^{\frac{1}{2}}$$

$$+ \left(\int \left(\left(\widehat{p}_{\boldsymbol{X}}^{b}(\boldsymbol{x})\right)^{\frac{1}{2}} - \left(p_{\boldsymbol{X}}^{0}(\boldsymbol{x})\right)^{\frac{1}{2}}\right)^{2}d\boldsymbol{x}\right)^{\frac{1}{2}} \leq h(p_{\boldsymbol{X}}^{0},\widehat{p}_{\boldsymbol{X}}^{b}) + \left(\mathbb{E}\left(h^{2}\left(p_{\boldsymbol{Y}|\boldsymbol{X}}^{0},\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{b}\right)\right)\right)^{\frac{1}{2}}.$$

Hence,  $d(\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^b, p_{\boldsymbol{X}|\boldsymbol{Y}}^0) \leq 2(h(\widehat{p}_{\boldsymbol{X}}^b, p_{\boldsymbol{X}}^0) + (\mathbb{E}(h^2(\widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}}^b, p_{\boldsymbol{Y}|\boldsymbol{X}}^0)))^{1/2})$ . Consequently,

$$P(d(\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{b}, p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) \geq 2(\eta_{b} + \eta_{m}))$$

$$\leq P(h(p_{\boldsymbol{X}}^{0}, \widehat{p}_{\boldsymbol{X}}) \geq \eta_{m}) + P((\mathbb{E}_{\boldsymbol{X}}(h^{2}(p_{\boldsymbol{Y}|\boldsymbol{X}}^{0}, \widehat{p}_{\boldsymbol{Y}|\boldsymbol{X}})))^{\frac{1}{2}} \geq \eta_{b}) \equiv I_{1} + I_{2}.$$

To bound  $I_1$ , let  $I_3 = P\left((n+\tilde{n})^{-1}\sum_{i=1}^{n+\tilde{n}}\left(\log(p_{\mathbf{X}}^0(\mathbf{X}^i)) - \log(p_{\mathbf{X}}^*(\mathbf{X}^i))\right) - \lambda_m J(p_{\mathbf{X}}^0) + \lambda_m J(p_{\mathbf{X}}^*)\right) \ge c_9 \eta_m^2 / 4$ ;  $I_4 = P\left(\sup_{d_m(p,p^0) \ge \eta_m} (n+\tilde{n})^{-1}\sum_{i=1}^{n+\tilde{n}}\left(\log(\frac{p_{\mathbf{X}}(\mathbf{X}^i)}{p_{\mathbf{X}}^0(\mathbf{X}^i)} - \lambda_m J(p_{\mathbf{X}}) + \lambda_m J(p_{\mathbf{X}}^0)\right) \ge -c_9 \eta_m^2 / 4\right)$ , where  $c_9 = 1 - 2\exp(-\tau/2)/(1 - \exp(-\tau/2))^2 > 0$  is a constant defined by the truncation constant  $\tau > 0$ . Then  $I_1$  is upper bounded by

$$P\Big(\sup_{d(p,p^0) \ge \eta_m} (n+\tilde{n})^{-1} \sum_{i=1}^{n+\tilde{n}} \Big(\log(p_{\mathbf{X}}(\mathbf{X}^i)/p_{\mathbf{X}}^*(\mathbf{X}^i))\Big) - \lambda_m J(p_{\mathbf{X}}) + \lambda_m J(p_{\mathbf{X}}^*)\Big) \ge 0\Big) \le I_3 + I_4.$$

By the Markov inequality,

$$I_{3} \leq P\left((n+\tilde{n})^{-1}\sum_{i=1}^{n+\tilde{n}}\left(\log(p_{\mathbf{X}}^{0}(\mathbf{X}^{i})) - \log(p_{\mathbf{X}}^{*}(\mathbf{X}^{i}))\right) \geq c_{9}\eta_{m}^{2}/4 - \lambda_{m}J(p_{\mathbf{X}}^{*})\right)$$

$$\leq \prod_{i=1}^{n+\tilde{n}}\mathbb{E}_{\mathbf{X}}\left(\frac{p_{\mathbf{X}}^{0}(\mathbf{X}^{i})}{p_{\mathbf{X}}^{*}(\mathbf{X}^{i})}\right)^{\alpha}\exp\left(-\frac{c_{9}\alpha}{8}(n+\tilde{n})\eta_{m}^{2}\right)$$

$$\leq \left(1+\alpha\gamma_{m}\right)^{n+\tilde{n}}\exp\left(-\frac{c_{9}\alpha}{8}(n+\tilde{n})\eta_{m}^{2}\right) \leq \exp\left(-\frac{c_{9}\alpha}{8}(n+\tilde{n})\eta_{m}^{2} + (n+\tilde{n})\alpha\gamma_{m}\right).$$

By Corollary 1 of [33],  $I_4 \leq 7 \exp(-c_8(n+\tilde{n})\eta_m^2/2)$ , implying  $I_1 \leq 7 \exp(-c_7(n+\tilde{n})\eta_m^2) + \exp\left(-\frac{c_9\alpha}{8}(n+\tilde{n})\eta_m^2 + (n+\tilde{n})\alpha\gamma_m\right)$  for some constant  $c_7 > 0$ . For  $I_2$ , a similar probabilistic bound can be established by applying the same argument of Theorem [2] and switching the role of  $\boldsymbol{X}$  and  $\boldsymbol{Y}$ . This leads to  $I_2 \leq 7 \exp\left(-c_8n\eta_b^2\right) + \exp\left(-\frac{c_9\alpha}{8}n\eta_b^2 + n\alpha\gamma_b\right)$  for some constant  $c_8 > 0$ . The desired result then follows.

**Proof of Theorem 2.** Denote

$$I_{5} = P\left(\sup_{d(p,p^{0}) \geq \eta_{f}} \left(n^{-1} \sum_{i=1}^{n} \left(\log\left(\frac{p_{\boldsymbol{X}|\boldsymbol{Y}}^{(\tau)}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})}{p_{\boldsymbol{X}|\boldsymbol{Y}}^{*}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})}\right)\right) - \lambda J(p_{\boldsymbol{X}|\boldsymbol{Y}}) + \lambda J(p_{\boldsymbol{X}|\boldsymbol{Y}}^{*})\right) \geq 0\right),$$

$$I_{6} = P\left(n^{-1} \sum_{i=1}^{n} \left(\log\left(\frac{p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})}{p_{\boldsymbol{X}|\boldsymbol{Y}}^{*}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})}\right)\right) - \lambda J(p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) + \lambda J(p_{\boldsymbol{X}|\boldsymbol{Y}}^{*}) \geq c_{9}\eta_{f}^{2}/4\right).$$

