
CARMS: Categorical-Antithetic-REINFORCE Multi-Sample Gradient Estimator

Alek Dimitriev

McCombs School of Business
The University of Texas at Austin
Austin, TX 78712
alek.dimitriev@mccombs.utexas.edu

Mingyuan Zhou

McCombs School of Business
The University of Texas at Austin
Austin, TX 78712
mingyuan.zhou@mccombs.utexas.edu

Abstract

Accurately backpropagating the gradient through categorical variables is a challenging task that arises in various domains, such as training discrete latent variable models. To this end, we propose CARMS, an unbiased estimator for categorical random variables based on multiple mutually negatively correlated (jointly antithetic) samples. CARMS combines REINFORCE with copula based sampling to avoid duplicate samples and reduce its variance, while keeping the estimator unbiased using importance sampling. It generalizes both the ARMS antithetic estimator for binary variables, which is CARMS for two categories, as well as LOORF/VarGrad, the leave-one-out REINFORCE estimator, which is CARMS with independent samples. We evaluate CARMS on several benchmark datasets on a generative modeling task, as well as a structured output prediction task, and find it to outperform competing methods including a strong self-control baseline. The code is publicly available.¹

1 Introduction

When optimizing an expectation based objective of the form $\mathbb{E}_{z \sim q_\phi(z)}[f(z)]$, we sometimes require the gradients with respect to the parameters ϕ of the distribution. This is challenging for discrete variables, because the commonly used reparameterization gradient does not directly work, unless the discrete distribution is approximated by a continuous and reparameterizable one [Jang et al., 2017, Maddison et al., 2017]. A significant part of this field has thus been score function (REINFORCE) based estimators [Glynn, 1990, Williams, 1992, Fu, 2006], which are general and do not require f to be the differentiable. In this paper, we focus on the case when z is a high dimensional categorical variable with logits ϕ . An common example of this form is the evidence lower bound (ELBO) [Jordan et al., 1998], which arises in variational inference and is used for training variational autoencoders [Kingma and Welling, 2014, Rezende et al., 2014]. Because their latent space consists of a large number of categorical variables, they require Monte Carlo gradients with respect to the parameters of the stochastic distribution, but have excellent performance, as a categorical VAE has in practice achieved state-of-the-art zero-shot image generation [Ramesh et al., 2021].

Our main contribution is a novel unbiased and low-variance gradient estimator for categorical variables. The Categorical-Antithetic-REINFORCE-Multi-Sample (CARMS) estimator uses a copula to generate any number of antithetic (mutually negatively correlated) [Owen, 2013] categorical samples, and constructs an unbiased estimator by combining them into a baseline for variance reduction. This approach is inspired by the ARMS estimator [Dimitriev and Zhou, 2021], which also uses multiple antithetic samples, but only works for binary variables. For two categories, CARMS

¹<https://github.com/alekdimi/carms>

reduces to it, while for independent samples, CARMS reduces to the leave-one-out-REINFORCE (LOORF) estimator. Our approach achieves higher ELBO with VAE models than the state of the art, as well as higher log likelihood for conditional image completion.

Related work One widely used group of gradient estimators for categorical variables is based on trading off bias for lower variance. This includes the straight through (ST) estimator [Bengio et al., 2013], the direct argmax [Lorberbom et al., 2019], as well as the concurrently developed equivalent Gumbel-Softmax (GS) [Jang et al., 2017] or Concrete [Maddison et al., 2017] that uses a continuous relaxation. These were further improved by combining them with REINFORCE to obtain unbiased estimators, which includes REBAR [Tucker et al., 2017], as well as RELAX [Grathwohl et al., 2018], which uses a free-form neural network.

In practice, REINFORCE is almost always augmented by baselines, *e.g.*, in variational inference [Mnih and Gregor, 2014, Ranganath et al., 2014, Paisley et al., 2012, Ruiz et al., 2016, Kucukelbir et al., 2017]. MuProp [Gu et al., 2016] uses a first order mean-field Taylor approximation, but requires f to be differentiable. Other estimators apply Rao-Blackwellization [Liu et al., 2019, Kool et al., 2020], *e.g.*, to one latent variable at a time [Titsias and Lázaro-Gredilla, 2015], to the ST estimator [Paulus et al., 2021], or to the reparameterized Dirichlet vector in ARSM [Dong et al., 2021]. Besides REINFORCE and reparameterization type gradients, there is also the measure valued gradient [Rosca et al., 2019] and finite differences [Fu, 2006], but are less common in practice. A comprehensive review can be found in Mohamed et al. [2020].

The more recent approaches for categorical variables are unbiased and REINFORCE based, building on the idea of using multiple samples to construct a baseline for variance reduction. One such baseline is LOORF [Kool et al., 2019a], originally introduced in Salimans and Knowles [2014] and also known as VarGrad [Richter et al., 2020], where its theoretical properties are further analyzed. VIMCO [Mnih and Rezende, 2016] has a similar form, but is specific to the importance weighted multi sample bound [Burda et al., 2016]. A different approach is ARSM [Yin et al., 2019], which reparameterizes the gradient with a Dirichlet distribution and uses swaps to obtain multiple correlated samples. However, it adds some amount of variance due to the continuous reparameterization, and can require up to $C(C - 1)/2$ function evaluations per step, which can be computationally expensive. More recently, the unordered set estimator (UNORD) [Kool et al., 2020] uses the Gumbel top-k trick to sample without replacement, and then constructs a baseline using a multiplicative term to preserve its unbiasedness.

2 Background

Let $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_D)$ denote D independent categorical variables, where \mathbf{z}_d is a one hot encoded categorical sample from $\mathbf{z}_d \sim q_{\phi_d}(\mathbf{z}_d) = \text{Cat}(\sigma(\phi_d))$, with $\sigma(\mathbf{x})_i = e^{\mathbf{x}_i} / \sum_j e^{\mathbf{x}_j}$ being the softmax function. Let also $\hat{\mathbf{i}}$ be the basis vector with its i^{th} coordinate set to one, and otherwise zero, such that $P(\mathbf{z} = \hat{\mathbf{i}}) = \sigma(\phi)_i$, or equivalently $\mathbb{E}[\mathbf{z}] = \sigma(\phi)$. Unless otherwise stated, superscripts denote the dimension, and subscripts denote different samples, with vectors and matrices being bold lowercase and bold uppercase symbols, respectively. We are interested in optimizing the following objective with respect to the logits $\phi = (\phi_1, \dots, \phi_D)$:

$$\mathcal{L}(\phi) = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z})}[f(\mathbf{z})], \quad q_{\phi}(\mathbf{z}) = \prod_{d=1}^D q_{\phi_d}(\mathbf{z}_d).$$

Although the score function gradient contains two terms:

$$\nabla_{\phi} \mathcal{L}(\phi) = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z})} \left[\nabla_{\phi} f(\mathbf{z}) + f(\mathbf{z}) \nabla_{\phi} \ln q_{\phi}(\mathbf{z}) \right]$$

the first term is easily estimated, so we omit the subscript in f for notational clarity, and we focus on the latter term in this work. Since CARMS generalizes the multisample LOORF estimator beyond independent samples, we review LOORF here. We also review the ARMS estimator, which uses a copula to generate antithetic samples, but is restricted to binary variables.

2.1 LOORF

Leave-one-out-REINFORCE (LOORF) [Salimans and Knowles, 2014, Kool et al., 2019a], also known as VarGrad [Richter et al., 2020], is a general score function based gradient estimator that uses N i.i.d. samples to construct a baseline for variance reduction, and is competitive with state of the art estimators. Its theoretical properties have recently been analyzed in Richter et al. [2020], where the same estimator results from minimizing the variance of the log ratio between the posterior and approximating distribution in variational inference. The only requirements for LOORF are the ability to sample $\mathbf{z} \sim q_\phi(\mathbf{z})$ and evaluate $f(\mathbf{z})$. Given N samples $\mathbf{z}_1, \dots, \mathbf{z}_N \stackrel{iid}{\sim} \text{Cat}(\sigma(\phi))$, it has the form:

$$g_{\text{LOORF}} = \frac{1}{N-1} \sum_{n=1}^N \left(f(\mathbf{z}_n) - \bar{f}(\mathbf{z}) \right) \nabla_\phi \ln q_\phi(\mathbf{z}_n) = \frac{1}{N-1} \sum_{n=1}^N \left(f(\mathbf{z}_n) - \bar{f}(\mathbf{z}) \right) (\mathbf{z}_n - \sigma(\phi_n)), \quad (1)$$

where $\bar{f}(\mathbf{z}) = \frac{1}{N} \sum_{n=1}^N f(\mathbf{z}_n)$. It is commonly used due to its simplicity and strong performance.

2.2 ARMS

The Antithetic-REINFORCE-MultiSample (ARMS) [Dimitriev and Zhou, 2021] estimator is a recent work that uses any number N of mutually negatively correlated samples. However, although unbiased, it is limited to binary variables, which motivated us to extend it to the categorical case. Since any antithetic copula in two dimensions reduces to the pair $(u, 1-u)$, it generalizes DisARM [Dong et al., 2020], independently discovered as unbiased uniform gradient (U2G) [Yin et al., 2020], which use $N = 2$ samples. ARMS achieves this generalization by using a copula [Trivedi and Zimmer, 2007], which is any multivariate distribution whose marginals are uniform random variables:

$$\mathbf{u} = (u_1, \dots, u_N) \sim \mathcal{C}_N, \quad \forall i : u_i \sim \text{Unif}(0, 1).$$

ARMS uses a Dirichlet or Gaussian copula with strong negative dependence between each dimension of \mathbf{u} . However any copula can be used instead, with the only two requirements being the ability to generate samples easily, as well as being able to evaluate the bivariate CDF $\Phi(u_i, u_j)$ of the copula, so that the debiasing term can be calculated. Given a copula sample \mathbf{u} , it uses inverse CDF sampling to convert it into N antithetic Bern(p) samples, which are simply $b_i = \mathbb{1}_{u_i < p}, \forall i$. For a set of N antithetic D -dimensional Bernoulli variables $\mathbf{b}_1, \dots, \mathbf{b}_N$ with probabilities $\sigma(\phi)$, the N -sample unbiased estimator has the following simple form:

$$g_{\text{ARMS}} = \frac{1}{N-1} \sum_{n=1}^N \left(f(\mathbf{b}_n) - \frac{1}{n} \sum_{m=1}^N f(\mathbf{b}_m) \right) \frac{\mathbf{b}_n - \sigma(\phi)}{1 - \rho}, \quad (2)$$

with $\rho = (\rho_1, \dots, \rho_D)$ and $\rho_d = \text{corr}(\mathbf{b}_i^d, \mathbf{b}_j^d)$, which is simple to compute given the bivariate CDF of the copula.

