

# Learned Lossless Compression via an Extension of the Bayes Codes

Yuta Nakahara  
Center for Data Science  
Waseda University  
Tokyo, Japan  
y.nakahara@waseda.jp

Shota Saito  
Gunma University  
Faculty of Informatics  
Gunma, Japan  
shota.s@gunma-u.ac.jp

Koshi Shimada and Toshiyasu Matsushima  
Waseda University  
Dept. of Pure and Applied Math.  
Tokyo, Japan  
i\_am\_koshi@suou.waseda.jp and toshimat@waseda.jp

**Abstract**—Learned compression is an emerging scheme of data compression where the encoder and decoder are learned from data. In learned compression, various deep neural network-based methods are proposed. However, the idea of data-based code construction does not necessarily require deep neural networks. In this paper, we introduce a framework to incorporate this idea into classical codes via an extension of the Bayes codes.

## I. INTRODUCTION

Recently, a new scheme called learned compression has been used for data compression. In contrast to most classical data compression methods, in this scheme, the encoder and decoder are constructed based on a large amount of training data. This scheme is rapidly expanding due to the development of deep generative models. Various deep neural network (DNN) based coding methods have been proposed, e.g., variational auto-encoders [1] and recurrent neural networks [2].

However, the idea of data-based code construction does not necessarily require DNNs. It is also possible to incorporate this idea into classical codes. In particular, lossless compression using the Bayes codes [3] is suitable for incorporating this idea, because the Bayes codes assume a Bayesian statistical model as the data generation model, which can be hierarchically extended. In fact, we have already reported some examples in [4], [5]. In this paper, we generalize and summarize these results as a unified framework of learned compression.

We expect that statistical-model-based methods can be computationally more efficient and more easily avoid overfitting than DNN-based methods. There is also a possibility that statistical-model-based methods may be more suitable than DNN-based methods for lossless compression, where no distortion is allowed and theoretical limits of compression are well studied. In fact, many of the existing learned compression methods have been applied to lossy compression. Thus, we consider the DNN-based methods and the statistical-model-based methods should be studied with careful comparison. We hope this paper encourages the interplay between the two.

## II. PRELIMINARIES: THE BAYES CODES [3]

We consider lossless compression of a sequence  $\mathbf{x} = x_1 \dots x_n$ , which is generated from a probability distribution  $p(\mathbf{x}|\boldsymbol{\theta})$ . Here, we assume  $\boldsymbol{\theta}$  is unknown and we cannot use the true distribution  $p(\mathbf{x}|\boldsymbol{\theta})$  as a coding probability of entropy

codes such as arithmetic codes [6]. Therefore, we estimate it by  $\hat{p}(\mathbf{x})$ . In the Bayes codes [3], we assume  $\boldsymbol{\theta}$  follows a prior distribution  $p(\boldsymbol{\theta}|\boldsymbol{\eta})$ , where  $\boldsymbol{\eta}$  is a given hyper-parameter. This enables us to adopt an estimation criterion for  $p(\mathbf{x}|\boldsymbol{\theta})$ . That is the Bayes risk function (see, e.g., [7]) based on the Kullback-Leibler (KL) information between the true distribution  $p(\mathbf{x}|\boldsymbol{\theta})$  and an estimated distribution  $\hat{p}(\mathbf{x})$ . It is known that the optimal distribution minimizing this criterion is given as follows.

$$\hat{p}(\mathbf{x}) = \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\eta})d\boldsymbol{\theta}. \quad (1)$$

The code where (1) is used as a coding probability of the arithmetic code is called the Bayes code. The expected code length of the Bayes code converges to the entropy of  $p(\mathbf{x}|\boldsymbol{\theta})$  with the true  $\boldsymbol{\theta}$  for sufficiently large data length, and its convergence speed achieves a theoretical limit [8].