By the definition of a minimizer, for any  $\eta_f > 0$ ,

$$P(d(\widehat{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{\mathrm{f}}, p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) \ge \eta_{\mathrm{f}}) \le I_5 + I_6,$$

where  $\log(\mathbb{F}_i^{(\tau)}) = \log(p_{\boldsymbol{X}|\boldsymbol{Y}}^{(\tau)}(\boldsymbol{X}^i|\boldsymbol{Y}^i)) - \log(p_{\boldsymbol{X}|\boldsymbol{Y}}^0(\boldsymbol{X}^i|\boldsymbol{Y}^i))$  and

$$p_{\boldsymbol{X}|\boldsymbol{Y}}^{(\tau)}(\boldsymbol{x}|\boldsymbol{y}) = \begin{cases} \exp(-\tau)p_{\boldsymbol{X}|\boldsymbol{Y}}^*(\boldsymbol{x}|\boldsymbol{y}), & \text{if } p_{\boldsymbol{X}|\boldsymbol{Y}}(\boldsymbol{x}|\boldsymbol{y}) < \exp(-\tau)p_{\boldsymbol{X}|\boldsymbol{Y}}^*(\boldsymbol{x}|\boldsymbol{y}), \\ p_{\boldsymbol{X}|\boldsymbol{Y}}(\boldsymbol{x}|\boldsymbol{y}), & \text{otherwise,} \end{cases}$$

is the left truncation of  $p_{\boldsymbol{X}|\boldsymbol{Y}}(\boldsymbol{x}|\boldsymbol{y})$ .

Next, we bound  $I_5$  and  $I_6$  separately. An application of the same argument as in [38] yields that

$$I_{5} \leq P\left(n^{-1}\sum_{i=1}^{n}\left(\log(p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})) - \log(p_{\boldsymbol{X}|\boldsymbol{Y}}^{*}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i}))\right) \geq c_{9}\eta_{f}^{2}/4 - \lambda_{f}J(p_{\boldsymbol{X}|\boldsymbol{Y}}^{*})\right)$$

$$\leq P\left(\prod_{i=1}^{n}\left(\frac{p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})}{p_{\boldsymbol{X}|\boldsymbol{Y}}^{*}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})}\right)^{\alpha} \geq \exp\left(\frac{c_{9}\alpha}{8}n\eta_{f}^{2}\right)\right)$$

$$\leq \prod_{i=1}^{n}\mathbb{E}_{\boldsymbol{Y}}\mathbb{E}_{\boldsymbol{X}|\boldsymbol{Y}}\left(\frac{p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})}{p_{\boldsymbol{X}|\boldsymbol{Y}}^{*}(\boldsymbol{X}^{i}|\boldsymbol{Y}^{i})}\right)^{\alpha}\exp\left(-\frac{c_{9}\alpha}{8}n\eta_{f}^{2}\right) \leq \left(1 + \alpha\gamma_{f}\right)^{n}\exp\left(-\frac{c_{9}\alpha}{8}n\eta_{f}^{2}\right)$$

$$\leq \exp\left(-\frac{\alpha}{8}c_{9}n\eta_{f}^{2} + n\log(1 + \alpha\gamma_{f})\right) \leq \exp\left(-\frac{\alpha}{8}c_{9}n\eta_{f}^{2} + n\alpha\gamma_{f}\right), \tag{27}$$

where the second inequality follows from  $\lambda_f J(p_{X|Y}^*) \le c_9 \eta_f^2/8$  and the third inequality follows from Markov's inequality.

Our treatment of bounding  $I_6$  relies on a chaining argument over a suitable partition of  $\mathcal{F}_f$  and the left-truncation of likelihood ratios as in [38, 33]. Now, consider a partition of  $\mathcal{F}_f = \bigcup_{k=1}^{\infty} \bigcup_{j=0}^{\infty} \mathcal{F}_{kj}$ :

$$\mathcal{F}_{kj} = \left\{ p \in \mathcal{F}_{f} : 2^{i-1} \eta_{n}^{2} \leq d^{2}(p^{0}, p) \leq 2^{i} \eta_{n}^{2}, 2^{j-1} J^{0} \leq J(p) \leq 2^{j} J(p^{0}) \right\},$$

$$\mathcal{F}_{k0} = \left\{ p \in \mathcal{F}_{f} : 2^{i-1} \eta_{n}^{2} \leq d^{2}(p^{0}, p) \leq 2^{i} \eta_{n}^{2}, J(p) \leq J(p^{0}) \right\}; \quad k = 1, \dots, j = 0, \dots,$$

where  $\log(\mathbb{F}_i^{(\tau)}) = \log(p_{\mathbf{X}|\mathbf{Y}}^{(\tau)}(\mathbf{X}^i|\mathbf{Y}^i)) - \log(p_{\mathbf{X}|\mathbf{Y}}^0(\mathbf{X}^i|\mathbf{Y}^i))$ . Then for any  $\eta_f > 0$ ,

$$I_{6} \leq P\left(\sup_{d(p,p^{0})\geq\eta_{f}}\left(n^{-1}\sum_{i=1}^{n}\log\left(\mathbb{F}_{i}^{(\tau)}\right) - \lambda_{f}J(p_{\boldsymbol{X}|\boldsymbol{Y}}) + \lambda_{f}J(p_{\boldsymbol{X}|\boldsymbol{Y}}^{0})\right) \geq -c_{9}\eta_{f}^{2}/4\right)$$

$$\leq \sum_{k=1}^{\infty}\sum_{j=0}^{\infty}P\left(\sup_{p\in\mathcal{F}_{kj}}\left(n^{-1}\sum_{i=1}^{n}\log\left(\mathbb{F}_{i}^{(\tau)}\right) - \lambda_{f}J(p_{\boldsymbol{X}|\boldsymbol{Y}}) + \lambda_{f}J(p_{\boldsymbol{X}|\boldsymbol{Y}}^{0})\right) \geq -c_{9}\eta_{f}^{2}/4\right)$$

$$\equiv \sum_{k=1}^{\infty}\sum_{i=0}^{\infty}I_{kj},$$

$$(28)$$

where  $I_{kj} = P\left(\sup_{f \in \mathcal{F}_{kj}} \left(n^{-1} \sum_{i=1}^{n} \log \left(\mathbb{F}_{i}^{(\tau)}\right) - \lambda_{\mathrm{f}} J(p_{\mathbf{X}|\mathbf{Y}}) + \lambda_{\mathrm{f}} J(p_{\mathbf{X}|\mathbf{Y}}^{0})\right) \ge -c_{9} \eta_{\mathrm{f}}^{2}/4\right)$ . To treat  $I_{kj}$ , we control the mean and variance of  $\log \left(\mathbb{F}^{(\tau)}\right)$ . By Lemma 4 of [38],

$$-\sup_{p\in\mathcal{F}_{kj}} \mathbb{E}\left(\log(\mathbb{F}^{(\tau)})\right) = -\sup_{p\in\mathcal{F}_{kj}} \mathbb{E}_{\boldsymbol{Y}}\left(\mathbb{E}_{\boldsymbol{X}|\boldsymbol{Y}}\left(\log(\mathbb{F}^{(\tau)})\right)\right) \ge c_9 \inf_{p\in\mathcal{F}_{kj}} d^2(p, p^*) \ge c_9 (2^{k-1}\eta_f)^2,$$
(29)