3 CARMS

The CARMS estimator has a similar form to LOORF, but requires a multiplicative term to remain unbiased. It has two requirements: an easy way to sample antithetic categorical variables, and being able to compute the bivariate probability mass function (PMF), which should be identical for any pair. For clarity, we begin by assuming the ability to do this, and derive the univariate version for two samples, which we then extend to N samples. Next, we generalize CARMS to any number of categorical variables. Lastly, in Section 3.2, we show two different ways of sampling that satisfy both conditions: inverse CDF sampling with an easily computable analytical bivariate PMF, and the Gumbel max trick combined with an empirical estimate of the PMF.

The two sample version of CARMS can be obtained by replacing the two *i.i.d.* samples in 2-LOORF with an arbitrarily correlated pair of categorical variables and an added debiasing term. For $N = 2$ samples $\mathbf{z}, \mathbf{z}' \sim \text{Cat}(\sigma(\phi))$, LOORF has the following simple form:

$$g_{\text{LOORF}}(\mathbf{z}, \mathbf{z}') = \frac{1}{2} (f(\mathbf{z}) - f(\mathbf{z}')) (\mathbf{z} - \mathbf{z}'), \quad (3)$$

which, importantly, is unbiased [Kool et al., 2019a, Dimitriev and Zhou, 2021]. This means that using a simple importance weight preserves its unbiasedness, for an arbitrary bivariate categorical distribution $(\mathbf{z}, \mathbf{z}') \sim \text{Cat}(\sigma(\phi), \sigma(\phi))$:

$$\begin{aligned}\mathbb{E} \left[g_{\text{LOORF}}(\mathbf{z}, \mathbf{z}') \frac{P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}})}{P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}})} \right] &= \sum_{i,j} P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}}) \frac{P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}})}{P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}})} g_{\text{LOORF}}(\hat{\mathbf{i}}, \hat{\mathbf{j}}) \\ &= \sum_{i,j} P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}})g_{\text{LOORF}}(\hat{\mathbf{i}}, \hat{\mathbf{j}}) = \mathbb{E}_{\mathbf{z}, \mathbf{z}' \stackrel{iid}{\sim} \text{Cat}(\sigma(\phi))} \left[g_{\text{LOORF}}(\mathbf{z}, \mathbf{z}') \right] = \nabla_{\phi} \mathcal{L}(\phi).\end{aligned}$$

We summarize the derivation of the two sample version of CARMS, which we denote as the Categorical-Antithetic-REINFORCE-Two-Sample (CARTS) estimator in the following theorem.

Theorem 1 Let $(\mathbf{z}, \mathbf{z}')$ be a sample from an arbitrary bivariate Categorical distribution with marginal distributions $\mathbf{z}, \mathbf{z}' \sim \text{Cat}(\sigma(\phi))$, and a known bivariate PMF. An unbiased estimator of $\nabla_{\phi} \mathbb{E}[f(\mathbf{z})]$ is:

$$g_{\text{CARTS}}(\mathbf{z}, \mathbf{z}') = \frac{1}{2} \left(f(\mathbf{z}) - f(\mathbf{z}') \right) (\mathbf{z} - \mathbf{z}') \mathbf{z}^T \mathbf{R} \mathbf{z}', \quad \mathbf{R}_{ij} = \frac{\sigma(\phi)_i \sigma(\phi)_j}{P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}})}. \quad (4)$$

The intuition behind using negatively correlated variables is that we want to avoid the case when $\mathbf{z} = \mathbf{z}'$, because the sample is then “wasted.” We formalize this intuition below, and defer the proof to the Appendix.

Theorem 2 Let $(\mathbf{z}, \mathbf{z}') \sim \text{Cat}(\sigma(\phi), \sigma(\phi))$. If the bivariate PMF satisfies:

$$\forall i \neq j : P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}}) \geq P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}}),$$

with strict inequality for at least one pair, then $\text{Var}[g_{\text{CARTS}}(\mathbf{z}, \mathbf{z}')] < \text{Var}[g_{\text{LOORF}}(\mathbf{z}, \mathbf{z}')]$.

We now extend CARTS to N samples, using the following identity, the proof of which can be found in Dimitriev and Zhou [2021]. It states that N -sample LOORF is equivalent to averaging 2-sample LOORF over all $\binom{N}{2}$ pairs:

$$\begin{aligned}g_{\text{LOORF}}(\mathbf{z}_1, \dots, \mathbf{z}_N) &= \frac{1}{N} \sum_{n=1}^N \left(f(\mathbf{z}_n) - \frac{1}{N} \sum_{m=1}^N f(\mathbf{z}_m) \right) (\mathbf{z}_n - \sigma(\phi)) \\ &= \frac{1}{N(N-1)} \sum_{n \neq m} \frac{1}{2} \left(f(\mathbf{z}_n) - f(\mathbf{z}_m) \right) (\mathbf{z}_n - \mathbf{z}_m) = \frac{1}{N(N-1)} \sum_{n \neq m} g_{\text{LOORF}}(\mathbf{z}_n, \mathbf{z}_m).\end{aligned}$$

With the above identity it easily follows that given N antithetic categorical samples $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]^T$, applying CARTS to all pairs results in an unbiased estimator, due to linearity of expectations:

$$\begin{aligned}\mathbb{E}[g_{\text{CARMS}}(\mathbf{Z}, \mathbf{R})] &= \mathbb{E} \left[\frac{1}{N(N-1)} \sum_{n \neq m} g_{\text{CARTS}}(\mathbf{z}_n, \mathbf{z}_m, \mathbf{R}) \right] = \mathbb{E} \left[\frac{1}{N(N-1)} \sum_{n \neq m} g_{\text{LOORF}}(\mathbf{z}_n, \mathbf{z}_m) \right] \\ &= \mathbb{E}[g_{\text{LOORF}}(\mathbf{Z})] = \nabla_{\phi} \mathcal{L}(\phi).\end{aligned}$$

We summarize CARMS in the next theorem and also rewrite it in a simpler matrix form used in our implementation. The matrix form can be obtained after some algebra, which can be found in the Appendix.

Theorem 3 Let $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]^T$ be a sample from an arbitrary N -variate Categorical distribution with identical marginal and bivariate distributions, such that $\mathbf{z}_i \sim \text{Cat}(\sigma(\phi))$, and $\mathbf{R}_{ij} = \sigma(\phi)_i \sigma(\phi)_j / P(\mathbf{z}_n = \hat{\mathbf{i}}, \mathbf{z}_m = \hat{\mathbf{j}})$. An unbiased estimator of $\nabla_{\phi} \mathbb{E}[f(\mathbf{z})]$ is:

$$g_{\text{CARMS}}(\mathbf{Z}) = \frac{1}{N(N-1)} \sum_{n \neq m} \frac{1}{2} \left(f(\mathbf{z}_n) - f(\mathbf{z}_m) \right) (\mathbf{z}_n - \mathbf{z}_m) \mathbf{z}_n^T \mathbf{R} \mathbf{z}_m'$$

Lemma 4 Let $f(\mathbf{Z}) = [f(\mathbf{z}_1), \dots, f(\mathbf{z}_N)]^T$, \circ denote the Hadamard product, $\mathbf{1}_{N \times N}$ a matrix of ones, and \mathbf{I}_N the identity matrix. Define $\mathcal{O} = \frac{1}{N-1} (\mathbf{1}_{N \times N} - \mathbf{I}_N) \circ (\mathbf{Z} \mathbf{R} \mathbf{Z}^T)$, and $\mathcal{D} = \text{diag}(\mathcal{O} \mathbf{1}_N)$,

to be a diagonal matrix. The CARMS estimator can equivalently be written in the following form:

$$g_{\text{CARMS}}(\mathbf{Z}) = \frac{1}{N} f(\mathbf{Z})^T (\mathcal{D} - \mathcal{O}) (\mathbf{Z} - \mathbf{1}_N \sigma(\phi)^T). \quad (5)$$

Furthermore, for independent samples, $\mathbf{Z}\mathcal{R}\mathbf{Z}^T = \mathbf{1}_{N \times N}$ and LOORF has the form:

$$g_{\text{LOORF}}(\mathbf{Z}) = \frac{1}{N} f(\mathbf{Z})^T \left(\mathbf{I}_{N \times N} - \frac{1}{N-1} (\mathbf{1}_{N \times N} - \mathbf{I}_N) \right) (\mathbf{Z} - \mathbf{1}_N \sigma(\phi)^T) \quad (6)$$

3.1 Multivariate CARMS

In the univariate case, Monte Carlo estimators are not necessary, as the expectation has C terms and can simply be analytically summed, but for many categorical variables, stochastic gradients are required. For D dimensions, we can sample $\mathbf{Z}^d \sim \text{Cat}(\sigma(\phi^d))$ independently and combine them into a $D \times N \times C$ tensor \mathbf{Z} , with the corresponding $D \times C \times C$ importance ratio tensor \mathcal{R} . Below, we use superscripts and subscripts to index the first and second dimension of the tensors, with \mathbf{Z}^{-d} denoting excluding the d^{th} row of \mathbf{Z} . Focusing on the d^{th} dimension, we have:

$$\begin{aligned} \nabla_{\phi^d} \mathbb{E}[f(\mathbf{Z})] &= \mathbb{E}_{\mathbf{Z}^{-d}} [\nabla_{\phi^d} \mathbb{E}_{\mathbf{Z}^d} [f(\mathbf{Z}^{-d}, \mathbf{Z}^d)]] \\ &= \mathbb{E}_{\mathbf{Z}^{-d}} \left[\mathbb{E}_{\mathbf{Z}^d} \left[\frac{1}{N(N-1)} \sum_{n \neq m} \frac{1}{2} (f(\mathbf{Z}_n^{-d}, \mathbf{Z}_n^d) - f(\mathbf{Z}_m^{-d}, \mathbf{Z}_m^d)) (\mathbf{Z}_n^d - \mathbf{Z}_m^d) \mathbf{Z}_n^{dT} \mathcal{R}^d \mathbf{Z}_m^d \right] \right] \\ &= \mathbb{E} \left[\frac{1}{N(N-1)} \sum_{n \neq m} \frac{1}{2} (f(\mathbf{Z}_n) - f(\mathbf{Z}_m)) (\mathbf{Z}_n^d - \mathbf{Z}_m^d) \mathbf{Z}_n^{dT} \mathcal{R}^d \mathbf{Z}_m^d \right]. \end{aligned}$$

Just like the univariate case, we only need N evaluations of f regardless of the dimensionality D .

