

## A IMPLEMENTING GRADIENT ESTIMATORS BY MODIFYING BACKPROPAGATION

An advantage of the GRMC-K estimator is the ease with which it can be implemented using automatic differentiation software. Here, we provide a pseudo code template for such an implementation.

```
class GRMCK(Function):
    def forward(logits, tau, k):
        sample = sampleOnehotCategorical(logits)
        save_for_backward(sample, logits, tau, k)
        return sample

    def backward(grad_output):
        sample, logits, tau, k = self.saved_objects
        logZ = logsumexp(logits)
        maxgumbel = getGumbel(logZ, k)
        tgumbels = getTruncatedGumbel(
            logits, k, sample, maxgumbel)
        gumbels = mergeGumbels(
            maxgumbel, tgumbels, sample)
        J = getSmaxJacobian(gumbels + logits).mean(0)
        return grad_output.matmul(J)
```

## B IMPLEMENTING GRADIENT ESTIMATORS WITH THE SURROGATE LOSS FRAMEWORK

In this section, we consider an alternative framework for implementing the gradient estimators presented in the main body. This framework is due to (Schulman et al., 2015) and known as the surrogate loss framework. The key idea is that after the forward pass through a stochastic computation graph, all sampling decisions have been taken. Therefore, any gradient can be written as resulting from the differentiation of a surrogate objective in a deterministic computation graph.

Our exposition in the main body only considered a simplified scenario with a single discrete random variable. Therefore, we present here two cases, involving a layer of multiple and a cascade of discrete random variables. These two cases are general, because any case can be reduced to either of these two or a combination of them.

For ease of exposition, we again do not consider any direct dependence of  $f$  on the parameters of interest  $\theta$ . The extension to this case is straight-forward and follows from basic calculus.

We also introduce the following notation to denote the stop of gradient flow. For  $X^* = \text{stop\_gradient}(X)$  indicates that the gradient flow is interrupted at  $X$  and no gradient information is passed backward.

### B.1 PARALLEL CASE

Let  $D^1, \dots, D^m$  be a sequence of independent random variables. For  $j \leq m$ , let  $D^j$  be a discrete random variable  $D^j \in \{0, 1\}^n$  in a one-hot encoding,  $\sum D_i^j = 1$ , with distribution given by  $p_{\theta^j}(D^j) \propto \exp(D^j \theta)$  where  $\theta^j \in \mathbb{R}^n$ . Further, let  $S_\tau^j$  be defined analogously to equation (3). Given a continuously differentiable  $f: \mathbb{R}^{mn} \rightarrow \mathbb{R}$ , we wish to minimize

$$\min_{\theta} \mathbb{E} [f(D^1, \dots, D^m)], \quad (14)$$

where the expectation is taken over all  $m$  random variables.

In this setting,  $\nabla_{\text{REINF}}$  can be computed by differentiating the following surrogate objec-

tive,

$$f(D^{1*}, \dots, D^{m*}) \sum_{j=1}^m \log p_{\theta^j}(D^j) \quad (15)$$

In this setting,  $\nabla_{\text{GS}}$  can be computed by differentiating the following surrogate objective,

$$f(S_\tau^1, \dots, S_\tau^m) \quad (16)$$

In this setting,  $\nabla_{\text{ST}}$  can be computed by differentiating the following surrogate objective,

$$\sum_{j=1}^m \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^j} \right)^* \text{softmax}_\tau(\theta^j) \quad (17)$$

In this setting,  $\nabla_{\text{STGS}}$  can be computed by differentiating the following surrogate objective,

$$\sum_{j=1}^m \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^j} \right)^* S_\tau^j \quad (18)$$

In this setting,  $\nabla_{\text{GRMCK}}$  can be computed by differentiating the following surrogate objective,

$$\sum_{j=1}^m \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^j} \right)^* \left[ \frac{1}{K} \sum_{k=1}^K S_\tau^{jk} \right] \quad (19)$$

## B.2 SEQUENTIAL CASE

Let  $D^1, \dots, D^m$  be a sequence of non-independent random variables. For  $j \leq m$ , let  $D^j$  be a discrete random variable  $D^j \in \{0, 1\}^n$  in a one-hot encoding,  $\sum D_i^j = 1$ , with distribution given by  $p_{\theta^j}(D^j) \propto \exp(D^{jT} \theta^j)$  where  $\theta^j \in \mathbb{R}^n$ . For  $2 \leq j \leq m$ , let  $\theta^j = h(D^{j-1})$ , where  $h: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a continuously differentiable function. Given a continuously differentiable  $f: \mathbb{R}^{mn} \rightarrow \mathbb{R}$ , we wish to minimize

$$\min_{\theta} \mathbb{E} [f(D^1, \dots, D^m)] \quad (20)$$

In this setting,  $\nabla_{\text{REINF}}$  and  $\nabla_{\text{GS}}$  can be computed by differentiating the surrogate objective given in the parallel case.