and the variance is bounded by

$$\sup_{p \in \mathcal{F}_{kj}} \operatorname{Var}\left(\log(\mathbb{F}^{(\tau)})\right) \leq \sup_{p \in \mathcal{F}_{kj}} \mathbb{E}_{\boldsymbol{Y}}\left(\mathbb{E}_{\boldsymbol{X}|\boldsymbol{Y}}\left(\log(\mathbb{F}^{(\tau)})^{2}\right)\right)$$

$$\leq 4 \exp(\tau) \sup_{p \in \mathcal{F}_{kj}} \mathbb{E}_{\boldsymbol{Y}} h^{2}\left(p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}, p_{\boldsymbol{X}|\boldsymbol{Y}}\right) \leq 4 \exp(\tau)(2^{k}\eta_{f})^{2} \leq 8 \exp(\tau)\delta_{kj}/c_{9}, \tag{30}$$

where the second inequality follows from Lemma 3 of [38]. Then,  $I_{kj}$  is upper-bounded by

$$I_{kj} \leq P\left(\sup_{f \in \mathcal{F}_{kj}} \left(n^{-1} \sum_{i=1}^{n} \log\left(\mathbb{F}_{i}^{(\tau)}\right) - \mathbb{E}\log\left(\mathbb{F}^{(\tau)}\right)\right)$$

$$\geq -\sup_{f \in \mathcal{F}_{kj}} \left(\mathbb{E}\log\left(\mathbb{F}^{(\tau)}\right) + \lambda \left(J(p_{\mathbf{X}|\mathbf{Y}}^{0}) - J(p_{\mathbf{X}|\mathbf{Y}})\right)\right) - c_{9}\eta_{\mathrm{f}}^{2}/4\right)$$

$$\leq P\left(\sup_{f \in \mathcal{F}_{kj}} \left(n^{-1} \sum_{i=1}^{n} \log\left(\mathbb{F}_{i}^{(\tau)}\right) - \mathbb{E}\log\left(\mathbb{F}^{(\tau)}\right)\right) \geq \delta_{kj}\right) \leq 3 \exp\left(-a_{3}n\delta_{kj}\right), \tag{31}$$

where  $a_3 > 0$  is a constant,  $\delta_{kj} = c_9 2^{k-1} \eta_n^2 / 2 + \lambda (2^{j-1} - 1) J(p_{X|Y}^0)$ ,  $\delta_{k0} = c_9 2^{k-2} \eta_n^2 / 2$ , the second inequality follows from the assumption that  $\lambda J(p_{X|Y}^0) \leq c_9 \eta_{\rm f}^2 / 4$  and (29), and the last inequality follows from Lemma 2 and the fact that the *j*-th  $(j \geq 2)$  moment  $\mathbb{E}(|\log(\frac{p_{X|Y}^{(\tau)}(X|Y)}{p_{X|Y}^0(X|Y)})|^j)$  is bounded by

$$\mathbb{E}_{\boldsymbol{Y}} \mathbb{E}_{\boldsymbol{X}|\boldsymbol{Y}} \Big( \exp \Big( \Big| \log \Big( \frac{p_{\boldsymbol{X}|\boldsymbol{Y}}^{(\tau)}(\boldsymbol{X}|\boldsymbol{Y})}{p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}(\boldsymbol{X}|\boldsymbol{Y})} \Big) \Big| \Big) - 1 - \frac{1}{2} \Big| \log \Big( \frac{p_{\boldsymbol{X}|\boldsymbol{Y}}^{(\tau)}(\boldsymbol{X}|\boldsymbol{Y})}{p_{\boldsymbol{X}|\boldsymbol{Y}}^{0}(\boldsymbol{X}|\boldsymbol{Y})} \Big) \Big| \Big)$$

$$\leq j! 2^{j} a_{1} \mathbb{E}_{\boldsymbol{Y}} \| (p_{\boldsymbol{X}|\boldsymbol{Y}})^{1/2} - (p_{\boldsymbol{X}|\boldsymbol{Y}}^{0})^{1/2} \|_{2}^{2},$$

where  $a_1 = \left(\exp(\tau/2) - 1 - \tau/2\right)/\left(1 - \exp(-\tau/2)\right)^2 > 0$  is a constant and the last inequality follows from Lemma 5 in [34]. It suffices to verify the condition (2.4) of [38]. A combination of (28) and (31) yield that  $I_6 \leq \sum_{k=1}^{\infty} \sum_{j=0}^{\infty} 3 \exp(-c_{13}n\delta_{kj}^2) \leq 7 \exp(-c_{13}n\eta_{\rm f}^2)$ , which, together with (27) yields that  $P\left(d(\hat{p}_{X|Y}^{\rm f}, p_{X|Y}^0) \geq \eta_{\rm f}\right) \leq I_5 + I_6 \leq 7 \exp(-c_{13}n\eta_{\rm f}^2) + \exp\left(-\frac{\alpha}{8}c_9n\eta_{\rm f}^2 + n\alpha\gamma_{\rm f}\right)$ . The desired result then follows.

**Proof of Theorem 3.** Let  $(\tilde{X}^i, \tilde{Y}^i)_{i=1}^N$  be a cross-validation sample. By (5),

$$-\frac{1}{N}\sum_{i=1}^{N}\log\widehat{p}_{\boldsymbol{x}|\boldsymbol{Y}}^{c}(\tilde{\boldsymbol{X}}^{i}|\tilde{\boldsymbol{Y}}^{i})\leq\min\Big(-\frac{1}{N}\sum_{i=1}^{N}\log\widehat{p}_{\boldsymbol{x}|\boldsymbol{Y}}^{f}(\tilde{\boldsymbol{X}}^{i}|\tilde{\boldsymbol{Y}}^{i}),-\frac{1}{N}\sum_{i=1}^{N}\log\widehat{p}_{\boldsymbol{x}|\boldsymbol{Y}}^{b}(\tilde{\boldsymbol{X}}^{i}|\tilde{\boldsymbol{Y}}^{i})\Big),$$

then the desired result follows from the law of large number, by taking the limit for both the sides as  $N \to \infty$ . This completes the proof.

**Proof of Corollary** 1. For the direct sequential generation, we apply the same argument in the proof of Theorem 2. Denote

$$I_{7} = P\left(\sup_{d(p,p^{0}) \geq \eta_{f}} \sum_{i=1}^{n} \frac{1}{nT} \sum_{t=1}^{T} \log\left(\frac{p^{0}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i},\boldsymbol{Y}^{i})}{p^{*}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i},\boldsymbol{Y}^{i})}\right) - \lambda_{f} J_{f}(\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) + \lambda_{f} J_{f}(\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{*}) \geq \frac{c_{9}\eta_{f}^{2}}{4}\right),$$

$$I_{8} = P\left(\sup_{d(p,p^{0}) \geq \eta_{f}} \sum_{i=1}^{n} \frac{1}{nT} \sum_{t=1}^{T} \log\left(\frac{p^{(\tau)}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i},\boldsymbol{Y}^{i})}{p^{0}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i},\boldsymbol{Y}^{i})}\right) - \lambda_{f} J_{f}(\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}) + \lambda_{f} J_{f}(\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) \geq -\frac{c_{9}\eta_{f}^{2}}{4}\right).$$

Then  $P(d(\widehat{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}}^{f}, \boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) \geq \eta_{f}) \leq I_{7} + I_{8}$ , where  $p^{(\tau)}(X_{t+1}|\boldsymbol{X}_{1:t}, \boldsymbol{Y})$  is the left truncation of  $p(X_{t+1}|\boldsymbol{X}_{1:t}, \boldsymbol{Y})$  as defined in the proof of Theorem 2.