3.2 Antithetic categorical variables

To use CARMS in practice, we describe two ways to generate antithetic categorical variables in this section. Both methods are based on transformations of uniform random variables, which makes them amenable to antithetic copulas, such as the Gaussian or Dirichlet copula [Dimitriev and Zhou, 2021].

3.2.1 Inverse CDF sampling

The inverse transform sampling is commonly used to transform uniform random variables, which are easy to generate, to some other distribution:

$$x \sim F_x(x) \iff u \sim \text{Unif}(0, 1), x = F_x^{-1}(u),$$

where $F_x(x)$ denotes the CDF of the desired distribution. For a categorical variable and some ordering of its probability vector $\mathbf{p} = (p_1, \dots, p_C)$ the inverse CDF approach amounts to the following algorithm. Define the left and right boundaries:

$$\mathbf{l}_i = \sum_{k=1}^{i-1} p_k, \quad \mathbf{r}_i = \sum_{k=1}^i p_k = \mathbf{l}_i + \mathbf{p}_i. \quad (7)$$

Then setting $\mathbf{z} = \hat{\mathbf{j}}$ if $u \in [\mathbf{l}_j, \mathbf{r}_j]$, where $u \sim \text{Unif}(0, 1)$ produces a categorical variable because $P(\mathbf{z} = \hat{\mathbf{j}}) = P(u \in [\mathbf{l}_j, \mathbf{r}_j]) = \mathbf{r}_j - \mathbf{l}_j = p_j$. To obtain N antithetic categorical variables, we sample $\mathbf{u} \sim \mathcal{C}_N$ from a copula and set \mathbf{z}_n in a vectorized manner. Importantly, this sampling approach allows us to easily analytically evaluate the bivariate PMF (proof in the Appendix):

$$P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}}) = \Phi_{\mathcal{C}}(\mathbf{r}_i, \mathbf{r}_j) - \Phi_{\mathcal{C}}(\mathbf{r}_i, \mathbf{l}_j) - \Phi_{\mathcal{C}}(\mathbf{l}_i, \mathbf{r}_j) + \Phi_{\mathcal{C}}(\mathbf{l}_i, \mathbf{l}_j), \quad (8)$$

where $\Phi_{\mathcal{C}}$ denotes the bivariate CDF of the copula. However, we are not guaranteed a non-zero probability for any (i, j) pair within a particular ordering. For example, if the ordering is $\mathbf{p} = (0.1, 0.2, 0.7)$ and we use a bivariate copula $(u, 1-u)$, then the probability of obtaining the pair (1,2) is zero:

$$P(\mathbf{z} = 1, \mathbf{z}' = 2) = P(u \in [0, 0.1], 1-u \in [0.1, 0.3]) = P(u \in [0, 0.1], u \in [0.7, 0.9]) = 0.$$

Algorithm 1 Antithetic inverse CDF categorical sampling

Input: Number of samples N , probabilities $\mathbf{p} = \sigma(\phi)$, copula \mathcal{C} .
 Sample $\mathbf{u} = (u_1, \dots, u_N) \sim \mathcal{C}_N$.
 Shuffle \mathbf{p} according to Eq. 9, and define the left and right boundaries \mathbf{l}, \mathbf{r} according to Eq. 7.
 Set $\forall n: z_n = \hat{j}$, where j is such that $u_n \in [\mathbf{l}_j, \mathbf{r}_j]$.
 For $i, j \in \{1, \dots, C\}$: compute $\mathcal{R}_{ij} = \mathbf{p}_i \mathbf{p}_j / P(z_n = \hat{i}, z_{n'} = \hat{j})$ according to Eq. 10.
return: $(z_1, \dots, z_N), \mathcal{R}$.

Algorithm 2 Antithetic Gumbel categorical sampling

Input: Number of samples N , probabilities $\mathbf{p} = \sigma(\phi)$, copula \mathcal{C} .
 Sample $\mathbf{u}_c = (u_{c1}, \dots, u_{cN}) \stackrel{iid}{\sim} \mathcal{C}_N$ for each category c .
 Set $\forall c, n: g_{cn} = -\ln(-\ln(u_{cn}))$.
 Convert to categorical $\forall n: z_{nc} = 1$ if $c = \text{argmax}_d g_{dn}$ else 0.
 Approximate the bivariate PMF $P(z = \hat{i}, z' = \hat{j})$ according to Eq 11.
 Set $\forall i, j \in \{1, \dots, C\}: \mathcal{R}_{ij} = \mathbf{p}_i \mathbf{p}_j / P(z = \hat{i}, z' = \hat{j})$.
return: $(z_1, \dots, z_N), \mathcal{R}$.

However, we are always guaranteed a non-zero probability for the pair $(1, C)$ (proof for the Dirichlet copula in the appendix). Therefore, we can average over $C(C - 1)/2$ orderings of \mathbf{p} , where the ordering σ_{ij} rearranges \mathbf{p} such that categories i and j are the first and last category of the vector, i.e. $\sigma_{ij}(i) = 1, \sigma_{ij}(j) = C$. Thus, any pair of categories (i, j) has non-zero probability in at least one ordering, namely the ordering σ_{ij} . There are many sets of orderings that satisfy this, but one simple way to do this is to translate the original vector \mathbf{p} i elements to the right, then swap the last element with the j^{th} , or more precisely:

$$\sigma_{ij}(k) = \begin{cases} C, & k = j \\ j, & k = C \\ (k - i + 1) \bmod C, & \text{otherwise} \end{cases}. \quad (9)$$

The bivariate PMF given the set of orderings is:

$$\begin{aligned} P(z = \hat{i}, z' = \hat{j}) &= \frac{2}{C(C-1)} \sum_{k < l} P(z = \sigma_{kl}(\hat{i}), z' = \sigma_{kl}(\hat{j}) \mid \sigma_{kl}(\mathbf{p})) \\ &\approx \frac{1}{|O|} \sum_{o \in O} P(z = \sigma_o(\hat{i}), z' = \sigma_o(\hat{j}) \mid \sigma_o(\mathbf{p})). \end{aligned} \quad (10)$$

If the number of categories C is large, we can approximate the sum by sampling a random subset of orderings O , $|O| \ll C(C - 1)/2$, as shown above. Note that we don't need to calculate Eq. 10 for all possible pairs of categories, just those that were in a sample of N . This results in $O(N^2 |O|)$ total computation per gradient, and is easily parallelizable. We summarize the inverse CDF sampling method in Algorithm 1.

3.2.2 Gumbel max sampling

It is not strictly necessary to analytically calculate the bivariate PMF. If the variance is not too large, a Monte Carlo estimation suffices, and we can use the well known Gumbel max trick [Gumbel, 1954, Jang et al., 2017, Maddison et al., 2017]:

$$z \sim \text{Cat}(\sigma(\theta)) \iff g_c \stackrel{iid}{\sim} \text{Gumbel}(0, 1), \quad z_d = \begin{cases} 1, & d = \text{argmax}_c g_c + \theta_c \\ 0, & \text{otherwise} \end{cases}$$

where, as before, z denotes the one hot encoding of the categorical variable. This is also equivalent to the following exponential racing [Ross, 2014, Caron and Teh, 2012, Zhang and Zhou, 2018, Yin et al., 2019] sampling: set $z = \hat{j}$ iff $j = \text{argmin}_i \epsilon_i$, where $\epsilon_i \sim \text{Exp}(e^{\phi_i})$.