Further, according to [3], we can derive a sequential coding algorithm where the following probability is used as the coding probability of the arithmetic code for the  $j$ th symbol  $x_j$ .

$$\hat{p}(x_j|x^{j-1}) = \int p(x_j|x^{j-1}, \boldsymbol{\theta})p(\boldsymbol{\theta}|x^{j-1}, \boldsymbol{\eta})d\boldsymbol{\theta}, \quad (2)$$

where  $x^{j-1}$  denotes  $x_1 \dots x_{j-1}$ . It is known that the code length of the Bayes code using this coding probability coincides with that using (1). If we assume  $p(\boldsymbol{\theta})$  is a conjugate prior for  $p(\mathbf{x}|\boldsymbol{\theta})$ , we can analytically solve the integral in (2) and repeat encoding of  $x_j$  and updating of  $p(\boldsymbol{\theta}|x^j, \boldsymbol{\eta})$ .

*Remark 1:* In [9],  $p(\mathbf{x}|\boldsymbol{\theta})$  is assumed to be a context-tree model, which includes arbitrary (infinite) order Markov models as sub-models. Moreover, [9] has proposed an algorithm to calculate (2) for the context-tree model without any approximation, and its complexity is only  $O(n)$ . Therefore, the Bayes codes are not only for a trivial model like the Bernoulli model but applicable to a broad range of models.

## III. PROBLEM SETTING

In this section, we extend the Bayes codes and mathematically formulate the problem of learned compression, i.e., constructing an encoder and decoder from training sequences for a target sequence. Let  $n_i$  denote the length of the  $i$ th sequence and we define  $\mathbf{x}_i := x_{i,1} \dots x_{i,n_i}$ . Let  $\mathbf{x}^m := (\mathbf{x}_1, \dots, \mathbf{x}_m)$  be the training sequences and  $\mathbf{x}_{m+1}$  be the target sequence. We assume each sequence  $\mathbf{x}_i$  is independently generated according to a parameter  $\boldsymbol{\theta}_i$ . Let  $\boldsymbol{\theta}^m := (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m)$  be

the tuple of them. In addition, we assume the probability distribution of  $\mathbf{x}_i$  has another parameter  $\alpha_i$ . For  $\alpha_i$ , we assume it becomes difficult to calculate integrals in the coding probability (1) or (2) for  $\mathbf{x}_i$  if we assume any prior distribution on it. Therefore, we treat  $\alpha_i$  as an unknown constant. Using such a parameter, the model becomes more flexible. Consequently,  $\mathbf{x}_i$  independently follows  $p(\mathbf{x}_i|\theta_i; \alpha_i)$ . We also define a tuple  $\alpha^m := (\alpha_1, \dots, \alpha_m)$ . Moreover, we assume a prior distribution on  $\theta_i$ . It includes an additional parameter  $\beta$  that is assumed to be an unknown constant because of a similar reason to  $\alpha_i$ , and we assume  $\theta_i$  independently follows  $p(\theta_i|\eta; \beta)$ . Lastly, we assume  $\eta$  is unknown and assume a hyper-prior  $p(\eta)$ . An idea of a hyper-prior distribution also plays an important role in the DNN-based methods [1].

Under this condition, we define the problem of learned compression as a problem to construct a coding algorithm for the target sequence  $\mathbf{x}_{m+1}$  using training sequences  $\mathbf{x}_1, \dots, \mathbf{x}_m$ .

*Example 1:* In our previous study on lossless text compression [4],  $p(\mathbf{x}_i|\theta_i; \alpha_i)$  was the context-tree model,  $p(\theta_i|\eta; \beta)$  was Dirichlet distributions and a tree distribution [10], and  $p(\eta)$  was beta distributions. In our previous study on lossless image compression [5],  $p(\mathbf{x}_i|\theta_i; \alpha_i)$  was a two-dimensional autoregressive model,  $p(\theta_i|\eta; \beta)$  was a mixture of Gauss-gamma distributions, and  $p(\eta)$  was a Dirichlet distribution and Gauss-Wishart distributions. Therefore, we could assume a multimodal pixel distribution represented as a mixture of autoregressive models in [5]. Non-linear autoregressive models are also used in the DNN-based methods [2].