In this setting,  $\nabla_{\text{ST}}$  can be computed by differentiating the following surrogate objective,

$$L_m := \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^m} \right)^* \text{softmax}_\tau(\theta^m) \quad (21)$$

$$L_j := \left( \left( \frac{dL_{j+1}(D^1, \dots, D^m)}{dD^j} \right)^* + \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^j} \right)^* \right) \text{softmax}_\tau(\theta^j) \quad (22)$$

In this setting,  $\nabla_{\text{STGS}}$  can be computed by differentiating the following surrogate objective,

$$L_m := \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^m} \right)^* \text{softmax}_\tau(\theta^m + G^m) \quad (23)$$

$$L_j := \left( \left( \frac{dL_{j+1}(D^1, \dots, D^m)}{dD^j} \right)^* + \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^j} \right)^* \right) \text{softmax}_\tau(\theta^j + G^j) \quad (24)$$

In this setting,  $\nabla_{\text{GRMCK}}$  can be computed by differentiating the following surrogate objective,

$$L_m := \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^m} \right)^* \left[ \frac{1}{K} \sum_{k=1}^K (\text{softmax}_\tau(\theta^m + G^{mk})) \right] \quad (25)$$

$$L_j := \left( \left( \frac{dL_{j+1}(D^1, \dots, D^m)}{dD^j} \right)^* + \left( \frac{\partial f(D^1, \dots, D^m)}{\partial D^j} \right)^* \right) \left[ \frac{1}{K} \sum_{k=1}^K (\text{softmax}_\tau(\theta^j + G^{jk})) \right] \quad (26)$$

## C PROOFS FOR THE PROPOSITIONS

In this section, we provide derivations for all the propositions given in the main body.

### C.1 PROPOSITION 1

The derivation is based on Jensen’s inequality and the law of iterated expectations.

*Proof.*

$$\mathbb{E} [\|\nabla_{\text{GR}} - \nabla_{\theta}\|^2] = \mathbb{E} [\|\mathbb{E} [\nabla_{\text{STGS}}|D] - \nabla_{\theta}\|^2] \quad (27)$$

$$= \mathbb{E} [\|\mathbb{E} [\nabla_{\text{STGS}} - \nabla_{\theta}|D]\|^2] \quad (28)$$

$$\leq \mathbb{E} [\mathbb{E} [\|\nabla_{\text{STGS}} - \nabla_{\theta}\|^2|D]] \quad (29)$$

$$= \mathbb{E} [\|\nabla_{\text{STGS}} - \nabla_{\theta}\|^2] \quad (30)$$

□

The inequality is strict whenever  $\text{var} [\nabla_{\text{STGS}}|D] > 0$ , which is the case if  $\tau < \infty$  and  $|\theta_i| < \infty$  for all  $i \leq n$ .

### C.2 PROPOSITION 2

The derivation is based on Jensen’s inequality and the linearity of expectations. For ease of exposition, denote by  $\nabla_{\text{STGS}}(S^k|D)$  a particular realization of the ST-GS estimator for a given  $D$ .

*Proof.*

$$\mathbb{E} [\|\nabla_{\text{GRMCK}} - \nabla_{\theta}\|^2] = \mathbb{E} \left[ \left\| \frac{1}{K} \sum_{k=1}^K \nabla_{\text{STGS}}(S^k|D) - \nabla_{\theta} \right\|^2 \right] \quad (31)$$

$$= \mathbb{E} \left[ \left\| \frac{1}{K} \sum_{k=1}^K (\nabla_{\text{STGS}}(S^k|D) - \nabla_{\theta}) \right\|^2 \right] \quad (32)$$

$$\leq \mathbb{E} \left[ \frac{1}{K} \sum_{k=1}^K \|\nabla_{\text{STGS}}(S^k|D) - \nabla_{\theta}\|^2 \right] \quad (33)$$

$$= \frac{1}{K} \sum_{k=1}^K \mathbb{E} [\|\nabla_{\text{STGS}}(S^k|D) - \nabla_{\theta}\|^2] \quad (34)$$

$$= \mathbb{E} [\|\nabla_{\text{STGS}} - \nabla_{\theta}\|^2] \quad (35)$$

□

The inequality is strict whenever  $K > 1$  and  $\text{var} [\nabla_{\text{STGS}}|D] > 0$ , which is the case if  $\tau < \infty$  and  $|\theta_i| < \infty$  for all  $i \leq n$ .