Since a Gumbel variable can be generated using a uniform: $g = -\ln(-\ln(u))$, $u \sim \text{Unif}(0, 1)$, we use a copula to generate $\mathbf{u}_c = (u_{c1}, \dots, u_{cN})$. This produces a negative correlation between all pairs of samples for category c , and we do the same for all categories. If we let $\mathbf{Z} = [z_1, \dots, z_N]^T$ as

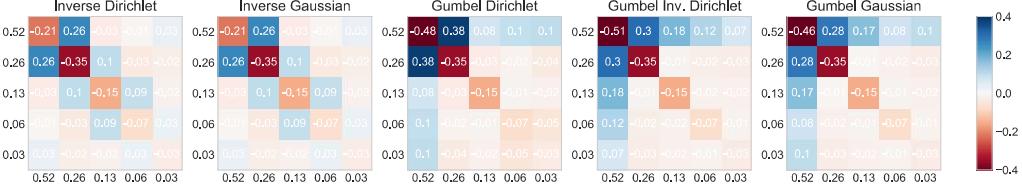


Figure 1: For an antithetic pair of one-hot encoded categorical variables \mathbf{z}, \mathbf{z}' we show the correlation between each pair \mathbf{z}_i and \mathbf{z}'_j , which are marginally Bernoulli variables with correlation $\rho_{ij} = \text{corr}(\mathbf{z}_i, \mathbf{z}_j)$. The correlation matrices shown use different categorical sampling methods or different copulas, estimated using $N = 10^3$ Monte Carlo samples.

before, a simple empirical estimate of the bivariate PMF computed using all pairs and only matrix operations is:

$$P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}}) \approx \frac{1}{N(N-1)} \mathbf{Z}^T (\mathbf{1}_{N \times N} - I_N) \mathbf{Z}. \quad (11)$$

We summarize the antithetic Gumbel sampling method in Algorithm 2, and in Fig. 1 we show an example of a bivariate PMF produced by either the Gumbel or inverse CDF method using different copulas. The results are similar using either a Dirichlet or Gaussian copula, with larger differences between the categorical sampling method. For CARMS to be unbiased we need a good estimate of the ratio $p(z_i)p(z_j)/p(z_i, z_j)$, which in turn requires a good estimate of the bivariate probability mass function (PMF) in the denominator. If we draw N samples to do this empirically, there is a small probability that some pair (i, j) does not occur at all, so in experiments we clip the maximum value of the ratio to 10 to avoid infinities.

4 Experimental results

In this section, we first illustrate the variance reduction that CARMS offers on a toy example. We also optimize a categorical variational autoencoder (VAE) [Kingma and Welling, 2014, Rezende et al., 2014], and a stochastic network for structured output prediction, which are standard tasks [Jang et al., 2017] for categorical variables, done on three different benchmark datasets. Each experiment uses both antithetic categorical approaches for CARMS: inverse CDF sampling and the Gumbel max trick, denoted as CARMS-I and CARMS-G, respectively, and we use the Dirichlet copula for both. We compare our approach to three state-of-the-art unbiased estimators: LOORF/VarGrad [Kool et al., 2019a, Richter et al., 2020], the unordered set estimator (UNORD) [Kool et al., 2020], and ARSM [Yin et al., 2019]. The code for all experiments is freely available².

4.1 Toy example

We first showcase the variance reduction over other methods in a simple toy example, where we take the gradient with respect to the logits ϕ of:

$$\mathcal{L}(\phi) = \mathbb{E}[f(\mathbf{z}_1, \dots, \mathbf{z}_D)], \quad \mathbf{z}_d \sim \text{Cat}(\sigma(\phi_d)), \quad f(\mathbf{z}_1, \dots, \mathbf{z}_D) = \sum_{d=1}^D \sum_{c=1}^C d \cdot c \cdot z_{dc}.$$

For simplicity, let the number of categories, dimensions, and samples be $C = D = N = 3$. The probabilities are randomly sampled from a Dirichlet distribution: $\sigma(\phi) \sim \text{Dir}(\mathbf{1}_C \cdot \alpha)$, but are identical for all methods for a given α . We vary the entropy of the probabilities from high ($\alpha = 1$) to low ($\alpha = 1000$). In a high entropy setting, there is little difference between the estimators, as the variance itself is very high, but differences emerge as we increase α and lower the entropy. The combined log variance of the gradient of each logit, for different methods and different α , is shown in Fig 2.

²<https://github.com/alekdimi/carms>

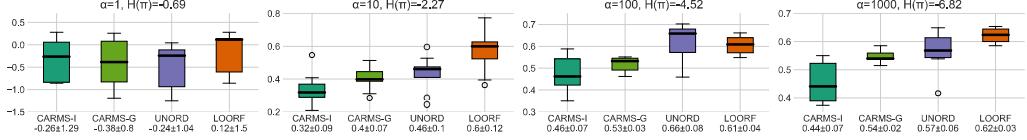


Figure 2: Log variance of the gradient of different estimators with respect to the logits on a toy problem. Columns correspond to different entropy levels of the logits.

4.2 Categorical variational autoencoder

For this task, we follow the experimental setting from ARMS [Dimitriev and Zhou, 2021], except we use categorical instead of binary latent variables, and maximize the ELBO:

$$\text{ELBO}(\phi) = \mathbb{E} \left[\ln \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \approx \sum_{n=1}^N \ln \frac{p(\mathbf{x}|\mathbf{z}_n)p(\mathbf{z}_n)}{q_\phi(\mathbf{z}_n|\mathbf{x})}, \quad z_1, \dots, z_N \sim \text{Cat}_N(\sigma(\phi)),$$

where $\text{Cat}_N(\sigma(\phi))$ denotes an N -variate categorical distribution with identical marginals. The number of categories is $C \in \{3, 5, 10\}$ with $D = \lfloor 200/C \rfloor$ latent variables, respectively, to make the total computational effort similar. In the binary case $C = 2$, CARMS reduces to ARMS, for which a thorough comparison has already been produced. The task is training a categorical VAE using either a linear or nonlinear encoder/decoder pair on three different datasets: Dynamic(ally binarized) MNIST [LeCun et al., 2010], Fashion MNIST [Xiao et al., 2017], and Omniglot [Lake et al., 2015]. All datasets are freely available under the MIT license, and do not contain any personally identifiable information or offensive content. For a fair comparison, all methods use the same learning rate, optimizer, model architecture, and number of samples. Since ARSM uses a variable number of function evaluations per step, we use one sample per step, for which ARSM uses around twice as many evaluations as the other methods. The results are combined from five independent runs for each experimental configuration.

The VAE consists of a stochastic layer with $\lfloor 200/C \rfloor$ units, each of which is a C -way categorical variable. For the nonlinear case, there are additionally two layers of 200 units with LeakyReLU [Maas et al., 2013] activations. The prior logits are optimized using SGD with a learning rate of 10^{-2} , whereas the encoder and decoder are optimized using Adam [Kingma and Ba, 2015] with a learning rate of 10^{-4} , following Yin et al. [2019]. The optimization is run for 10^6 steps, with a batch size of 50, from which the global dataset mean is subtracted. The models are trained on an Intel Xeon Platinum 8280 2.7GHz CPU, and an individual run takes approximately 16 hours on one core of the machine, with a total carbon emissions estimated to be 28.19 kg of CO₂ [Lacoste et al., 2019]. An exception is UNORD for 10 samples, which was significantly more heavy in computation. Although the paper states that a step can be performed in $O(2^C)$, the provided code requires $O(C!)$ evaluations per step, forcing us to use $N = 8$ samples at most.

In Fig 3, we plot the training and test log likelihood using 100 samples, and gradient variance w.r.t. the logits over time, for a nonlinear network on dynamic MNIST. Similar plots for other datasets and network types can be found in the Appendix. Also shown in Table 1 is the final training log likelihood using 100 samples for all three datasets, network types, categories, and gradient estimators. The corresponding table with the final test log likelihood can be found in the Appendix. In general, both versions of CARMS perform comparably, and result in slightly higher log likelihood than other methods. Because Gumbel CARMS uses an empirical estimate of the debiasing ratio, it has higher variance and slightly worse performance. When limiting the number of function evaluations, there is a gap between ARSM (which used on average $2C$ evaluations per step, out of a maximum of $C(C - 1)/2$ per step) compared to the other methods, which used C evaluations. This is possibly due to the continuous reparameterization that ARSM uses, which adds variance.

4.3 Structured prediction with stochastic categorical networks

We also compare all methods on the standard benchmark task of predicting the lower half of an image from the upper half, *i.e.*, the conditional distribution $p(\mathbf{x}_l | \mathbf{x}_u)$, where \mathbf{x}_u and \mathbf{x}_l denote the upper and lower half of an image, respectively. We use a stochastic categorical network to estimate this

Table 1: Final training 100 sample log likelihood of VAEs using different estimators, where the stochastic layer contains $C=3, 5$, or 10 categories, with $\lfloor 200/C \rfloor$ latent variables and C samples per gradient step, respectively. Results are reported on three datasets: Dynamic MNIST, Fashion MNIST, and Omniglot over 5 runs, with the best performing methods in bold.

	Categories	CARMS-I	CARMS-G	LOORF	UNORD	ARSM
Dynamic MNIST	Linear	3 -105.34 ± 0.25	-105.36 ± 0.24	-105.64 ± 0.23	-105.23 ± 0.23	-107.35 ± 0.56
		5 -103.53 ± 0.13	-103.35 ± 0.18	-103.54 ± 0.15	-103.50 ± 0.13	-106.13 ± 0.53
		10 -103.22 ± 0.05	-103.12 ± 0.06	-103.48 ± 0.06	-103.56 ± 0.06	-106.71 ± 0.58
	Nonlinr	3 -94.85 ± 0.28	-94.60 ± 0.28	-95.12 ± 0.21	-95.21 ± 0.22	-99.62 ± 0.50
		5 -93.05 ± 0.14	-92.60 ± 0.12	-92.91 ± 0.16	-92.98 ± 0.12	-98.89 ± 0.43
		10 -92.13 ± 0.05	-92.42 ± 0.10	-92.44 ± 0.04	-93.05 ± 0.14	-97.76 ± 0.41
Fashion MNIST	Linear	3 -245.44 ± 0.22	-245.69 ± 0.19	-245.8 ± 0.19	-245.90 ± 0.21	-247.51 ± 0.45
		5 -242.06 ± 0.13	-241.90 ± 0.10	-242.17 ± 0.09	-242.43 ± 0.12	-244.63 ± 0.47
		10 -240.44 ± 0.03	-240.52 ± 0.04	-240.91 ± 0.04	-241.08 ± 0.06	-243.29 ± 0.36
	Nonlinr	3 -233.13 ± 0.16	-233.20 ± 0.16	-233.75 ± 0.10	-233.34 ± 0.12	-237.93 ± 0.20
		5 -231.72 ± 0.09	-231.67 ± 0.13	-232.01 ± 0.05	-232.19 ± 0.09	-237.29 ± 0.34
		10 -230.77 ± 0.05	-231.16 ± 0.05	-231.35 ± 0.03	-231.74 ± 0.02	-235.88 ± 0.17
Omniglot	Linear	3 -114.53 ± 0.12	-114.73 ± 0.15	-114.90 ± 0.13	-114.77 ± 0.15	-116.34 ± 0.42
		5 -114.01 ± 0.10	-113.93 ± 0.11	-114.01 ± 0.10	-114.00 ± 0.09	-115.72 ± 0.35
		10 -114.97 ± 0.06	-114.99 ± 0.06	-115.17 ± 0.05	-115.55 ± 0.08	-117.48 ± 0.35
	Nonlinr	3 -110.19 ± 0.23	-110.16 ± 0.21	-110.33 ± 0.27	-110.31 ± 0.21	-115.13 ± 0.51
		5 -109.14 ± 0.09	-109.35 ± 0.13	-109.39 ± 0.16	-109.55 ± 0.14	-115.07 ± 0.37
		10 -108.66 ± 0.08	-108.62 ± 0.07	-108.98 ± 0.06	-109.54 ± 0.15	-115.11 ± 0.34