#### IV. LEARNED COMPRESSION ALGORITHMS

If  $\alpha_{m+1}$  and  $\beta$  are given, in a similar manner to [3], the optimal coding probability for  $\mathbf{x}_{m+1}$  is given as follows:

$$p(\mathbf{x}_{m+1}|\mathbf{x}^m; \alpha_{m+1}, \beta) = \iint p(\mathbf{x}_{m+1}|\theta_{m+1}; \alpha_{m+1}) \times p(\theta_{m+1}|\eta; \beta) p(\eta|\mathbf{x}^m; \beta) d\theta_{m+1} d\eta. \quad (3)$$

However, we have some difficulties in calculating this. First,  $\alpha_{m+1}$  and  $\beta$  are unknown. Second, the integral for  $\eta$  is often computationally expensive. Therefore, we build in pre-determined point estimators  $\hat{\alpha}_{m+1}$ ,  $\hat{\beta}$ , and  $\hat{\eta}$ .  $\hat{\beta}$  and  $\hat{\eta}$  are estimated from  $\mathbf{x}^m$  and shared between the encoder and decoder beforehand.  $\hat{\alpha}_{m+1}$  is estimated from  $\mathbf{x}_{m+1}$  and sent to the decoder as a header information. Therefore, we use the following approximate coding probability.

$$\hat{p}(\mathbf{x}_{m+1}) = \int p(\mathbf{x}_{m+1}|\theta_{m+1}; \hat{\alpha}_{m+1}) p(\theta_{m+1}|\hat{\eta}; \hat{\beta}) d\theta_{m+1}. \quad (4)$$

This method can be considered a kind of learned compression algorithm. It is because estimating  $\eta$  and  $\beta$  from  $\mathbf{x}^m$  corresponds to constructing a coding algorithm for  $\mathbf{x}_{m+1}$ . Hereafter, we call an estimation phase for  $\eta$  and  $\beta$  a learning phase and we call a phase to estimate  $\alpha_{m+1}$  and coding  $\mathbf{x}_{m+1}$  a compression phase. We will describe them in order.

##### A. Learning Phase

To estimate  $\eta$  and  $\beta$ , we can use any Bayesian estimation methods such as variational Bayesian (VB) methods and Markov chain Monte Carlo (MCMC) methods. Here, we use

the VB methods because it is usually more scalable than the MCMC methods. In particular, the VB method simultaneously estimating a deterministic parameter such as  $\alpha^m$  and  $\beta$  is called the empirical variational Bayesian (EVB) method. For more scalability, we can also use the stochastic variational inference (SVI) method [11].

In EVB method, we approximate the joint posterior distribution  $p(\theta^m, \eta|\mathbf{x}^m; \alpha^m, \beta)$  by a *variational distribution*  $q(\theta^m, \eta)$ , which satisfies the factorization property:  $q(\theta^m, \eta) = q(\theta^m)q(\eta)$ .

It is known (see, e.g., [12]) that minimizing the KL information  $\text{KL}(q(\theta^m, \eta)||p(\theta^m, \eta|\mathbf{x}^m; \alpha^m, \beta))$  is equivalent to maximizing a function called *variational lower bound*  $\text{VL}(q; \alpha^m, \beta)$ . Further, it is also known (see, e.g., [12]) that the optimal variational distribution, which maximize  $\text{VL}(q; \alpha^m, \beta)$ , fulfills the following equations:

$$\ln q^*(\theta^m) = \mathbb{E}_{q^*(\eta)} [\ln p(\theta^m, \eta, \mathbf{x}^m; \alpha^m, \beta)] + \text{const.}, \quad (5)$$

$$\ln q^*(\eta) = \mathbb{E}_{q^*(\theta^m)} [\ln p(\theta^m, \eta, \mathbf{x}^m; \alpha^m, \beta)] + \text{const.} \quad (6)$$