### C.3 PROPOSITION 3

The derivation is based on the law of total variance.

*Proof.*

$$\text{var} [\bar{\nabla}_{\text{GRMCK}}^{1:B}] = \mathbb{E} [\text{var} [\bar{\nabla}_{\text{GRMCK}}^{1:B} | D_{1:B}, X_{1:B}]] + \text{var} [\mathbb{E} [\bar{\nabla}_{\text{GRMCK}}^{1:B} | D_{1:B}, X_{1:B}]] \quad (36)$$

$$= \mathbb{E} \left[ \text{var} \left[ \frac{1}{B} \sum_{b=1}^B \nabla_{\text{GRMCK}}^b \middle| D_b, X_b \right] \right] + \text{var} \left[ \mathbb{E} \left[ \frac{1}{B} \sum_{b=1}^B \nabla_{\text{GRMCK}}^b \middle| D_b, X_b \right] \right] \quad (37)$$

$$= \mathbb{E} \left[ \frac{1}{B^2} \sum_{b=1}^B \text{var} [\nabla_{\text{GRMCK}}^b | D_b, X_b] \right] + \text{var} \left[ \mathbb{E} \left[ \frac{1}{B} \sum_{b=1}^B \nabla_{\text{GRMCK}}^b \middle| D_b, X_b \right] \right] \quad (38)$$

$$= \frac{1}{B} \mathbb{E} [\text{var} [\nabla_{\text{GRMCK}} | D, X]] + \text{var} \left[ \frac{1}{B} \sum_{b=1}^B \mathbb{E} [\nabla_{\text{GRMCK}}^b | D_b, X_b] \right] \quad (39)$$

$$= \frac{1}{B} \mathbb{E} \left[ \frac{1}{K} \text{var} [\nabla_{\text{STGS}} | D, X] \right] + \frac{1}{B} \text{var} [\mathbb{E} [\nabla_{\text{GRMCK}} | D, X]] \quad (40)$$

$$= \frac{1}{BK} \mathbb{E} [\text{var} [\nabla_{\text{STGS}} | D, X]] + \frac{1}{B} \text{var} [\nabla_{\text{GR}}] \quad (41)$$

□

## D EXPERIMENTAL DETAILS

### D.1 UNSUPERVISED PARSING ON LISTOPS

For our unsupervised parsing experiment on ListOps, we use the basic version of the model described in Choi et al. (2017) with an embedding dimension and hidden dimension of 128. We do not use the *leaf-rnn*. We do not use the *intra-attention module*. We do not use dropout, but set weight decay to be  $1e-4$ . Because our interest is in using this experiment primarily as a testbed to evaluate the effectiveness of different gradient estimators for this model at different temperatures and for trees of different depth, we use a very simple experimental set-up. We rely on stochastic gradient descent without momentum to train all models. We use grid search to determine an optimal learning rate from  $\{0.1, 0.2, \dots, 1.0\}$  and set the temperature  $\tau$  to be in  $\{0.01, 0.1, 1.0\}$ . We repeat five independent random runs at each setting and report the mean over the five runs. We train for ten epochs and set the batch size to be equal to the maximum sequence length  $L$ .

Havrylov et al. (2019) also consider unsupervised parsing on ListOps with a variant of the model in Choi et al. (2017). They achieve near perfect accuracy, albeit in a highly customized experimental set-up. We list the most important differences below:

- Havrylov et al. (2019) does not use single-evaluation estimators, we do: They report near perfect accuracy only when using the self-critical baseline. This baseline requires an additional forward pass. All their single-evaluation results are in a similar ballpark as ours accounting for the additional differences below.
- Havrylov et al. (2019) uses extensive hyperparameter tuning, we do not: They tune learning rate, learning rate schedule, weight decay, entropy regularisation, variance reduction hyperparameters, optimizer (Adadelata), number of updates for PPO, leaf transformations and train for 300 epochs. In contrast, we only tune the (constant) learning rate via gridsearch and use SGD (see above) for each temperature and train for ten epochs.
- Havrylov et al. (2019) uses customized training procedures, we are simply plug-in: They use PPO, gradient normalization, different control variates and entropy regularization. We simply plug our estimator into the model from Choi et al. (2017).
- Havrylov et al. (2019) uses a model with more parameters than us: We do not use any leaf LSTM. It improves performance (Choi et al., 2017), but may also confound tree learning (e.g., leaf-LSTM may learn to solve the task, making tree obsolete), so we do not use it.
- Havrylov et al. (2019) uses more training data than us: We reserved 10% of the training set for validation, while they use less than 2%.