For  $I_7$ ,

$$I_{7} \leq P\left(\prod_{i=1}^{n} \left(\prod_{t=1}^{T} \frac{p^{0}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i}, \boldsymbol{Y}^{i})}{p^{*}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i}, \boldsymbol{Y}^{i})}\right)^{\frac{\alpha}{T-1}} \geq \exp\left(\frac{c_{9}\alpha}{8}n\eta_{f}^{2}\right)\right)$$

$$\leq \prod_{i=1}^{n} \mathbb{E}\left(\prod_{t=1}^{T} \frac{p^{0}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i}, \boldsymbol{Y}^{i})}{p^{*}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i}, \boldsymbol{Y}^{i})}\right)^{\frac{\alpha}{T-1}} \exp\left(-\frac{c_{9}\alpha}{8}n\eta_{f}^{2}\right)$$

$$\leq \prod_{i=1}^{n} \mathbb{E}\left(\left(\frac{p^{0}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i}, \boldsymbol{Y}^{i})}{p^{*}(X_{t+1}^{i}|\boldsymbol{X}_{1:t}^{i}, \boldsymbol{Y}^{i})}\right)^{\alpha}\right) \exp\left(-\frac{c_{9}\alpha}{8}n\eta_{f}^{2}\right)$$

$$\leq (1 + \alpha\gamma_{f})^{n} \exp\left(-\frac{c_{9}\alpha}{8}n\eta_{f}^{2}\right) \leq \bar{r} \exp\left(-\frac{\alpha}{8}c_{9}n\eta_{f}^{2} + n\alpha\gamma_{f}\right).$$

For  $I_8$ , let  $\mathbb{F}_t^{(\tau)} = \log \left( p^{(\tau)}(X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y})/p^0(X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y}) \right)$ . For the first moment,

$$\mathbb{E}\left(T^{-1}\sum_{t=1}^{T}\mathbb{F}_{t}^{(\tau)}\right) = T^{-1}\sum_{t=1}^{T}\mathbb{E}_{\boldsymbol{X}_{1:t},\boldsymbol{Y}}\left(\mathbb{E}_{X_{t+1}}\left(\mathbb{F}_{t}^{(\tau)}|\boldsymbol{X}_{1:t},\boldsymbol{Y}\right)\right) \\
\leq -c_{9}T^{-1}\sum_{t=1}^{T}\mathbb{E}_{\boldsymbol{X}_{1:t},\boldsymbol{Y}}\left(\left\|\left(p_{X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y}}^{(\tau)}\right)^{\frac{1}{2}}-\left(p_{X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y}}^{0}\right)^{\frac{1}{2}}\right\|_{2}^{2}\right) = -c_{9}d^{2}(\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{(\tau)},\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{0}).$$

For the j-th moment with  $j \geq 2$ ,

$$\mathbb{E} \left| T^{-1} \sum_{t=1}^{T} \mathbb{F}_{t}^{(\tau)} \right|^{j} \leq T^{-1} \sum_{t=1}^{T} \mathbb{E}_{\boldsymbol{X}_{1:t},\boldsymbol{Y}} \left( \mathbb{E}_{X_{t+1}} \left( |\mathbb{F}_{t}^{(\tau)}|^{j} |\boldsymbol{X}_{1:t}, \boldsymbol{Y} \right) \right) \\
\leq j! 2^{j} T^{-1} \sum_{t=1}^{T} \mathbb{E}_{\boldsymbol{X}_{1:t},\boldsymbol{Y}} \left( \mathbb{E}_{X_{t+1}} \left( (\exp(|\mathbb{F}_{t}^{(\tau)}|/2) - 1 - |\mathbb{F}_{t}^{(\tau)}|/2) |\boldsymbol{X}_{1:t}, \boldsymbol{Y} \right) \right) \\
\leq j! 2^{j} a_{3} \left( \sum_{t=1}^{T} \mathbb{E}_{\boldsymbol{X}_{1:t},\boldsymbol{Y}} \left( \| (p_{X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y}}^{(\tau)})^{\frac{1}{2}} - (p_{X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y}}^{0})^{\frac{1}{2}} \|_{2}^{2} \right) \right) \\
\leq j! 2^{j} a_{1} d^{2} (\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{(\tau)}, \boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{0}),$$

where the first inequality follows from the Jensen's inequality. Then

$$P(d(\widehat{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}}^{f}, \boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) \ge \eta_{f}) \le 6 \exp(-c_{7}n\eta_{f}^{2}) + \exp(-\frac{\alpha}{8}c_{9}n\eta_{f}^{2} + n\alpha\gamma_{f}), \tag{32}$$

follows the same arguments as in the proof of Theorem 2

For the indirect generation, let  $p_t(\cdot) = p(\cdot|\boldsymbol{X}_{1:t})$  and  $\mathbb{E}_t(\cdot) = \mathbb{E}(\cdot|\boldsymbol{X}_{1:t})$ . Then,

$$d(\widehat{p}_{X|Y}^{b}, p_{X|Y}^{0}) = \left(T^{-1} \sum_{t=1}^{T} \mathbb{E}_{X_{1:t}} \mathbb{E}_{Y|X_{1:t}} h^{2}(\widehat{p}_{X_{t+1}|X_{1:t},Y}^{b}, p_{X_{t+1}|X_{1:t},Y}^{0})\right)^{\frac{1}{2}}$$

$$\leq 2 \left(T^{-1} \sum_{t=1}^{T} \left(\mathbb{E}\left(h(\widehat{p}_{X_{t+1}|X_{1:t}}, p_{X_{t+1}|X_{1:t}}^{0}) + \left(\mathbb{E}_{t} h^{2}(\widehat{p}_{Y|X_{1:t+1}}, p_{Y|X_{1:t+1}}^{0})\right)^{\frac{1}{2}}\right)^{2}\right)\right)^{\frac{1}{2}}$$

$$\leq 2 \left(2T^{-1} \sum_{t=1}^{T} \left(\mathbb{E} h^{2}(\widehat{p}_{X_{t+1}|X_{1:t}}, p_{X_{t+1}|X_{1:t}}^{0}) + \mathbb{E} h^{2}(\widehat{p}_{Y|X_{1:t+1}}, p_{Y|X_{1:t+1}}^{0})\right)\right)^{\frac{1}{2}}$$