Simultaneously, we maximize the variational lower bound  $\text{VL}(q; \alpha^m, \beta)$  with respect to  $\alpha^m$  and  $\beta$ , and define

$$(\alpha^*)^m, \beta^* := \arg \max_{\alpha^m, \beta} \text{VL}(q; \alpha^m, \beta). \quad (7)$$

However,  $q^*(\theta^m)$ ,  $q^*(\eta)$ ,  $(\alpha^*)^m$ , and  $\beta^*$  depend on each other. Therefore, we update them in turn from any initial values until the convergence. Let  $q^{(t)}(\theta^m)$ ,  $q^{(t)}(\eta)$ ,  $(\alpha^{(t)})^m$  and  $\beta^{(t)}$  denote approximate solutions at the  $t$ th iteration. After convergence, we use the expectation or the mode of  $q^{(\infty)}(\eta)$  and  $\beta^{(\infty)}$  as  $\hat{\eta}$  and  $\hat{\beta}$ . In the following, we describe some discussions for calculating (5), (6), and (7).

1) *Update of  $q^{(t)}(\theta^m)$  and  $q^{(t)}(\eta)$ :* By calculating (5), an additional factorization is induced, and  $q^{(t)}(\theta^m) = \prod_{i=1}^m q^{(t)}(\theta_i)$  holds. Each factor is represented as follows:

$$\ln q^{(t)}(\theta_i) = \ln p(\mathbf{x}_i|\theta_i; \alpha_i^{(t-1)}) + \mathbb{E}_{q^{(t-1)}(\eta)} [\ln p(\theta_i|\eta; \beta^{(t-1)})] + \text{const.} \quad (8)$$

By calculating (6),  $q^{(t)}(\eta)$  is represented as follows:

$$\ln q^{(t)}(\eta) = \sum_{i=1}^m \mathbb{E}_{q^{(t)}(\theta_i)} [\ln p(\theta_i|\eta; \beta^{(t-1)})] + \ln p(\eta) + \text{const.} \quad (9)$$

*Remark 2:* When the components of the assumed model are exponential families and they are locally conjugate each other, in most cases, (8) and (9) have closed-form parametric representations, and their updates are represented as the updates of their parameters. So, we can efficiently calculate them.

2) *Update of  $(\alpha^{(t)})^m$  and  $\beta^{(t)}$ :* In a similar manner to general cases (e.g., [12]), the variational lower bound  $\text{VL}(q^{(t)}; (\alpha^{(t-1)})^m, \beta^{(t-1)})$  is described as follows.

$$\begin{aligned} \text{VL}(q^{(t)}; (\alpha^{(t-1)})^m, \beta^{(t-1)}) &= \sum_{i=1}^m \mathbb{E}_{q^{(t)}} [\ln p(\mathbf{x}_i|\theta_i; \alpha_i^{(t-1)})] \\ &+ \sum_{i=1}^m \mathbb{E}_{q^{(t)}} [\ln p(\theta_i|\eta; \beta^{(t-1)})] + \mathbb{E}_{q^{(t)}} [\ln p(\eta)] \\ &- \sum_{i=1}^m \mathbb{E}_{q^{(t)}} [\ln q(\theta_i)] - \mathbb{E}_{q^{(t)}} [\ln q(\eta)]. \end{aligned} \quad (10)$$

Therefore, regarding  $(\alpha^{(t-1)})^m$ , we only have to maximize  $\mathbb{E}_{q^{(t)}} [\ln p(\mathbf{x}_i|\theta_i; \alpha_i^{(t-1)})]$  for each  $\alpha_i^{(t-1)}$ , and for  $\beta^{(t-1)}$ , we

only have to maximize  $\sum_{i=1}^m \mathbb{E}_{q^{(t)}}[\ln p(\theta_i|\eta; \beta^{(t-1)})]$ , i.e., we use the following updating formulas.

$$\alpha_i^{(t)} = \arg \max_{\alpha_i} \mathbb{E}_{q^{(t)}}[\ln p(\mathbf{x}_i|\theta_i; \alpha_i)] \quad (11)$$

$$\beta^{(t)} = \arg \max_{\beta} \sum_{i=1}^m \mathbb{E}_{q^{(t)}}[\ln p(\theta_i|\eta; \beta)] \quad (12)$$

*Remark 3:* In most cases, the objective functions in (11) and (12) have similar forms to log-likelihood functions of  $\alpha_i$  and  $\beta$ . Therefore, we can apply various methods for maximum likelihood estimation.