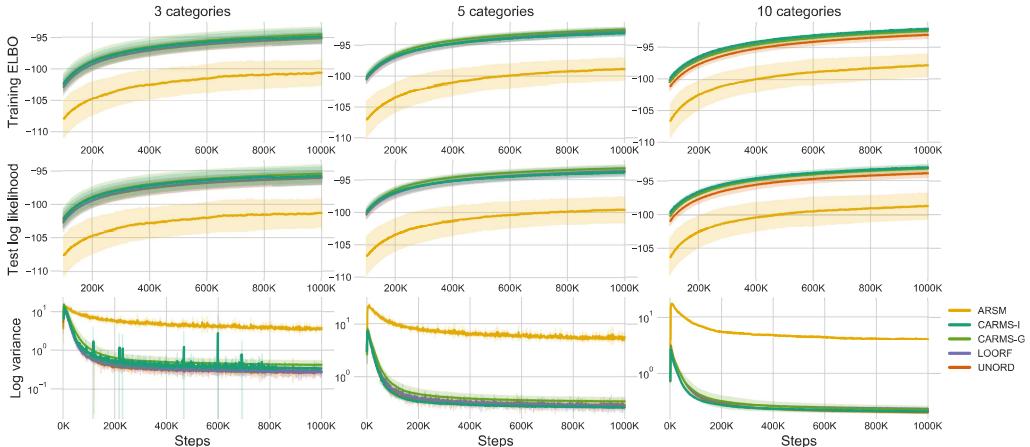


Figure 3: Training a nonlinear categorical VAE with different estimators on Dynamic MNIST using ELBO. Columns correspond to $C \in \{3, 5, 10\}$ categories with C samples per gradient step, respectively. Rows correspond to the 100 sample training and test log likelihood, and the variance of the gradient with respect to the logits of the encoder network. Results for different datasets and other networks can be found in the Appendix.

distribution, with the objective: $\mathbb{E}_{z \sim p(z_m | x_u)} \left[\frac{1}{M} \sum_{m=1}^M \ln p(x_l | z_m) \right]$, where z denotes a stochastic categorical layer. The encoder/decoder pair each contains one hidden layer with $\lfloor 200/C \rfloor$ latent variables and a LeakyReLU activation, with the optimization performed for $C \in \{3, 5, 10\}$ categories, on all three datasets, with identical settings for each gradient estimator for a fair comparison. We use $M = 1$ for training, and $M = 1000$ for evaluation on both the train and test set. In Table 2, we show the final training set log likelihood, with the corresponding test log likelihood table in the appendix, though the results are qualitatively similar. The results are similar to the VAE experiment, with the inverse CDF CARMS having a slightly higher log likelihood. However, the differences between estimators are less pronounced, with the unordered set estimator is being mostly on par with CARMS, and ARSM only slightly trailing the other methods, with a much smaller gap.

Table 2: Final training log likelihood of a categorical network for conditional estimation using different gradient estimators, where the stochastic layer contains $C = 3, 5$, or 10 categories, with $\lfloor 200/C \rfloor$ latent variables and C samples per gradient step, respectively. Results are reported on three datasets: Dynamic MNIST, Fashion MNIST, and Omniglot over 5 runs, with the best performing methods in bold.

	Categories	CARMS-I	CARMS-G	LOORF	UNORD	ARSM
Dynamic MNIST	3	57.98 ± 0.10	58.35 ± 0.14	58.19 ± 0.14	58.06 ± 0.04	60.22 ± 0.19
	5	57.57 ± 0.05	57.85 ± 0.12	57.78 ± 0.08	57.6 ± 0.05	59.38 ± 0.13
	10	58.17 ± 0.26	58.33 ± 0.13	58.20 ± 0.14	58.18 ± 0.17	59.32 ± 0.11
Fashion MNIST	3	132.83 ± 0.08	132.90 ± 0.08	133.10 ± 0.05	133.06 ± 0.06	134.56 ± 0.35
	5	132.68 ± 0.05	132.81 ± 0.12	132.91 ± 0.07	132.94 ± 0.14	134.09 ± 0.12
	10	133.32 ± 0.15	133.43 ± 0.23	133.54 ± 0.17	133.38 ± 0.10	134.02 ± 0.21
Omniglot	3	65.57 ± 0.10	66.05 ± 0.14	65.92 ± 0.26	65.81 ± 0.09	68.00 ± 0.09
	5	65.65 ± 0.18	66.16 ± 0.32	65.92 ± 0.23	65.78 ± 0.06	67.99 ± 0.32
	10	66.76 ± 0.31	66.94 ± 0.24	66.87 ± 0.07	66.66 ± 0.24	68.35 ± 0.16

5 Discussion

We have presented a novel approach for training categorical variables, which extends the ARMS estimator for the binary case. It goes beyond i.i.d. samples by using a copula to generate antithetic categorical samples, but preserves unbiasedness by including a multiplicative term. For i.i.d. samples, the form of the estimator reduces to LOORF. We showcase its usefulness on several datasets, a different number of categories and types of deep neural networks. In variational inference tasks, and conditional estimation, CARMS outperforms other state of the art estimators. The main limitation of this work is its specificity to categorical (including binary) variables. We hope to extend this, e.g. to Plackett-Luce models for top-k sampling [Kool et al., 2019b, Grover et al., 2019], which has important applications in ranking [Dadaneh et al., 2020]. There is also a general limitation that CARMS shares with other state-of-the-art estimators, which is higher complexity than LOORF, a very simple but strong baseline. Future work includes investigating theoretical properties and large scale applications, as well as possible general antithetic gradient estimators, and we plan to investigate adaptive correlations that take into account the properties of f to further reduce the variance.

Potential societal impact This work is focused on better optimization for categorical variables, which includes any network containing a stochastic categorical layer. In particular, generative models such as VAEs are widely used, and are sometimes trained on more human centered datasets. These models can sometimes be used to impersonate a person’s face or voice. We only used non-human datasets such as MNIST, but with enough knowledge and using the available code, anyone, including bad actors, can train the same model on human content for malicious purposes. However, we strongly believe that wide knowledge dissemination and open source code is crucial for reproducibility in all of science.

Acknowledgments

The authors acknowledge the support of NSF IIS-1812699, the APX 2019 project sponsored by the Office of the Vice President for Research at The University of Texas at Austin, the support of a gift fund from ByteDance Inc., and the Texas Advanced Computing Center (TACC) for providing HPC resources that have contributed to the research results reported within this paper.

References

- Y. Bengio, N. Léonard, and A. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- Y. Burda, R. B. Grosse, and R. Salakhutdinov. Importance weighted autoencoders. In *4th International Conference on Learning Representations, ICLR*, 2016.
- F. Caron and Y. W. Teh. Bayesian nonparametric models for ranked data. In *Advances in Neural Information Processing Systems*, volume 25, 2012.

- S. Z. Dadaneh, S. Boluki, M. Zhou, and X. Qian. Arsm gradient estimator for supervised learning to rank. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3157–3161. IEEE, 2020.
- A. Dimitriev and M. Zhou. Arms: Antithetic-reinforce-multi-sample gradient for binary variables. *International Conference on Machine Learning*, 2021.
- Z. Dong, A. Mnih, and G. Tucker. Disarm: An antithetic gradient estimator for binary latent variables. In *Advances in Neural Information Processing Systems 33*, 2020.
- Z. Dong, A. Mnih, and G. Tucker. Coupled gradient estimators for discrete latent variables. *Third Symposium on Advances in Approximate Bayesian Inference*, 2021.
- M. C. Fu. Gradient estimation. *Handbooks in operations research and management science*, 13: 575–616, 2006.
- P. W. Glynn. Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 33(10):75–84, 1990.
- W. Grathwohl, D. Choi, Y. Wu, G. Roeder, and D. Duvenaud. Backpropagation through the void: Optimizing control variates for black-box gradient estimation. In *6th International Conference on Learning Representations, ICLR*, 2018.
- A. Grover, E. Wang, A. Zweig, and S. Ermon. Stochastic optimization of sorting networks via continuous relaxations. *arXiv preprint arXiv:1903.08850*, 2019.
- S. Gu, S. Levine, I. Sutskever, and A. Mnih. Muprop: Unbiased backpropagation for stochastic neural networks. In *4th International Conference on Learning Representations, ICLR*, 2016.
- E. J. Gumbel. *Statistical theory of extreme values and some practical applications: a series of lectures*, volume 33. US Government Printing Office, 1954.
- E. Jang, S. Gu, and B. Poole. Categorical reparameterization with gumbel-softmax. In *5th International Conference on Learning Representations, ICLR*. OpenReview.net, 2017.
- M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. In *Learning in graphical models*, pages 105–161. Springer, 1998.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR*, 2015.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR*, 2014.
- W. Kool, H. van Hoof, and M. Welling. Buy 4 REINFORCE samples, get a baseline for free! In *Workshop, Deep Reinforcement Learning Meets Structured Prediction, ICLR*, 2019a.
- W. Kool, H. Van Hoof, and M. Welling. Stochastic beams and where to find them: The gumbel-top-k trick for sampling sequences without replacement. In *International Conference on Machine Learning*, pages 3499–3508. PMLR, 2019b.
- W. Kool, H. van Hoof, and M. Welling. Estimating gradients for discrete random variables by sampling without replacement. In *International Conference on Learning Representations*, 2020.
- A. Kucukelbir, D. Tran, R. Ranganath, A. Gelman, and D. M. Blei. Automatic differentiation variational inference. *The Journal of Machine Learning Research*, 18(1):430–474, 2017.
- A. Lacoste, A. Lucioni, V. Schmidt, and T. Dandres. Quantifying the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*, 2019.
- B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. *ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist*, 2, 2010.

