

## Are score function estimators an underestimated approach to learning with $k$ -subset sampling?

We revisit score function estimators for  $k$ -subset sampling and find results competitive with popular approaches based on pathwise gradient estimation and relaxed sampling, despite weaker assumptions. Furthermore, our method produces exact samples, unbiased gradients, and introduces no hyperparameters in need of tuning.

### Keywords

- Score function estimators
- $k$ -Subset sampling
- Variance reduction

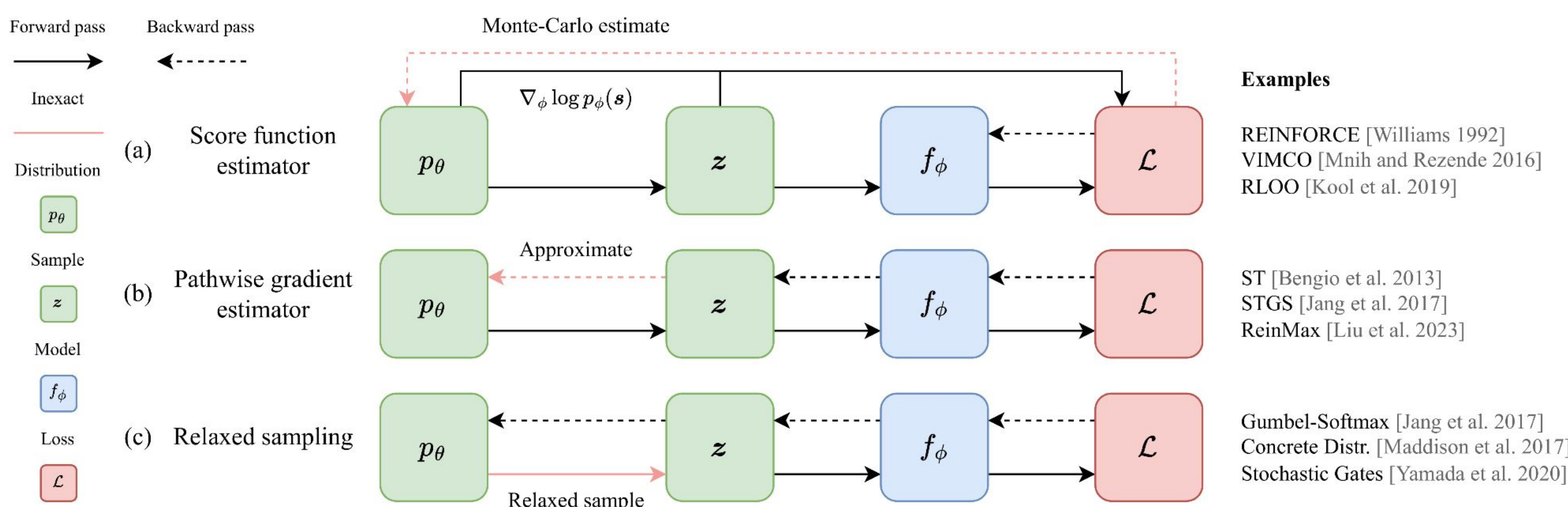


Figure 1: Three prominent approaches to learning by sampling.

### The $k$ -Subset Distribution

We model a subset using  $n$  parameters as the following conditional distribution:

$$p_{\theta,k}(z) = p_{\theta}(b \mid \sum_{i=1}^n b_i = k) = \frac{\prod_{i=1}^n p_{\theta}(b_i)}{p_{\theta}(\sum_{i=1}^n b_i = k)} \mathbb{1}[\sum_{i=1}^n b_i = k]$$

The distribution has a combinatorially large support: all subsets.

### Computing the Score Function

To construct our estimator, we need the score function, which involves computing the density of a Poisson binomial distribution. This can be efficiently done using a DFT:

$$p_{\theta}(\sum_{i=1}^n b_i = k) = \frac{1}{n+1} \text{DFT} \left( \prod_{i=1}^n p_{\theta}(b_i) e^{C + (1 - p_{\theta}(b_i))} \right)$$

With this, the score function is easily obtained using automatic differentiation.