## B. Compression Phase

1) *Calculating Header Information:* Before calculating (4), we have to estimate  $\alpha_{m+1}$  by  $\hat{\alpha}_{m+1}$ , and send it to the decoder. We maximize the following formula. This is equivalent to minimizing the code length of the Bayes code for  $\mathbf{x}_{m+1}$ .

$$\ln \int p(\mathbf{x}_{m+1}|\theta_{m+1}; \alpha_{m+1})p(\theta_{m+1}; \hat{\eta}, \hat{\beta})d\theta_{m+1} \quad (13)$$

We may use any optimization algorithm, but the expectation maximization (EM) algorithm works well for this problem. For any probability distribution  $r(\theta_{m+1})$  such that the right-hand side of (14) can be defined, it is known that the following inequality holds (see e.g., [12]).

$$\begin{aligned} & \ln \int p(\mathbf{x}_{m+1}|\theta_{m+1}; \alpha_{m+1})p(\theta_{m+1}; \hat{\eta}, \hat{\beta})d\theta_{m+1} \\ & \geq \int r(\theta_{m+1}) \ln \frac{p(\mathbf{x}_{m+1}, \theta_{m+1}; \alpha_{m+1}, \hat{\eta}, \hat{\beta})}{r(\theta_{m+1})} d\theta_{m+1} \end{aligned} \quad (14)$$

We maximize this lower bound by repeating the following E-step and M-step from an initial value.

*E-step:* Fix a tentative value  $\alpha_{m+1}^{\text{old}}$ , then (14) holds as an equality when the following holds.

$$r(\theta_{m+1}) = p(\theta_{m+1}|\mathbf{x}_{m+1}; \alpha_{m+1}^{\text{old}}, \hat{\eta}, \hat{\beta}) \quad (15)$$

Therefore, we update  $r(\theta_{m+1})$  by this formula. If we assume  $p(\theta_{m+1}|\eta; \beta)$  is a conjugate prior for  $p(\mathbf{x}_{m+1}|\theta_{m+1}; \alpha_{m+1})$ , this posterior distribution has a closed-form parametric representation, and we can efficiently calculate it in a similar manner to Remark 2.

*M-step:* Fix  $r(\theta_{m+1}) = p(\theta_{m+1}|\mathbf{x}_{m+1}; \alpha_{m+1}^{\text{old}}, \hat{\eta}, \hat{\beta})$ , then maximization of (14) is achieved when  $\alpha_{m+1}^{\text{new}}$  is as follows.

$$\begin{aligned} \alpha_{m+1}^{\text{new}} &= \arg \max_{\alpha_{m+1}} \int p(\theta_{m+1}|\mathbf{x}_{m+1}; \alpha_{m+1}^{\text{old}}, \hat{\eta}, \hat{\beta}) \\ & \quad \times \ln p(\mathbf{x}_{m+1}|\theta_{m+1}; \alpha_{m+1})d\theta_{m+1} \end{aligned} \quad (16)$$

If  $q^{(t)}(\theta_i)$  and  $p(\theta_{m+1}|\mathbf{x}_{m+1}; \alpha_{m+1}^{\text{old}}, \hat{\eta}, \hat{\beta})$  has the same form, which is often holds when we assume a conjugate prior, then (16) is equivalent to (11), and we can easily solve it.