- W. Lee and J. Y. Ahn. On the multidimensional extension of countermonotonicity and its applications. *Insurance: Mathematics and Economics*, 56:68–79, 2014.
- R. Liu, J. Regier, N. Tripuraneni, M. Jordan, and J. McAuliffe. Rao-blackwellized stochastic gradients for discrete distributions. In *International Conference on Machine Learning*, pages 4023–4031. PMLR, 2019.
- G. Lorberbom, T. S. Jaakkola, A. Gane, and T. Hazan. Direct optimization through arg max for discrete variational auto-encoder. In *Advances in Neural Information Processing Systems 32*, pages 6200–6211, 2019.
- A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML Workshop on Deep Learning for Audio, Speech and Language Processing*. Citeseer, 2013.
- C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation of discrete random variables. In *5th International Conference on Learning Representations, ICLR*. OpenReview.net, 2017.
- A. J. McNeil and J. Nešlehová. Multivariate archimedean copulas, d-monotone functions and l_1 -norm symmetric distributions. *The Annals of Statistics*, 37(5B):3059–3097, 2009.
- A. Mnih and K. Gregor. Neural variational inference and learning in belief networks. In *International Conference on Machine Learning*, pages 1791–1799. PMLR, 2014.
- A. Mnih and D. J. Rezende. Variational inference for monte carlo objectives. In *Proceedings of the 33rd International Conference on Machine Learning, ICML*, volume 48, pages 2188–2196. JMLR.org, 2016.
- S. Mohamed, M. Rosca, M. Figurnov, and A. Mnih. Monte carlo gradient estimation in machine learning. *J. Mach. Learn. Res.*, 21:132:1–132:62, 2020.
- A. B. Owen. Monte carlo theory, methods and examples, 2013.
- J. W. Paisley, D. M. Blei, and M. I. Jordan. Variational bayesian inference with stochastic search. In *Proceedings of the 29th International Conference on Machine Learning, ICML*, 2012.
- M. B. Paulus, C. J. Maddison, and A. Krause. Rao-blackwellizing the straight-through gumbel-softmax gradient estimator. In *International Conference on Learning Representations*, 2021.
- A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*, 2021.
- R. Ranganath, S. Gerrish, and D. Blei. Black box variational inference. In *Artificial intelligence and statistics*, pages 814–822. PMLR, 2014.
- D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.
- L. Richter, A. Boustati, N. Nüsken, F. J. Ruiz, and Ö. D. Akyildiz. Vargrad: A low-variance gradient estimator for variational inference. In *Advances in Neural Information Processing Systems*, volume 33, pages 13481–13492, 2020.
- M. Rosca, M. Figurnov, S. Mohamed, and A. Mnih. Measure-valued derivatives for approximate bayesian inference. *4th workshop on Bayesian Deep Learning, NeurIPS*, 2019.
- S. M. Ross. *Introduction to probability models*. Academic press, 2014.
- F. J. R. Ruiz, M. K. Titsias, and D. M. Blei. The generalized reparameterization gradient. In *Advances in Neural Information Processing Systems 29*, pages 460–468, 2016.
- T. Salimans and D. A. Knowles. On using control variates with stochastic approximation for variational bayes and its connection to stochastic linear regression. *arXiv preprint arXiv:1401.1022*, 2014.

- M. Titsias and M. Lázaro-Gredilla. Local expectation gradients for black box variational inference. In *Advances in neural information processing systems*, pages 2620–2628. Citeseer, 2015.
- P. K. Trivedi and D. M. Zimmer. *Copula modeling: an introduction for practitioners.* ., 2007.
- G. Tucker, A. Mnih, C. J. Maddison, D. Lawson, and J. Sohl-Dickstein. REBAR: low-variance, unbiased gradient estimates for discrete latent variable models. In *Advances in Neural Information Processing Systems 30*, pages 2627–2636, 2017.
- R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arxiv.org/abs/1708.07747*, 2017.
- M. Yin, Y. Yue, and M. Zhou. ARSM: augment-reinforce-swap-merge estimator for gradient backpropagation through categorical variables. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pages 7095–7104. PMLR, 2019.
- M. Yin, N. Ho, B. Yan, X. Qian, and M. Zhou. Probabilistic Best Subset Selection by Gradient-Based Optimization. *arXiv e-prints*, 2020.
- Q. Zhang and M. Zhou. Nonparametric bayesian lomax delegate racing for survival analysis with competing risks. *arXiv preprint arXiv:1810.08564*, 2018.

Appendix for CARMS: Categorical-Antithetic-REINFORCE Multi-Sample Gradient Estimator

A Background

A.1 Copulas

Although any multivariate random variable with uniform marginal distributions is commonly referred to as a copula, strictly speaking, a copula is specifically its CDF. We are interested in copulas with strong negative dependence between any pair of its variables, which we refer to as antithetic in this work. They are also referred to as countermonotonic copulas when investigated in other work [Lee and Ahn, 2014, McNeil and Nešlehová, 2009]. However, only in the bivariate case do these copulas achieve the theoretical lower bound for negative dependence, the Fréchet–Hoeffding bound:

$$\mathcal{C}(u_1, \dots, u_N) \geq \max \left\{ 1 - N + \sum_{i=1}^N u_i, 0 \right\}.$$

Furthermore, it has been shown [McNeil and Nešlehová, 2009] that no copula can achieve this lower bound for more than two dimensions. However, the Dirichlet copula used in this work is in the family of Archimedean copulas, and Lee and Ahn [2014] show that within this family, the Dirichlet copula has the strongest countermonotonicity. Its CDF is defined as:

$$\mathcal{C}(u_1, \dots, u_N) = \left(\max \left\{ 0, u_1^{\frac{1}{N-1}} + \dots + u_N^{\frac{1}{N-1}} - (N-1) \right\} \right)^{N-1}.$$

A.2 Copula sampling

We use the Dirichlet copula described in ARMS [Dimitriev and Zhou, 2021], which transforms a Dirichlet vector with concentration $\alpha = 1$ to a copula sample. Besides its countermonotonic properties, we choose this copula because both its univariate and bivariate marginal CDFs are analytically tractable. The former is required for transforming the Dirichlet vector $\mathbf{d} \sim \text{Dir}(\alpha)$ into uniform variables: $\mathbf{u} = 1 - (1 - \mathbf{d})^{N-1}$. For Eq. 8, we want to analytically calculate the bivariate joint for categorical sampling. This requires the bivariate CDF, which has the form:

$$\begin{aligned} P(u_i < p, u_j < q) &= P(d_i < 1 - (1-p)^{1/(N-1)}, d_j < 1 - (1-q)^{1/(N-1)}) \\ &= p + q - 1 + \max(0, (1-p)^{1/(N-1)} + (1-q)^{1/(N-1)}) \end{aligned}$$

A.3 Pair equivalent definition of LOORF

Because the LOORF for N samples for any distribution can be decomposed into all pairs, we can reuse the theorem, and we reproduce the short algebra needed to show it below:

$$\begin{aligned} g_{\text{loorf}}(\mathbf{b}) &= \frac{1}{n} \sum_{i=1}^n \left(f(\mathbf{b}_i) - \frac{1}{n-1} \sum_{j \neq i} f(\mathbf{b}_j) \right) \nabla_{\phi} \ln p(\mathbf{b}_i) \\ &= \frac{1}{n} \sum_{i=1}^n f(\mathbf{b}_i) \nabla_{\phi} \ln p(\mathbf{b}_i) - \frac{1}{n(n-1)} \sum_{i=1}^n \nabla_{\phi} \ln p(\mathbf{b}_i) \sum_{j \neq i} f(\mathbf{b}_j) \\ &= \frac{1}{n(n-1)} \sum_{i \neq j} \left(f(\mathbf{b}_i) \nabla_{\phi} \ln p(\mathbf{b}_i) - f(\mathbf{b}_i) \nabla_{\phi} \ln p(\mathbf{b}_j) \right) \\ &= \frac{1}{n(n-1)} \sum_{i \neq j} \frac{1}{2} (f(\mathbf{b}_i) - f(\mathbf{b}_j)) (\nabla_{\phi} \ln p(\mathbf{b}_i) - \nabla_{\phi} \ln p(\mathbf{b}_j)) \\ &= \frac{1}{n(n-1)} \sum_{i \neq j} g_{\text{pod}}(b_i, b_j). \end{aligned}$$

B Proofs

B.1 Proof of Theorem 2

Define $\Delta f_{ij} = f(\mathbf{z}_i) - f(\mathbf{z}_j)$. Using the assumption that $P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}}) \geq P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}})$, and starting with CARTS we have:

$$\begin{aligned}\text{Var}[g_{\text{CARTS}}] &= \sum_{i \neq j} P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}}) \left(\frac{1}{2} \Delta f_{ij} (\hat{\mathbf{i}} - \hat{\mathbf{j}}) \frac{P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}})}{P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}})} \right)^2 \\ &= \sum_{i \neq j} P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}}) \left(\frac{1}{2} \Delta f_{ij} (\hat{\mathbf{i}} - \hat{\mathbf{j}}) \right)^2 \frac{P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}})}{P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}})} \\ &< \sum_{i \neq j} P(\mathbf{z} = \hat{\mathbf{i}})P(\mathbf{z}' = \hat{\mathbf{j}}) \left(\frac{1}{2} \Delta f_{ij} (\hat{\mathbf{i}} - \hat{\mathbf{j}}) \right)^2 = \text{Var}[g_{\text{LOORF}}]\end{aligned}$$