$$\leq 2\sqrt{2} \left(d(\widehat{p}_{X}, p_{X}^{0}) + d(\widehat{p}_{Y|X}, p_{Y|X}^{0})\right),$$

where the first inequality follows from (26) by replacing  $p(\cdot)$  as  $p_t(\cdot)$ , and  $d^2(\widehat{p}_X, p_X^0) =$ 

$$\bar{\mathbb{E}}h^2(\widehat{p}_{X_{t+1}|X_{1:t}}, p_{X_{t+1}|X_{1:t}}^0), d^2(\widehat{p}_{Y|X}, p_{Y|X}^0) = \bar{\mathbb{E}}h^2(\widehat{p}_{Y|X_{1:t+1}}, p_{Y|X_{1:t+1}}^0). \text{ Therefore,}$$

$$P(d(\widehat{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}}^{b}, \boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}^{0}) \geq 2\sqrt{2}(\eta_{b} + \eta_{m})) \leq P(d(\widehat{\boldsymbol{p}}_{\boldsymbol{X}}, \boldsymbol{p}_{\boldsymbol{X}}^{0}) \geq \eta_{m}) + P(d(\widehat{\boldsymbol{p}}_{\boldsymbol{Y}|\boldsymbol{X}}, \boldsymbol{p}_{\boldsymbol{Y}|\boldsymbol{X}}^{0}) \geq \eta_{b})$$

$$\leq 7 \exp(-c_{7}(n+\tilde{n})\eta_{m}^{2}) + \exp(-\frac{\alpha c_{9}}{8}(n+\tilde{n})\eta_{m}^{2} + \alpha(n+\tilde{n})\gamma_{m})$$

$$+ 7 \exp(-c_{8}n\eta_{b}^{2}) + \exp(-\frac{\alpha c_{9}}{8}n\eta_{b}^{2} + \alpha n\gamma_{b}),$$

where the last inequality follows from (32). Similarly, bounds for  $P(d(\hat{p}_X, p_X^0) \ge \eta_m)$  and  $d(\hat{p}_{Y|X}, p_{Y|X}^0)$  can be established. The desired result then follows.

**Proof of Corollary 2.** It suffices to verify the entropy conditions in Corollary 1. For the direct generation, let  $\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}} = \{p(X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y};\boldsymbol{\theta}_{\mathrm{f}})\}_{t=1}^{T}$  and  $\bar{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}} = \{p(X_{t+1}|\boldsymbol{X}_{1:t},\boldsymbol{Y};\bar{\boldsymbol{\theta}}_{\mathrm{f}})\}_{t=1}^{T}$  in  $\mathcal{F}_{\mathrm{f},k}$ . Then

$$\kappa^{2}(\boldsymbol{p}_{\boldsymbol{X}|\boldsymbol{Y}}, \bar{\boldsymbol{p}}_{\boldsymbol{X}|\boldsymbol{Y}}) \leq \mathbb{E} \|\boldsymbol{\sigma}^{\frac{1}{2}}(\boldsymbol{W}_{f}^{o}\boldsymbol{h}_{t-1}) - \boldsymbol{\sigma}^{\frac{1}{2}}(\bar{\boldsymbol{W}}_{f}^{o}\bar{\boldsymbol{h}}_{t-1})\|_{2}^{2} \\
\leq \frac{1}{2}\mathbb{E} \|\boldsymbol{W}_{f}^{o}\boldsymbol{h}_{t-1} - \bar{\boldsymbol{W}}_{f}^{o}\bar{\boldsymbol{h}}_{t-1}\|_{2}^{2} \\
\leq \mathbb{E} \Big( \|(\boldsymbol{W}_{f}^{o} - \bar{\boldsymbol{W}}_{f}^{o})\boldsymbol{h}_{t-1}\|_{2}^{2} + \|\bar{\boldsymbol{W}}_{f}^{o}(\boldsymbol{h}_{t-1} - \bar{\boldsymbol{h}}_{t-1})\|_{2}^{2} \Big) \\
\leq \mathbb{E} \Big( \|\boldsymbol{W}_{f}^{o} - \bar{\boldsymbol{W}}_{f}^{o}\|_{F}^{2} \|\boldsymbol{h}_{t-1}\|_{2}^{2} + \|\bar{\boldsymbol{W}}_{f}^{o}\|_{F}^{2} \|\boldsymbol{h}_{t-1} - \bar{\boldsymbol{h}}_{t-1}\|_{2}^{2} \Big) \\
\leq \Big( \frac{2k(4k)^{T+1}}{T(4k-1)^{2}} \Big) \Big( \|\boldsymbol{W}_{f}^{x} - \bar{\boldsymbol{W}}_{f}^{x}\|_{F}^{2} + 2r_{f} \|\boldsymbol{W}_{f}^{h} - \bar{\boldsymbol{W}}_{f}^{h}\|_{F}^{2} + 4kc_{\boldsymbol{Y}} \|\boldsymbol{W}_{f}^{y} - \bar{\boldsymbol{W}}_{f}^{y}\|_{F}^{2} \Big) \\
+ r_{f} \|\boldsymbol{W}_{f}^{o} - \bar{\boldsymbol{W}}_{f}^{o}\|_{F}^{2} \leq T^{-1} \max(2r_{f}, 4kc_{15})^{2} (4k)^{T} \|\boldsymbol{\theta}_{f} - \bar{\boldsymbol{\theta}}_{f}\|_{2}^{2},$$

where the last inequality uses the fact that  $\|\boldsymbol{h}_t - \bar{\boldsymbol{h}}_t\|_2^2 \le \frac{(4k)^t - 1}{4k - 1} (2\|\boldsymbol{W}_{\mathrm{f}}^x - \bar{\boldsymbol{W}}_{\mathrm{f}}^x\|_F^2 + 4r_{\mathrm{f}}\|\boldsymbol{W}_{\mathrm{f}}^h - \bar{\boldsymbol{W}}_{\mathrm{f}}^h\|_F^2) + (4k)^t \|\boldsymbol{W}_{\mathrm{f}}^y - \bar{\boldsymbol{W}}_{\mathrm{f}}^y\|_F^2 \|\boldsymbol{Y}\|_2^2$ , which uses the fact that

$$\|\boldsymbol{h}_{t} - \bar{\boldsymbol{h}}_{t}\|_{2}^{2} \leq \|(\boldsymbol{W}_{f}^{x} - \bar{\boldsymbol{W}}_{f}^{x})\boldsymbol{1}_{[X_{t}]} + \boldsymbol{W}_{f}^{h}\boldsymbol{h}_{t-1} - \bar{\boldsymbol{W}}_{f}^{h}\bar{\boldsymbol{h}}_{t-1}\|_{2}^{2}$$

$$\leq 2\|\boldsymbol{W}_{f}^{x} - \bar{\boldsymbol{W}}_{f}^{x}\|_{F}^{2} + 4r_{f}\|\boldsymbol{W}_{f}^{h} - \bar{\boldsymbol{W}}_{f}^{h}\|_{F}^{2} + 4k\|\boldsymbol{h}_{t-1} - \bar{\boldsymbol{h}}_{t-1}\|_{2}^{2},$$

and  $\|\boldsymbol{h}_0 - \bar{\boldsymbol{h}}_0\|_2^2 \le \|\boldsymbol{W}_f^y - \bar{\boldsymbol{W}}_f^y\|_F^2 \|\boldsymbol{Y}\|_2^2$ .