### Reducing the Variance with Control Variates

The vanilla score function estimator suffers from high variance. Control variates of different sorts offer a simple remedy. We opt for a multi-sample control variate:

$$\begin{aligned} & \nabla_{\theta} \mathbb{E}_{p(x)} \mathbb{E}_{p_{\theta,k}(z)} [f_{\phi}(z, x)] \\ & \approx \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \nabla_{\theta} \log p_{\theta,k}(z^{(j)}) \\ & \cdot \left( f_{\phi}(z^{(j)}, x^{(i)}) - \frac{1}{M-1} \sum_{k \neq j} f_{\phi}(z^{(k)}, x^{(i)}) \right) \end{aligned}$$

## Results

We evaluate our method in a feature selection setting where a subset of features is sampled, passed through a downstream network, and optimized end-to-end. We use  $k = 30$  selections.

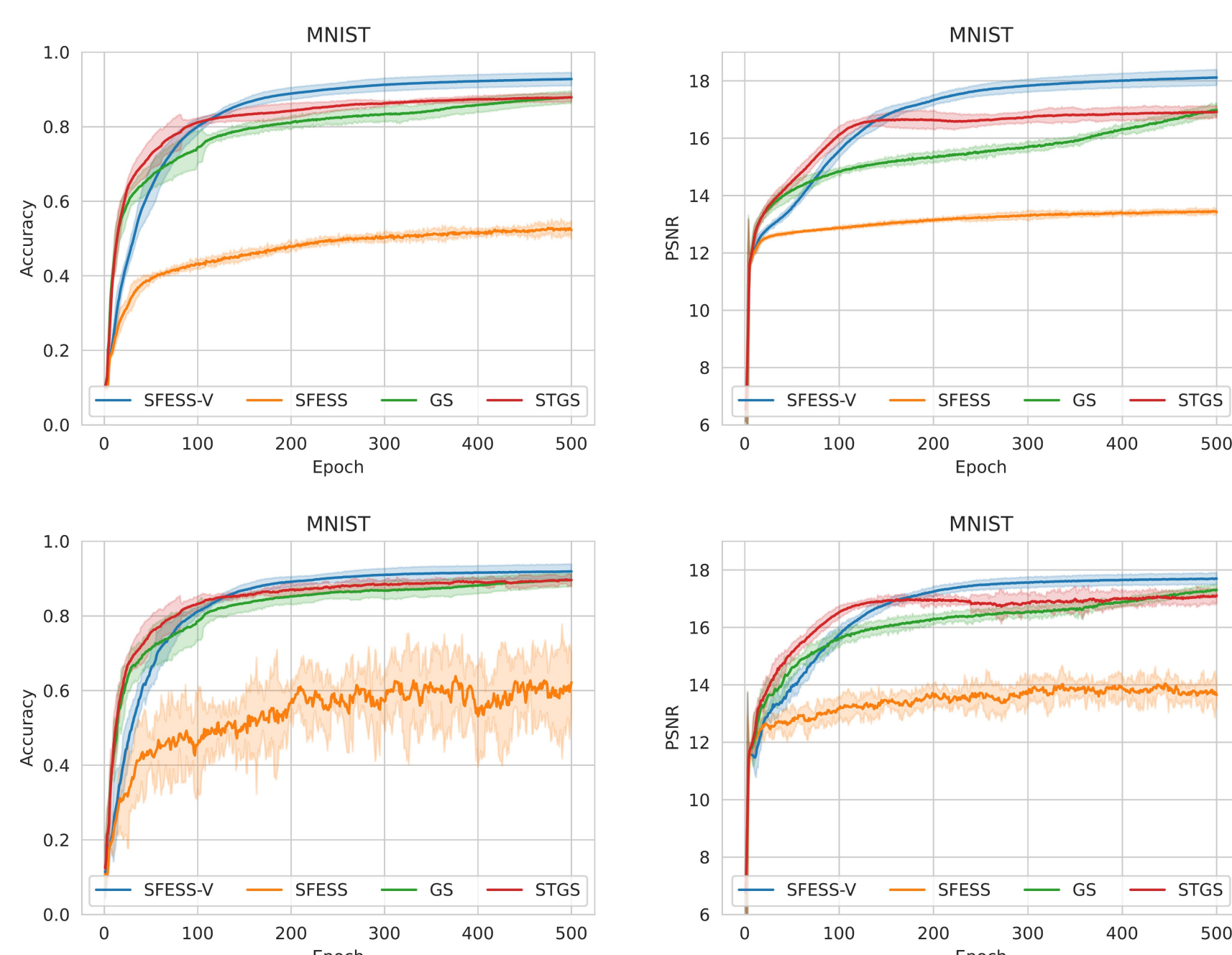
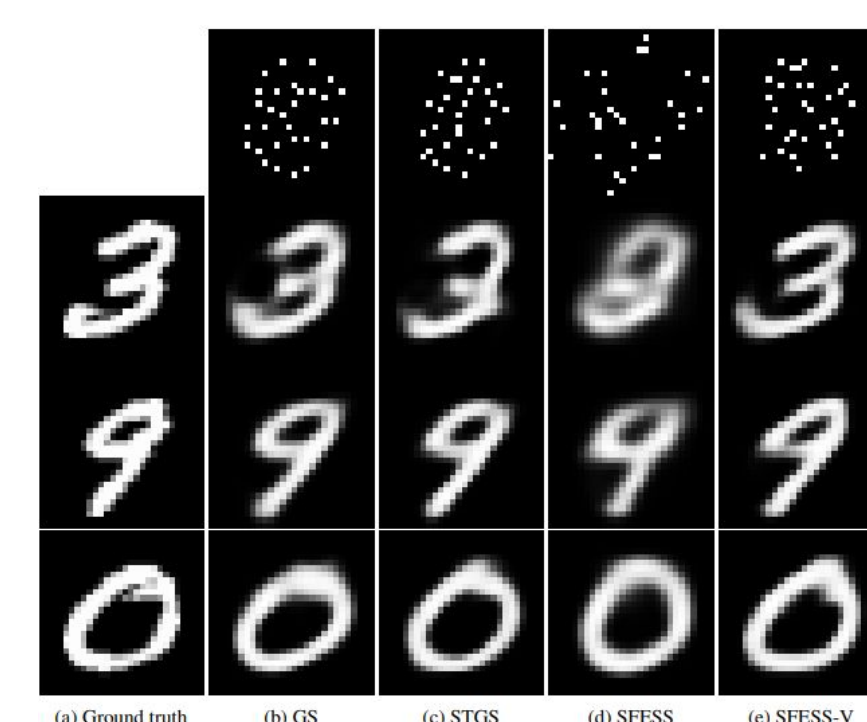