Below is an example that satisfies the conditions. Let $\mathbf{p} = (0.6, 0.3, 0.1)$, for which we show the independent vs one possible antithetic pmf:

$$\text{pmf}_{\text{indep}} = \mathbf{p}\mathbf{p}^T = \begin{bmatrix} 0.36 & 0.18 & 0.06 \\ 0.18 & 0.09 & 0.03 \\ 0.06 & 0.03 & 0.01 \end{bmatrix} \quad \text{pmf}_{\text{anti}} = \begin{bmatrix} 0.30 & 0.24 & 0.06 \\ 0.24 & 0.02 & 0.04 \\ 0.06 & 0.04 & 0.01 \end{bmatrix}.$$

Because every off-diagonal element of pmf_{anti} is larger or equal to the corresponding independent pmf value, the variance is guaranteed to be no larger:

$$\begin{aligned}\text{Var}(g_{\text{CARTS}}) &= 2 \left(0.24 \Delta f_{12}[1, 1, 0]^T \left(\frac{0.18}{0.24} \right)^2 + 0.06 \Delta f_{13}[1, 0, 1]^T + 0.04 \Delta f_{23}[0, 1, 1]^T \left(\frac{0.03}{0.04} \right)^2 \right) \\ &= 2 \left(0.18 \Delta f_{12}[1, 1, 0]^T \frac{0.18}{0.24} + 0.06 \Delta f_{13}[1, 0, 1]^T + 0.03 \Delta f_{23}[0, 1, 1]^T \frac{0.03}{0.04} \right) \\ &\leq 2 (0.18 \Delta f_{12}[1, 1, 0]^T + 0.06 \Delta f_{13}[1, 0, 1]^T + 0.03 \Delta f_{23}[0, 1, 1]^T) = \text{Var}(g_{\text{LOORF}}).\end{aligned}$$

B.2 Proof of Lemma 4

First, note that the $(ij)^{\text{th}}$ element of $\mathcal{D} - \mathcal{O}$ is $\left(\sum_{j \neq i} \mathcal{R}_{ij} \right) - \mathcal{R}_{ij}$. Then we can explicitly write out the summation, and distribute it into all pairs, to arrive at the form of Theorem 3:

$$\begin{aligned}g_{\text{CARMS}} &= \frac{1}{N} f(\mathbf{Z})^T (\mathcal{D} - \mathcal{O}) (\mathbf{Z} - \mathbf{1}_N \sigma(\phi)^T) \\ &= \frac{1}{N} \sum_{i=1}^N f(\mathbf{Z})_i^T \left(\sum_{j \neq i} \mathcal{R}_{ij} \right) - \left(\sum_{j \neq i} f(\mathbf{Z})_j^T \mathcal{R}_{ij} \right) (\mathbf{Z}_i - \sigma(\phi)) \\ &= \frac{1}{N(N-1)} \sum_{i \neq j} \left(f(\mathbf{Z})_i - f(\mathbf{Z})_j \right) (\mathbf{z}_i - \mathbf{z}_j) \mathbf{z}_i^T \mathcal{R}_{ij} \mathbf{z}_j\end{aligned}$$

B.3 Proof of Equation 8

$$\begin{aligned}P(\mathbf{z} = \hat{\mathbf{i}}, \mathbf{z}' = \hat{\mathbf{j}}) &= P(u \in [\mathbf{l}_i, \mathbf{r}_i], u' \in [\mathbf{l}_j, \mathbf{r}_j]) \\ &= P(u \leq \mathbf{r}_i, u' \in [\mathbf{l}_j, \mathbf{r}_j]) - P(u \leq \mathbf{l}_i, u' \in [\mathbf{l}_j, \mathbf{r}_j]) \\ &= P(u \leq \mathbf{r}_i, u' \leq \mathbf{r}_j) - P(u \leq \mathbf{r}_i, u' \leq \mathbf{l}_j) \\ &\quad - P(u \leq \mathbf{l}_i, u' \leq \mathbf{r}_j) + P(u \leq \mathbf{l}_i, u' \leq \mathbf{l}_j) \\ &= \Phi_{\mathcal{C}}(\mathbf{r}_i, \mathbf{r}_j) - \Phi_{\mathcal{C}}(\mathbf{r}_i, \mathbf{l}_j) - \Phi_{\mathcal{C}}(\mathbf{l}_i, \mathbf{r}_j) + \Phi_{\mathcal{C}}(\mathbf{l}_i, \mathbf{l}_j).\end{aligned}$$

B.4 Proof that the bivariate PMF for the pair $(\mathbf{1}, \mathbf{C})$ is non-zero

We show that joint probability of the first and last category $P(\mathbf{z} = \hat{\mathbf{1}}, \mathbf{z}' = \hat{\mathbf{C}})$ is always positive, because the Dirichlet copula has an area of non-zero density in a subset of this region. If we take $\epsilon = \min(p_1, p_C) > 0$, then:

$$\begin{aligned} P(\mathbf{z} = \hat{\mathbf{1}}, \mathbf{z}' = \hat{\mathbf{C}}) &= P(u < p_1, u' > 1 - p_C) \geq P(u < \epsilon, u' > 1 - \epsilon) \\ &= P(u < \epsilon) - P(u < \epsilon, u' < 1 - \epsilon) = \epsilon - \max \left\{ 0, \epsilon^{\frac{1}{N-1}} + (1 - \epsilon)^{\frac{1}{N-1}} - 1 \right\}^{N-1}. \end{aligned}$$

If the second term is zero, it trivially holds that the probability is non-zero since $\epsilon > 0$. Otherwise, note that for any integer $N > 1$:

$$\begin{aligned} \epsilon > 0 &\iff 1 - \epsilon < 1 \iff (1 - \epsilon)^{\frac{1}{N-1}} < 1 \iff (1 - \epsilon)^{\frac{1}{N-1}} - 1 < 0 \\ &\iff \epsilon^{\frac{1}{N-1}} + (1 - \epsilon)^{\frac{1}{N-1}} - 1 < \epsilon^{\frac{1}{N-1}} \iff \left(\epsilon^{\frac{1}{N-1}} + (1 - \epsilon)^{\frac{1}{N-1}} - 1 \right)^{N-1} < \epsilon \\ &\iff \epsilon - \left(\epsilon^{\frac{1}{N-1}} + (1 - \epsilon)^{\frac{1}{N-1}} - 1 \right)^{N-1} > 0, \end{aligned}$$

which concludes the proof.

C Additional results

Table 3: Final test log likelihood of VAEs using different estimators, where the stochastic layer contains $C=3, 5$, or 10 categories, with $\lfloor 200/C \rfloor$ latent variables and C samples per gradient step, respectively. Results are reported on three datasets: Dynamic MNIST, Fashion MNIST, and Omniglot over 5 runs.

	Categories	CARMS-I	CARMS-G	LOORF	UNORD	ARSM
Dynamic MNIST	Linear	3 -104.82 ± 0.25	3 -104.9 ± 0.25	3 -105.17 ± 0.24	3 -104.73 ± 0.24	3 -106.85 ± 0.57
		5 -103.12 ± 0.13	5 -102.9 ± 0.18	5 -103.16 ± 0.16	5 -103.06 ± 0.13	5 -105.73 ± 0.54
		10 -103.0 ± 0.04	10 -102.88 ± 0.06	10 -103.19 ± 0.06	10 -103.35 ± 0.08	10 -106.41 ± 0.59
	Nonlinr	3 -95.73 ± 0.32	3 -95.44 ± 0.31	3 -95.98 ± 0.25	3 -96.09 ± 0.25	3 -101.33 ± 0.54
		5 -93.83 ± 0.17	5 -93.27 ± 0.13	5 -93.67 ± 0.18	5 -93.69 ± 0.13	5 -99.61 ± 0.48
		10 -92.97 ± 0.07	10 -93.26 ± 0.11	10 -93.22 ± 0.04	10 -93.94 ± 0.15	10 -98.73 ± 0.44
Fashion MNIST	Linear	3 -247.46 ± 0.22	3 -247.77 ± 0.18	3 -247.82 ± 0.19	3 -247.92 ± 0.21	3 -249.5 ± 0.46
		5 -244.16 ± 0.13	5 -244.02 ± 0.1	5 -244.27 ± 0.1	5 -244.55 ± 0.13	5 -246.69 ± 0.48
		10 -242.6 ± 0.03	10 -242.69 ± 0.04	10 -243.1 ± 0.04	10 -243.27 ± 0.06	10 -245.43 ± 0.37
	Nonlinr	3 -235.81 ± 0.18	3 -235.89 ± 0.17	3 -236.47 ± 0.11	3 -236.06 ± 0.13	3 -240.46 ± 0.21
		5 -234.37 ± 0.09	5 -234.27 ± 0.04	5 -234.7 ± 0.05	5 -234.81 ± 0.09	5 -239.83 ± 0.34
		10 -233.39 ± 0.05	10 -233.81 ± 0.05	10 -234.02 ± 0.04	10 -234.41 ± 0.03	10 -238.55 ± 0.18
Omniglot	Linear	3 -115.15 ± 0.13	3 -115.31 ± 0.15	3 -115.46 ± 0.14	3 -115.33 ± 0.16	3 -116.92 ± 0.43
		5 -114.88 ± 0.1	5 -114.86 ± 0.12	5 -114.89 ± 0.11	5 -114.86 ± 0.1	5 -116.56 ± 0.37
		10 -116.31 ± 0.05	10 -116.41 ± 0.05	10 -116.55 ± 0.06	10 -116.92 ± 0.11	10 -118.79 ± 0.39
	Nonlinr	3 -114.54 ± 0.31	3 -114.4 ± 0.23	3 -114.6 ± 0.35	3 -114.5 ± 0.27	3 -119.0 ± 0.57
		5 -113.18 ± 0.12	5 -113.39 ± 0.16	5 -113.44 ± 0.21	5 -113.69 ± 0.18	5 -119.41 ± 0.41
		10 -112.71 ± 0.13	10 -112.7 ± 0.03	10 -112.93 ± 0.07	10 -113.8 ± 0.17	10 -120.05 ± 0.42

Table 4: Final test log likelihood of a categorical network for conditional estimation using different gradient estimators, where the stochastic layer contains $C=3, 5$, or 10 categories, with $\lfloor 200/C \rfloor$ latent variables and C samples per gradient step, respectively. Results are reported on three datasets: Dynamic MNIST, Fashion MNIST, and Omniglot over 5 runs.