Hence,  $H(u, \mathcal{F}_{f,k}) \leq \Lambda_f \log \left(\frac{3 \max(2r_f, 4kc_{15})(4k)^{(T+1)/2}/T^{1/2}}{u}\right)$  and the entropy condition is met by setting  $\epsilon_f = \left(\frac{\Lambda_f}{n} \log \left(\frac{\max(r_f, 2c_{15})2^T/T^{1/2}n}{\Lambda_f}\right)\right)^{1/2}$ .

For the indirect generation, it suffices to verify the entropy conditions in Corollary 2. Let  $\mathbf{p}_{\mathbf{X}} = \{P_{\mathbf{X}}(X_{t+1}|\mathbf{X}_{1:t};\boldsymbol{\theta}_m)\}_{t=1}^{T-1}$  and  $\bar{\mathbf{p}}_{\mathbf{X}} = \{P_{\mathbf{X}}(X_{t+1}|\mathbf{X}_{1:t};\bar{\boldsymbol{\theta}}_m)\}_{t=1}^{T-1}$ . Note that  $\mathbf{h}_0 = \bar{\mathbf{h}}_0 = \mathbf{0}_{r_m}$  and

$$\kappa^{2}(\boldsymbol{p}_{\boldsymbol{X}}, \bar{\boldsymbol{p}}_{\boldsymbol{X}}) \leq \bar{\mathbb{E}}\left(\|\boldsymbol{W}_{m}^{o} - \bar{\boldsymbol{W}}_{m}^{o}\|_{F}^{2}\|\boldsymbol{h}_{t-1}\|_{2}^{2} + \|\bar{\boldsymbol{W}}_{m}^{o}\|_{F}^{2}\|\boldsymbol{h}_{t-1} - \bar{\boldsymbol{h}}_{t-1}\|_{2}^{2}\right) 
\leq \frac{2k(4k)^{T+1}}{T(4k-1)^{2}}\left(\|\boldsymbol{W}_{m}^{x} - \bar{\boldsymbol{W}}_{m}^{x}\|_{F}^{2} + 2r_{f}\|\boldsymbol{W}_{m}^{h} - \bar{\boldsymbol{W}}_{m}^{h}\|_{F}^{2}\right) + r_{f}\|\boldsymbol{W}_{m}^{o} - \bar{\boldsymbol{W}}_{m}^{o}\|_{F}^{2} 
\leq 2r_{m}T^{-1}(4k)^{T}\|\boldsymbol{\theta}_{m} - \bar{\boldsymbol{\theta}}_{m}\|_{2}^{2}.$$

Then,  $H(u, \mathcal{F}_{m,k}) \leq \Lambda_m \log \left(\frac{3r_m(4k)^{(T+1)/2}T^{-1/2}}{u}\right)$  and the entropy condition is met by setting  $\epsilon_m = O_p\left(\left(\frac{\Lambda_m}{n+\tilde{n}}\log\left(\frac{r_m(n+\tilde{n})2^TT^{-1/2}}{\Lambda_m}\right)\right)^{1/2}\right)$ .

Moreover, if  $\mathbf{Y} \in \{0, 1\}^K$ ,

$$\kappa^{2}(\boldsymbol{p}_{\boldsymbol{Y}|\boldsymbol{X}}, \bar{\boldsymbol{p}}_{\boldsymbol{Y}|\boldsymbol{X}}) \leq \bar{\mathbb{E}} \| (\boldsymbol{\sigma}(\boldsymbol{\theta}_{b}\boldsymbol{\mathcal{E}}(\boldsymbol{X}_{1:t})))^{\frac{1}{2}} - (\boldsymbol{\sigma}(\bar{\boldsymbol{\theta}}_{b}\boldsymbol{\mathcal{E}}(\boldsymbol{X}_{1:t})))^{\frac{1}{2}} \|_{2}^{2}, 
\leq \bar{\mathbb{E}} \| \boldsymbol{\sigma}^{\frac{1}{2}} (\boldsymbol{\theta}_{b}\boldsymbol{\mathcal{E}}(\boldsymbol{X}_{1:t})) - \boldsymbol{\sigma}^{\frac{1}{2}} (\bar{\boldsymbol{\theta}}_{b}\boldsymbol{\mathcal{E}}(\boldsymbol{X}_{1:t})) \|_{2}^{2} \leq \frac{1}{2} \bar{\mathbb{E}} \| (\boldsymbol{\theta}_{b} - \bar{\boldsymbol{\theta}}_{b}) \boldsymbol{\mathcal{E}}(\boldsymbol{X}_{1:t}) \|_{2}^{2} 
\leq \frac{1}{2} \| \boldsymbol{\theta}_{b} - \bar{\boldsymbol{\theta}}_{b} \|_{F}^{2} \bar{\mathbb{E}} \| \boldsymbol{\mathcal{E}}(\boldsymbol{X}_{1:t}) \|_{2}^{2}.$$

Similarly,  $H(u, \mathcal{F}_{b,k}) \leq \Lambda_b \log \left(\frac{3\sqrt{kc_{16}}}{u\sqrt{2}}\right)$  and the entropy condition is met by setting  $\epsilon_b = O_p\left(\left(\frac{\Lambda_b}{n}\log(\frac{\sqrt{c_{16}}n}{\Lambda_b})\right)^{1/2}\right)$ .

If  $\boldsymbol{y} \in \mathbb{R}^K$ , then

$$\kappa^{2}(\boldsymbol{p}_{\boldsymbol{Y}|\boldsymbol{X}}, \bar{\boldsymbol{p}}_{\boldsymbol{Y}|\boldsymbol{X}}) = \bar{\mathbb{E}}\left(1 - \exp\left(-\frac{1}{8}\|(\boldsymbol{\theta}_{b} - \bar{\boldsymbol{\theta}}_{b})\mathcal{E}(\boldsymbol{X}_{1:t})\|_{2}^{2}\right)\right),$$

$$\leq \frac{1}{8}\bar{\mathbb{E}}\|(\boldsymbol{\theta}_{b} - \bar{\boldsymbol{\theta}}_{b})\mathcal{E}(\boldsymbol{X}_{1:t})\|_{2}^{2} \leq \frac{1}{8}\|\boldsymbol{\theta}_{b} - \bar{\boldsymbol{\theta}}_{b}\|_{F}^{2}\bar{\mathbb{E}}\|\mathcal{E}(\boldsymbol{X}_{1:t})\|_{2}^{2},$$
(33)

implying that  $H(u, \mathcal{F}_{b,k}) \leq \Lambda_b \log \left(\frac{3\sqrt{kc_{16}}}{2\sqrt{2}u}\right)$ , and the entropy condition holds when  $\epsilon_b = O_p\left(\left(\frac{\Lambda_b}{n}\log(\frac{\sqrt{c_{16}n}}{\Lambda_b})\right)^{1/2}\right)$ . This completes the proof.