Figure 2: Training (top) and validation (bottom) metrics for classification (left) and reconstruction (right) on MNIST.

Table: Feature selection results.

Metric	Dataset	GS	STGS	SFESS	SFESS-V
PSNR $\uparrow$	MNIST	17.402 $\pm$ 0.194	17.186 $\pm$ 0.319	13.686 $\pm$ 0.412	<b>17.775 <math>\pm</math> 0.209</b>
	Fashion-MNIST	16.922 $\pm$ 0.289	16.642 $\pm$ 0.358	15.308 $\pm$ 0.303	<b>17.805 <math>\pm</math> 0.075</b>
	KMNIST	12.641 $\pm$ 0.083	12.561 $\pm$ 0.140	11.300 $\pm$ 0.281	<b>12.696 <math>\pm</math> 0.120</b>
SSIM $\uparrow$	MNIST	0.771 $\pm$ 0.011	0.759 $\pm$ 0.319	0.416 $\pm$ 0.412	<b>0.796 <math>\pm</math> 0.209</b>
	Fashion-MNIST	0.586 $\pm$ 0.017	0.578 $\pm$ 0.031	0.456 $\pm$ 0.016	<b>0.642 <math>\pm</math> 0.011</b>
	KMNIST	0.428 $\pm$ 0.021	<b>0.464 <math>\pm</math> 0.048</b>	0.230 $\pm$ 0.026	<b>0.460 <math>\pm</math> 0.022</b>
Accuracy $\uparrow$	MNIST	0.898 $\pm$ 0.014	0.898 $\pm$ 0.019	0.627 $\pm$ 0.077	<b>0.921 <math>\pm</math> 0.015</b>
	Fashion-MNIST	0.777 $\pm$ 0.012	0.774 $\pm$ 0.027	0.643 $\pm$ 0.112	<b>0.809 <math>\pm</math> 0.009</b>
	KMNIST	0.604 $\pm$ 0.029	0.591 $\pm$ 0.032	0.425 $\pm$ 0.023	<b>0.634 <math>\pm</math> 0.059</b>

Figure 3: Learned selections and reconstructions on the MNIST test set.



Paper