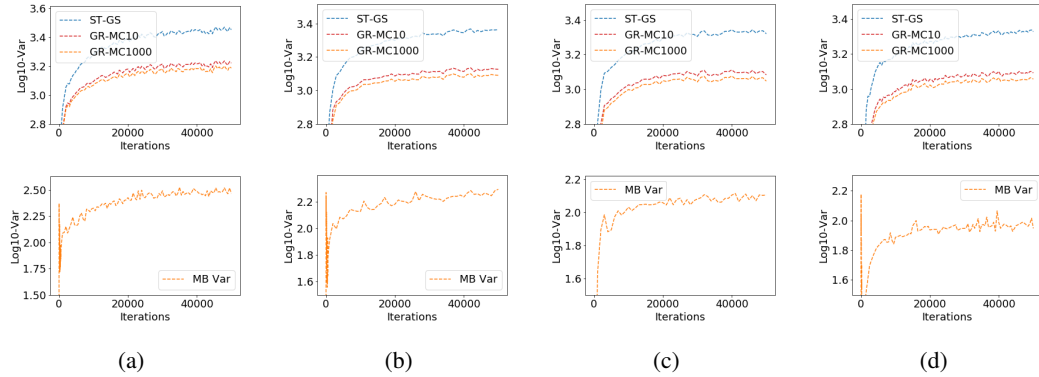


Figure 3: Our estimator (GR-MCK) effectively reduces the variance over the entire training trajectory at all arities. The variance reduction compares favorably to the minibatch variance. Columns correspond to arities, i.e. (a) binary, (b) 4-ary, (c) 8-ary, (d) 16-ary. First row, log10-trace of MC covariance matrix for various gradient estimators over iterations. Second row, log10-trace of MB covariance matrix over iterations (same for all gradient estimators).

## D.2 GENERATIVE MODELLING WITH VARIATIONAL AUTO-ENCODERS

We trained variational auto-encoders with  $n$ -ary discrete random variables with values on the corners of the hypercube  $\{-1, 1\}^{\log_2(n)}$ . The model with arity  $\{2, 4, 8, 16\}$  included  $\{240, 120, 80, 60\}$  random variables respectively.

All models were optimized using stochastic gradient descent with momentum for 50000 steps on minibatches of size 20 and 200 respectively. Hyperparameters were randomly sampled and the best setting was selected from twenty independent runs. Learning rate and momentum were randomly sampled from  $\{5, 6, \dots, 50\} \times 10^{-4}$  and  $(0, 1)$  respectively. We did not anneal the learning rate during training. For regularising the network, we used weight-decay, which was randomly sampled from  $\{0, 10^{-1}, 10^{-2}, \dots, 10^{-6}\}$ . The temperature was randomly sampled from  $[0.1, 1.0]$  and not annealed throughout training.

All models were evaluated on the validation and test set using the importance-weighted bound on the log-likelihood described in Burda et al. (2015) with 5000 samples.

To estimate the variance of a gradient estimator in the VAE experiment we used 5000 randomly sampled minibatches of size 20, for each of which we performed 100 independent forward passes and then computed the associated gradient for the parameters of the inference network. We then summed the variance to get a single scalar measurement.

To estimate the bias of a gradient estimator in the VAE experiment, we proceeded as above to approximate the expectation for a gradient estimator. We approximated the true gradient by following this procedure for the REINFORCE algorithm.

To assess training speed, we measured the average number of iterations needed to achieve a pre-specified loss threshold on the validation set. In particular, we ran multiple independent runs under the same experimental conditions for all gradient estimators. Among only runs that achieved the threshold within the total budget, we report the average number of iterations taken to cross the threshold.