**Lemma 2.** Let  $v_n(f) = n^{-1/2} \sum_{i=1}^n \left( v\left(f(\mathbf{X}^i), f^0(\mathbf{X}^i)\right) - \mathbb{E}v\left(f(\mathbf{X}^i), f^0(\mathbf{X}^i)\right) \right)$ , assume that there exist some generic constants  $a_2 > 0$  and  $a_3 > 0$ , for  $j \ge 2$ , such that

$$\mathbb{E}|v(f(X), f^0(X))|^j \le a_2 j! 2^j d^2(f, f^0),$$

and for any  $\delta > 0$ , if

$$\int_{\delta/2^8}^{2^{1/2}\delta^{1/2}} H^{1/2}(u, \mathcal{V}_{\delta}) du \le a_3 n^{1/2} \delta,$$

where  $V_{\delta} = \{v(f, f^0) : d^2(f, f^0) \leq \delta, f \in \mathcal{F}\}$ , then there exist some constants  $a_4 > 0$  and  $a_5 > 0$  depending on  $a_2$  and  $a_3$  such that

$$P^* \left( \sup_{d^2(f,f^0) < \delta; f \in \mathcal{F}} v_n(f) \ge a_4 n^{1/2} \delta \right) \le 3 \exp(-a_5 n \delta), \tag{34}$$

where  $P^*$  is the outer probability measure corresponding to  $p_X^0$ .

**Proof of Lemma 2.** The result follows from Lemmas 5 and Lemma 7 in [38], by replacing the Hellinger distance in Lemma 5 as a generic distance  $d(\cdot, \cdot)$ .

# References

- [1] C. M. Bishop. Pattern recognition and machine learning. springer, 2006.
- [2] L. Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings* of COMPSTAT'2010, pages 177–186. Springer, 2010.
- [3] L. Breiman. Random forests. Machine Learning, 45(1):5–32, 2001.

- [4] R. Caruana, S. Lawrence, and C. L. Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In Advances in Neural Information Processing Systems, pages 402–408, 2001.
- [5] J. Cheng and M. Lapata. Neural summarization by extracting sentences and words. arXiv preprint arXiv:1603.07252, 2016.
- [6] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078, 2014.
- [7] F. Cozman and I. Cohen. Risks of semi-supervised learning. Semi-supervised learning, pages 56–72, 2006.
- [8] F. G. Cozman, I. Cohen, and M. Cirelo. Unlabeled data can degrade classification performance of generative classifiers. In *Flairs conference*, pages 327–331, 2002.
- [9] F. G. Cozman, I. Cohen, and M. C. Cirelo. Semi-supervised learning of mixture models. In Proceedings of the 20th international conference on machine learning (ICML-03), pages 99–106, 2003.
- [10] B. Dai, X. Shen, and J. Wang. Embedding learning. *Journal of the American Statistical Association*, (in press), 2020.
- [11] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2018.
- [12] L. Dong, N. Yang, W. Wang, F. Wei, X. Liu, Y. Wang, J. Gao, M. Zhou, and H.-W. Hon. Unified language model pre-training for natural language understanding and generation. In *Advances in Neural Information Processing Systems*, pages 13042–13054, 2019.

- [13] V. Gjorgjioski, D. Kocev, and S. Džeroski. Comparison of distances for multi-label classification with pcts. In *Proceedings of the Slovenian KDD Conference on Data Mining and Data Warehouses (SiKDD'11)*, 2011.
- [14] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Infor*mation Processing Systems, pages 2672–2680, 2014.
- [15] D. He, Y. Xia, T. Qin, L. Wang, N. Yu, T.-Y. Liu, and W.-Y. Ma. Dual learning for machine translation. In Advances in Neural Information Processing Systems, pages 820–828, 2016.
- [16] S. Hochreiter and J. Schmidhuber. Long short-term memory. Neural Computation, 9(8):1735–1780, 1997.
- [17] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern* Recognition, pages 3128–3137, 2015.
- [18] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In Advances in Neural Information Processing Systems, pages 3581–3589, 2014.
- [19] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, pages 4743–4751, 2016.
- [20] I. Langkilde. Forest-based statistical sentence generation. In Proceedings of the 1st North American chapter of the Association for Computational Linguistics conference, pages 170–177. Association for Computational Linguistics, 2000.

- [21] J. D. Lee, M. Simchowitz, M. I. Jordan, and B. Recht. Gradient descent only converges to minimizers. In *Conference on Learning Theory*, pages 1246–1257, 2016.
- [22] S. Merity, N. S. Keskar, and R. Socher. Regularizing and optimizing LSTM language models. In *International Conference on Learning Representations*, 2018.
- [23] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119, 2013.
- [24] T. Mikolov, W.-T. Yih, and G. Zweig. Linguistic regularities in continuous space word representations. In Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 746–751, 2013.
- [25] V. Mullachery and V. Motwani. Image captioning. arXiv preprint arXiv:1805.09137, 2018.
- [26] R. Nallapati, I. Melnyk, A. Kumar, and B. Zhou. Sengen: Sentence generating neural variational topic model. arXiv preprint arXiv:1708.00308, 2017.
- [27] Y. Ollivier, C. Tallec, and G. Charpiat. Training recurrent networks online without backtracking. arXiv preprint arXiv:1507.07680, 2015.
- [28] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, pages 311–318. Association for Computational Linguistics, 2002.
- [29] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *International Conference on Machine Learning*, pages 1310–1318, 2013.

- [30] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9, 2019.
- [31] M. Schmidt, N. Le Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112, 2017.
- [32] M. Schuster and K. K. Paliwal. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11):2673–2681, 1997.
- [33] X. Shen. On the method of penalization. Statistica Sinica, 8(2):337–357, 1998.
- [34] X. Shen and W. H. Wong. Convergence rate of sieve estimates. *Annals of Statistics*, pages 580–615, 1994.
- [35] S. Smale and D.-X. Zhou. Estimating the approximation error in learning theory. *Analysis and Applications*, 1(01):17–41, 2003.
- [36] I. Sutskever, J. Martens, and G. E. Hinton. Generating text with recurrent neural networks. In *Proceedings of the 28th International Conference on Machine Learning* (ICML-11), pages 1017–1024, 2011.
- [37] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3156–3164, 2015.
- [38] W. H. Wong, X. Shen, et al. Probability inequalities for likelihood ratios and convergence rates of sieve mles. *Annals of Statistics*, 23(2):339–362, 1995.
- [39] T. Yang and C. E. Priebe. The effect of model misspecification on semi-supervised classification. *IEEE transactions on pattern analysis and machine intelligence*, 33(10):2093–2103, 2011.