	Categories	CARMS-I	CARMS-G	LOORF	UNORD	ARSM
Dynamic MNIST	3	59.61 ± 0.17	59.89 ± 0.14	59.78 ± 0.2	59.72 ± 0.06	60.92 ± 0.2
	5	59.43 ± 0.03	59.69 ± 0.1	59.59 ± 0.18	59.5 ± 0.09	60.32 ± 0.16
	10	59.93 ± 0.22	60.01 ± 0.16	59.95 ± 0.09	59.93 ± 0.19	60.25 ± 0.15
Fashion MNIST	3	134.92 ± 0.12	135.02 ± 0.07	135.16 ± 0.05	135.18 ± 0.08	136.17 ± 0.37
	5	134.81 ± 0.1	134.94 ± 0.13	135.05 ± 0.06	135.1 ± 0.11	135.78 ± 0.13
	10	135.35 ± 0.13	135.44 ± 0.24	135.61 ± 0.18	135.47 ± 0.1	135.81 ± 0.2
Omniglot	3	72.19 ± 0.11	72.45 ± 0.08	72.33 ± 0.15	72.34 ± 0.1	72.73 ± 0.12
	5	72.32 ± 0.14	72.56 ± 0.14	72.43 ± 0.11	72.59 ± 0.08	72.79 ± 0.24
	10	72.79 ± 0.12	72.88 ± 0.06	72.9 ± 0.06	72.94 ± 0.17	72.98 ± 0.15

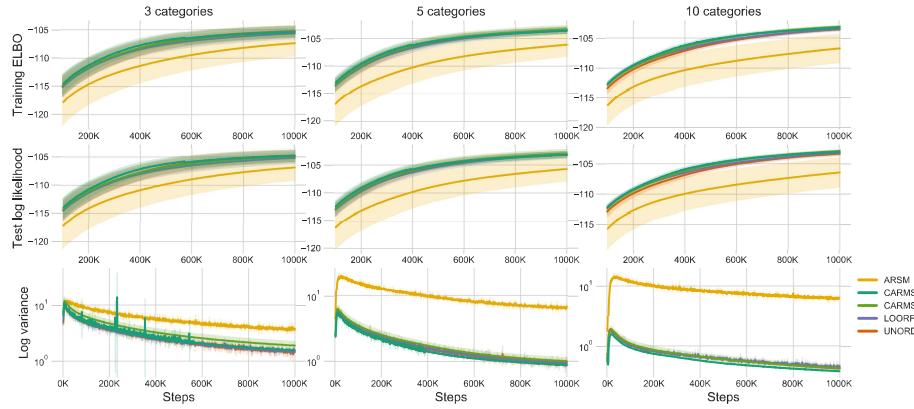


Figure 4: Training a linear categorical VAE with different estimators on Dynamic MNIST using ELBO. Columns correspond to $C \in \{3, 5, 10\}$ categories with C samples per gradient step, respectively. Rows correspond to the 100 sample training and test log likelihood, and the variance of the gradient with respect to the logits of the encoder network.

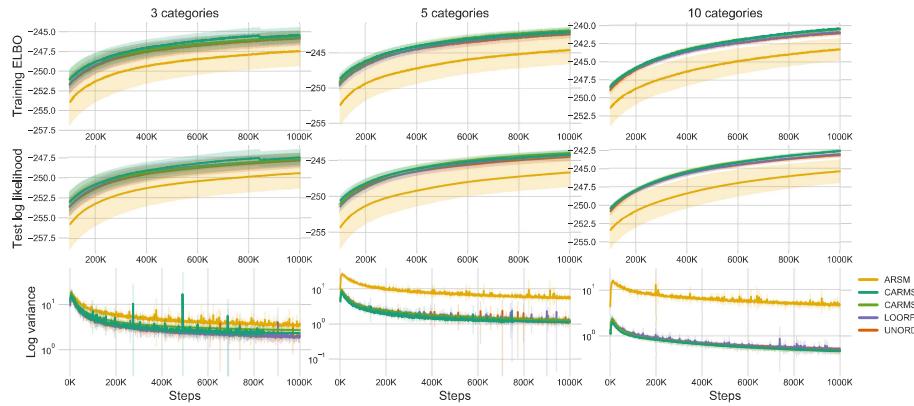


Figure 5: Training a linear categorical VAE with different estimators on Fashion MNIST using ELBO. Columns correspond to $C \in \{3, 5, 10\}$ categories with C samples per gradient step, respectively. Rows correspond to the 100 sample training and test log likelihood, and the variance of the gradient with respect to the logits of the encoder network.

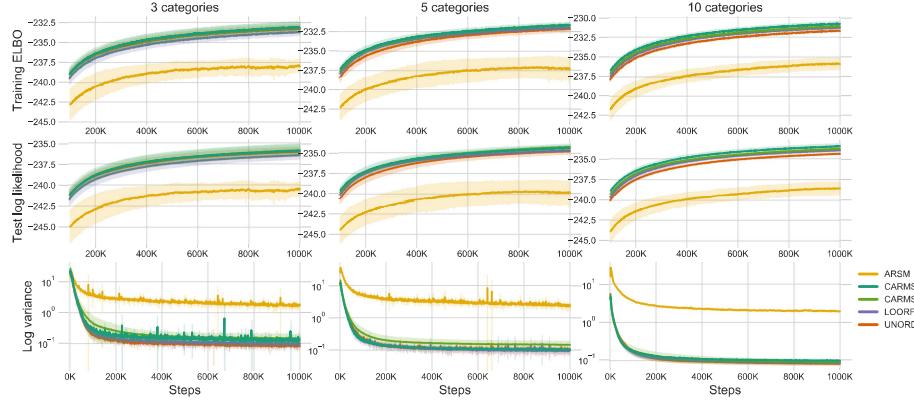


Figure 6: Training a nonlinear categorical VAE with different estimators on Fashion MNIST using ELBO. Columns correspond to $C \in \{3, 5, 10\}$ categories with C samples per gradient step, respectively. Rows correspond to the 100 sample training and test log likelihood, and the variance of the gradient with respect to the logits of the encoder network.

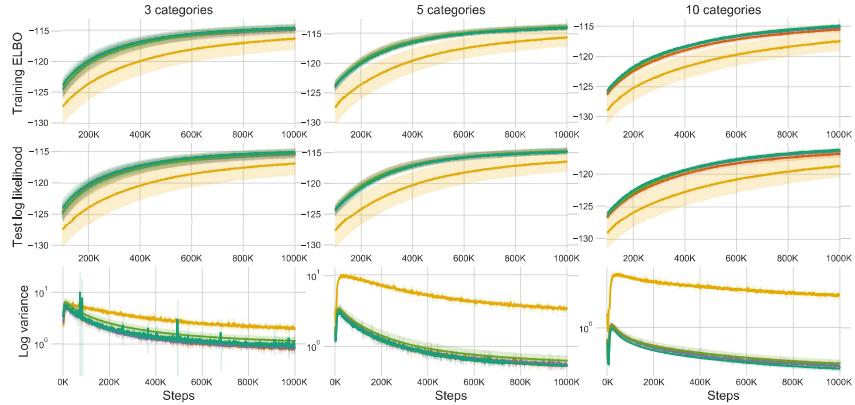


Figure 7: Training a linear categorical VAE with different estimators on Omniglot using ELBO. Columns correspond to $C \in \{3, 5, 10\}$ categories with C samples per gradient step, respectively. Rows correspond to the 100 sample training and test log likelihood, and the variance of the gradient with respect to the logits of the encoder network.

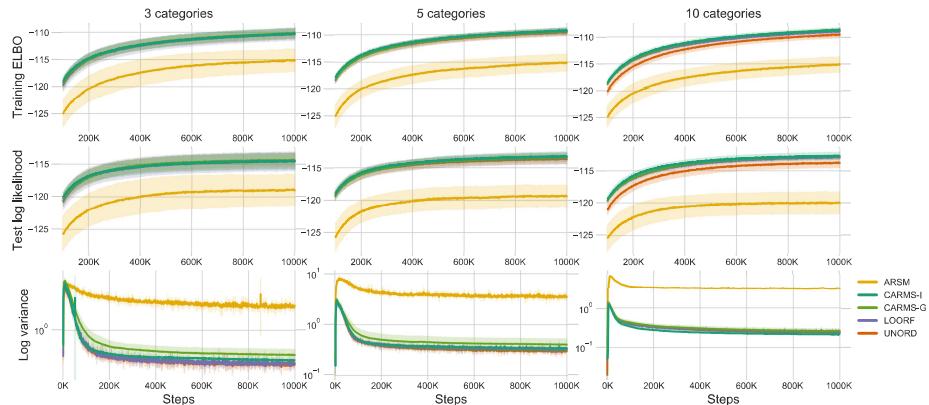


Figure 8: Training a nonlinear categorical VAE with different estimators on Omniglot using ELBO. Columns correspond to $C \in \{3, 5, 10\}$ categories with C samples per gradient step, respectively. Rows correspond to the 100 sample training and test log likelihood, and the variance of the gradient with respect to the logits of the encoder network.