2) *Entropy Coding:* Finally, we sequentially encode  $\mathbf{x}_{m+1}$  in a similar manner to Sec. II. More specifically, we use the following coding probability for  $x_{m+1,j}$  in the arithmetic code.

$$\begin{aligned} & \hat{p}(x_{m+1,j} | (x_{m+1})^{j-1}) \\ &= \int p(x_{m+1,j} | (x_{m+1})^{j-1}, \theta_{m+1}; \hat{\alpha}_{m+1}) \\ & \quad \times p(\theta_{m+1} | (x_{m+1})^{j-1}; \hat{\alpha}_{m+1}, \hat{\eta}, \hat{\beta})d\theta_{m+1}, \end{aligned} \quad (17)$$

where  $(x_{m+1})^{j-1}$  denote  $x_{m+1,1} \cdots x_{m+1,j-1}$ .

## V. NUMERICAL RESULTS IN PREVIOUS PAPERS

Here, we introduce numerical results obtained by applying our learned compression framework. Those results have been reported in our previous papers. In [4], we performed experiments on synthetic discrete sequences and real genome data. On synthetic data, the average code length was reduced by 2.08%. On real data, the average code length was reduced by 0.89%. In [5], we performed experiments on benchmark images, and the average code length was reduced by 5.02%.

## VI. CONCLUSION

We introduced a framework to incorporate the idea of learned compression, i.e., the idea to construct the codes from data, into classical data compression methods via an extension of the Bayes codes. We generalized some specific examples reported in our previous papers [4], [5] and summarized them as a unified framework of learned compression. We hope this framework encourages the interplay between the DNN-based methods and the statistical-model-based methods.

## ACKNOWLEDGMENT

This work was supported in part by JSPS KAKENHI Grant Number JP22K02811, JP22K14254, JP23K03863, JP23K04293, and by JST SPRING Grant Number JPMJSP2128.

## REFERENCES

- [1] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, "Variational image compression with a scale hyperprior," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=rkcQFMZrB>
- [2] G. Toderici, D. Vincent, N. Johnston, S. J. Hwang, D. Minnen, J. Shor, and M. Covell, "Full resolution image compression with recurrent neural networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5435–5443.
- [3] T. Matsushima, H. Inazumi, and S. Hirasawa, "A class of distortionless codes designed by Bayes decision theory," *IEEE Transactions on Information Theory*, vol. 37, no. 5, pp. 1288–1293, 1991.
- [4] Y. Nakahara, S. Saito, K. Shimada, and T. Matsushima, "Hyperparameter learning of bayesian context tree models," in *2023 IEEE International Symposium on Information Theory (ISIT)*, 2023, pp. 537–542.
- [5] Y. Nakahara and T. Matsushima, "Hyperparameter learning of stochastic image generative models with bayesian hierarchical modeling and its effect on lossless image coding," in *2021 IEEE Information Theory Workshop (ITW)*, 2021, pp. 1–6.
- [6] J. Rissanen and G. Langdon, "Universal modeling and coding," *IEEE Transactions on Information Theory*, vol. 27, no. 1, pp. 12–23, 1981.
- [7] J. O. Berger, *Statistical decision theory and Bayesian analysis*. Springer Science & Business Media, 2013.
- [8] B. S. Clarke and A. R. Barron, "Information-theoretic asymptotics of Bayes methods," *IEEE Transactions on Information Theory*, vol. 36, no. 3, pp. 453–471, May 1990.
- [9] T. Matsushima and S. Hirasawa, "Reducing the space complexity of a Bayes coding algorithm using an expanded context tree," in *2009 IEEE International Symposium on Information Theory*, June 2009, pp. 719–723.
- [10] Y. Nakahara, S. Saito, A. Kamatsuka, and T. Matsushima, "Probability distribution on full rooted trees," *Entropy*, vol. 24, no. 3, 2022. [Online]. Available: <https://www.mdpi.com/1099-4300/24/3/328>
- [11] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley, "Stochastic variational inference," *Journal of Machine Learning Research*, vol. 14, no. 40, pp. 1303–1347, 2013. [Online]. Available: <http://jmlr.org/papers/v14/hoffman13a.html>
- [12] C. Bishop, *Pattern Recognition and Machine Learning*. Springer, January 2006. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/pattern-recognition-machine-learning/>