- [40] D. Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:103–114, 2017.
- [41] H.-F. Yu, F.-L. Huang, and C.-J. Lin. Dual coordinate descent methods for logistic regression and maximum entropy models. *Machine Learning*, 85(1-2):41–75, 2011.
- [42] D. Zhou, T. Hofmann, and B. Schölkopf. Semi-supervised learning on directed graphs. In Advances in Neural Information Processing Systems, pages 1633–1640, 2004.
- [43] Y. Zhu, S. Lu, L. Zheng, J. Guo, W. Zhang, J. Wang, and Y. Yu. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 1097–1100, 2018.

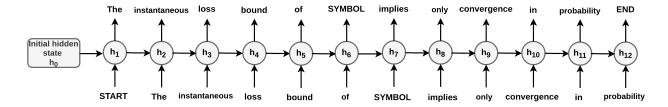


Figure 1: A generated sentence by indirect and direct RNN generators in (20) and (15), where the RNN architecture is displayed, in which sentence "The instantaneous loss bound of SYMBOL implies only convergence in probability" with topic "MISC" is consecutively generated by words,  $h_t$  is the hidden node of RNNs in (20) and (15), and  $h_0$  is the initial hidden state, which is zero under (15) and "MISC" under (20).

Table 1: Eleven next-word-prediction sequences with associated with a sentence.

Topic	Sentence
MISC	The loss bound of SYMBOL implies convergence in probability.
1.	Null Null Null Null Null Null Null Null
2.	Null Null Null Null Null Null Null START The $ ightarrow$ loss
3.	Null Null Null Null Null Null START The loss $\rightarrow$ bound
4.	Null Null Null Null Null START The loss bound $\rightarrow$ of
5.	Null Null Null Null Null START The loss bound of $\rightarrow$ SYMBOL
6.	Null Null Null Null START The loss bound of SYMBOL $\rightarrow$ implies
7.	Null Null Null START The loss bound of SYMBOL implies $\rightarrow$ convergence
8.	Null Null START The loss bound of SYMBOL implies convergence $\rightarrow$ in
9.	Null Null START The loss bound of SYMBOL implies convergence in $\rightarrow$ probability
10.	Null START The loss bound of SYMBOL implies only convergence in probability $\rightarrow$ .
11.	START The loss bound of SYMBOL implies only convergence in probability . $\rightarrow$ END

Table 2: Test errors in loss functions—Entropy,  $BLEU_l$ , and SF- $BLEU_l$  (standard errors in parentheses) of various generators based on 20 random partitions of the UCI sentence categorization text corpus. Here "Separate RNN", "Indirect", "Direct", "Direct-GPT2" and "Coupled" denote the separate RNN, indirect, and direct generators based on the RNN-LSTM architecture, the direct generator based on the RNN-GPT architecture, and the coupled generator, while Indirect-label or Coupled-label refers to the generation without unlabeled data.

Method	Entropy	$\mathrm{BLEU}_{1} ext{-loss}$	$\mathrm{BLEU}_2 ext{-loss}$	$\mathrm{BLEU_{3} ext{-}loss}$	$\mathrm{BLEU_{4} ext{-}loss}$	
Y: categorical la	bel					
Separate RNN	9.317(.040)	0.895(.010)	0.926(.008)	0.954(.007)	0.971(.005)	
Indirect	7.424(.049)	0.768(.003)	0.854(.002)	0.885(.002)	0.914(.002)	
Indirect-label	8.839(.060)	0.831(.008)	0.878(.005)	0.899(.004)	0.923(.003)	
Direct	9.537(.054)	0.823(.008)	0.872(.005)	0.895(.005)	0.919(.004)	
Direct-GPT2	8.684(.051)	0.900(.006)	0.954(.002)	0.970(.001)	0.981(.001)	
Coupled	7.424(.049)	0.768(.003)	0.854(.002)	0.885(.002)	0.914(.002)	
Coupled-label	8.644(.050)	0.880(.008)	0.932(.008)	0.949(.007)	0.963(.006)	
		$SF$ - $BLEU_1$ - $loss$	${\rm SF\text{-}BLEU_2\text{-}loss}$	$SF$ -BLEU $_3$ -loss	$SF$ - $BLEU_4$ - $loss$	
Separate RNN		0.076(.010)	0.208(.027)	0.271(.036)	0.303(.043)	
Indirect		0.105(.006)	0.296(.009)	0.416(.012)	0.502(.013)	
Indirect-label		0.138(.008)	0.363(.022)	0.472(.029)	0.545(.036)	
Direct		0.139(.006)	0.372(.019)	0.487(.026)	0.561(.032)	
Direct-GPT2		0.053(.006)	0.159(.019)	0.255(.031)	0.320(.040)	
Coupled		0.105(.006)	0.296(.009)	0.416(.012)	0.502(.013)	
Coupled-label		0.082(.011)	0.233(.028)	0.342(.038)	0.417(.045)	
Method	Entropy	$\mathrm{BLEU}_{1} ext{-loss}$	$\mathrm{BLEU}_{2} ext{-loss}$	$\mathrm{BLEU}_3$ -loss	$\mathrm{BLEU_{4} ext{-}loss}$	
Y: continuous label based on Doc2Vec 23 24						
Indirect	7.641(.036)	0.768(.005)	0.851(.003)	0.883(.003)	0.912(.003)	
Indirect-label	8.512(.041)	0.912(.010)	0.937(.008)	0.949(.007)	0.960(.005)	
Direct	9.102(.050)	0.916(.010)	0.939(.007)	0.950(.005)	0.961(.004)	
Coupled	7.641(.036)	0.768(.005)	0.851(.003)	0.883(.003)	0.912(.003)	
Coupled-label	8.512(.041)	0.912(.010)	0.937(.008)	0.949(.007)	0.960(.005)	
		$SF$ - $BLEU_1$ - $loss$	$SF$ -BLEU $_2$ -loss	$SF$ -BLEU $_3$ -loss	$SF-BLEU_4-loss$	
Indirect		0.097(.005)	0.261(.008)	0.361(.010)	0.440(.012)	
Indirect-labeled		0.064(.010)	0.165(.026)	0.211(.035)	0.232(.040)	
Direct		0.079(.014)	0.202(.037)	0.252(.046)	0.271(.050)	
Coupled		0.097(.005)	0.261(.008)	0.361(.010)	0.440(.012)	
Coupled-label		0.064(.010)	0.165(.026)	0.211(.035)	0.232(.040)	

Table 3: An abstract generated by the coupled generator based on one random partition of the UCI benchmark text corpus for sentence categorization. Here five sentences (1)-(5) correspond to five categories: AIM, OWN, CONTRAST, BASIS, MISC, with the first five words of each sentence prespecified. All sentences are grammatically legitimate except "improves" in (4) suffers from an error, and kolmogorov in (1) and israelis in (3) should be capitalized. These errors are correctable by a grammar checker.

(1) The paper extends research on the theory of choice rules. (2) We test our predictions using the ideas and the notion of kolmogorov complexity bound on the number of examples of the data sets. (3) The results demonstrate that israelis models can be used to provide new results for classification accuracy. (4) We show that implementation concerns the performance of the learning algorithm for improves the optimal predictor of the prediction. (5) The effect of balance is described by a high level of events and the objects that are shared.