## E ADDITIONAL FIGURES

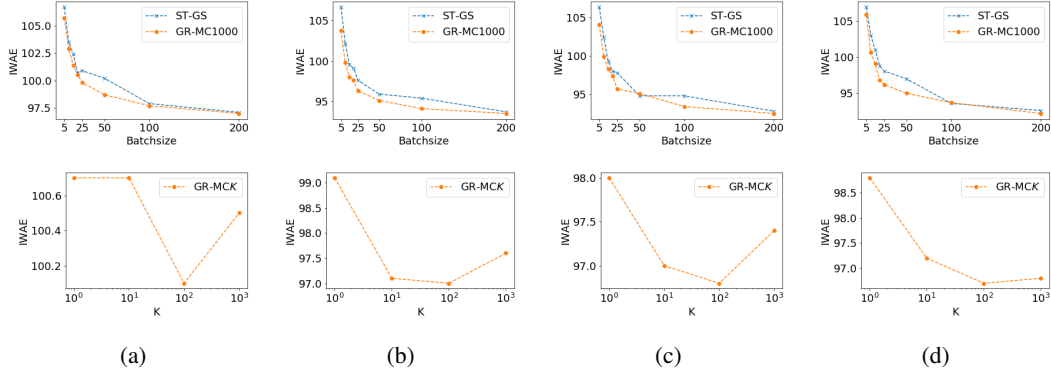


Figure 4: Increasing the number of Monte Carlo samples  $K$  to reduce variance in gradient estimation tends to improve performance. The performance difference tends to be larger at smaller batch sizes. Columns correspond to arities, i.e. (a) binary, (b) 4-ary, (c) 8-ary, (d) 16-ary. First row, IWAE on test set for best validated model trained at various batch sizes. Second row, IWAE on test set for best validated model trained at various  $K$  at batch size 20.

Table 3: Our estimator, GR-MCK, consistently achieves better performance across arities and batchsizes. The outperformance tends to be larger at smaller batchsizes. Best bound on the negative log-likelihood selected on the validation set from 20 independent runs at randomly searched hyperparameters.

	ESTIMATOR	BINARY		4-ARY		8-ARY		16-ARY	
		VALID.	TEST	VALID.	TEST	VALID.	TEST	VALID.	TEST
BATCH-SIZE 5	ST-GS	107.7	106.7	107.8	106.7	107.5	106.4	108.1	107.0
	GR-MC1000	<b>106.7</b>	<b>105.7</b>	<b>104.7</b>	<b>103.8</b>	<b>105.1</b>	<b>104.1</b>	<b>107.0</b>	<b>105.9</b>
BATCH-SIZE 10	ST-GS	104.4	103.5	103.2	102.2	103.5	102.4	104.1	103.1
	GR-MC1000	<b>103.7</b>	<b>102.9</b>	<b>100.8</b>	<b>99.8</b>	<b>100.9</b>	<b>99.9</b>	<b>101.8</b>	<b>100.7</b>
BATCH-SIZE 15	ST-GS	103.4	102.4	100.4	99.5	100.3	99.3	101.9	101.0
	GR-MC1000	<b>102.3</b>	<b>101.4</b>	<b>99.0</b>	<b>98.0</b>	<b>99.2</b>	<b>98.3</b>	<b>100.2</b>	<b>99.1</b>
BATCH-SIZE 20	ST-GS	101.5	100.7	100.0	99.1	99.0	98.0	99.8	98.8
	GR-MC1000	<b>101.3</b>	<b>100.5</b>	<b>98.4</b>	<b>97.6</b>	<b>97.5</b>	<b>96.5</b>	<b>97.8</b>	<b>96.8</b>
BATCH-SIZE 25	ST-GS	101.7	100.9	98.6	97.6	98.8	97.8	99.0	98.1
	GR-MC1000	<b>100.7</b>	<b>99.8</b>	<b>97.2</b>	<b>96.3</b>	<b>96.6</b>	<b>95.7</b>	<b>97.1</b>	<b>96.2</b>
BATCH-SIZE 50	ST-GS	101.2	100.2	96.7	95.9	<b>95.7</b>	<b>94.8</b>	98.0	97.0
	GR-MC1000	<b>99.5</b>	<b>98.7</b>	<b>96.0</b>	<b>95.1</b>	95.9	95.1	<b>95.9</b>	<b>95.0</b>
BATCH-SIZE 100	ST-GS	98.8	97.9	96.3	95.4	95.7	94.8	<b>94.4</b>	<b>93.6</b>
	GR-MC1000	<b>98.5</b>	<b>97.7</b>	<b>95.0</b>	<b>94.1</b>	<b>94.3</b>	<b>93.4</b>	94.6	93.7
BATCH-SIZE 200	ST-GS	97.9	97.1	94.5	93.7	93.6	92.8	93.4	92.6
	GR-MC1000	<b>97.8</b>	<b>97.0</b>	<b>94.3</b>	<b>93.5</b>	<b>93.2</b>	<b>92.5</b>	<b>93.1</b>	<b>92.2</b>