

A Appendix

A.1 Additional results

In this section, we include more detailed results of experiments in Section 5 evaluating the intermediate- q robustness of models trained on different objectives, across different datasets and perturbation distributions. Specifically, the intermediate- q robustness metrics we consider, \hat{Z}_{MC} and \hat{Z}_{PS} , correspond to estimations of the functional q -norm of the cross entropy loss function evaluated over the perturbation distribution using the Monte Carlo estimator and the path sampling estimator respectively. We also include the standard cross entropy loss and adversarial loss metrics for comparison purposes. The following results were obtained from multiple training runs using 3 different random seeds.

MNIST ℓ_∞ -norm ball In Table 4 we report the mean and standard deviation of all intermediate- q robustness estimates over the ℓ_∞ -norm ball with radius $\epsilon = 0.3$ on MNIST for models trained according to the cross entropy loss (“Standard”), the Monte Carlo estimate (“MC”), the path sampling estimate computed with Hamiltonian Monte Carlo (“PS”), and PGD adversarial loss. For training, we used $m = 50$ to compute the MC estimate, and $m = 25$ and $L = 2$ to compute the path sampling estimate. The adversarially trained model was trained using PGD with 50 iterations. The evaluation estimate \hat{Z}_{MC} was computed with $m = 2000$ samples and \hat{Z}_{PS+HMC} was computed with $m = 100$ samples and $L = 20$ leapfrog steps. The evaluation adversarial loss was computed with PGD with 100 iterations. While both estimators interpolate between loss over random samples (i.e. $q = 1$), and adversarial loss (i.e. $q = \infty$), the path sampling estimator with HMC consistently results in higher (better) estimates of the desired integral for the same number of iterations. These results also show promise for training according to the path sampling metric, as the models trained with the PS+HMC estimate of the objective for larger values of q get increasingly better performance on intermediate- q robustness metrics (i.e. lower \hat{Z}_{PS+HMC}) as compared to models trained with the MC estimate. However, the adversarially trained model still does better in terms of intermediate- q performance, likely due to the nature of MNIST. We also include test robust accuracy of each trained model (from the PGD-100 evaluation) in Table 5 which shows that training using these estimates with increasingly large q does indeed improve worst-case robust performance. We note that we didn’t include standard accuracy here, because on MNIST, an adversarially trained model doesn’t lose much in terms of standard performance, however on more challenging datasets, we would expect to see a similar decrease in standard accuracy as we increase q , and thus one could choose a value of q based on the desired trade-off between standard and robust accuracy.

CIFAR-10 ℓ_∞ -norm ball In Table 6 we report the mean and standard deviation of all intermediate- q robustness estimates over the ℓ_∞ -norm ball with radius $\epsilon = 0.03$ on CIFAR-10 in Table 4 for models trained according to the cross entropy loss (“Standard”), the Monte Carlo estimate (“MC”), and PGD adversarial loss. For training, we used $m = 10$ to compute the MC estimate, and the adversarially trained model was trained using PGD with 10 iterations. The evaluation estimate \hat{Z}_{MC} was computed with $m = 500$ samples and \hat{Z}_{PS+HMC} was computed with $m = 50$ samples and $L = 10$ leapfrog steps. The evaluation adversarial loss was computed with PGD with 50 iterations and 10 restarts. The results show that, given enough iterations, the intermediate- q robustness estimates do interpolate between loss over random samples and loss over worst-case samples. However, given the more challenging nature of CIFAR-10 vs. MNIST, more samples are needed to get good estimates of the desired integral, making the problem more computationally intensive. While for models trained with the MC estimator, we do see an improvement in intermediate- q robust performance from training with $q = 10$ vs. $q = 1$, training using values of q larger than 10 does not provide much, if any, additional benefit. However, given the same number of iterations, we were not able to improve upon these results by training using the path sampling estimator, suggesting that the number of samples is just not large enough to get a good estimate of the objective for either estimator.

CIFAR-10 spatial transformations In Table 7 we report the mean and standard deviation of all intermediate- q robustness estimates over nondifferentiable spatial transformations on CIFAR-10 for models trained according to the cross entropy loss (“Standard”) and the Monte Carlo estimate (“MC”). For training, we used $m = 10$ to compute the MC estimate. The evaluation estimates \hat{Z}_{MC} and \hat{Z}_{PS} were computed with $m = 500$ samples, where in this case the path sampling estimate \hat{Z}_{PS} is based on

Table 4: Mean and standard deviation over 3 training runs with different random seeds for experiments on MNIST for the ℓ_∞ -norm ball with radius 0.3. The results show that the intermediate- q robustness metrics, \hat{Z}_{MC} and \hat{Z}_{PS+HMC} , interpolate between loss over random samples (when $q = 1$) and adversarial loss (when $q = \infty$), with the path sampling estimator consistently resulting in higher (better) estimates for larger values of q . Models trained using the PS+HMC estimate with larger q have better intermediate- q performance (lower \hat{Z}_{PS+HMC}) for corresponding values of q than those trained using the MC estimate.

Train method	Standard	\hat{Z}_{MC}				\hat{Z}_{PS+HMC}				Adv. loss
		$q = 1$	$q = 10$	$q = 10^2$	$q = 10^3$	$q = 1$	$q = 10$	$q = 10^2$	$q = 10^3$	
Standard	0.028 ± 0.001	0.043 ± 0.004	0.140 ± 0.023	0.251 ± 0.041	0.268 ± 0.045	0.043 ± 0.004	0.160 ± 0.028	1.420 ± 0.202	4.456 ± 0.495	11.649 ± 0.893
MC $q = 1$	0.034 ± 0.001	0.032 ± 0.000	0.084 ± 0.001	0.143 ± 0.002	0.154 ± 0.003	0.032 ± 0.000	0.088 ± 0.002	0.692 ± 0.026	2.133 ± 0.087	7.363 ± 0.435
MC $q = 10$	0.027 ± 0.001	0.026 ± 0.002	0.058 ± 0.002	0.098 ± 0.002	0.105 ± 0.002	0.026 ± 0.002	0.058 ± 0.001	0.412 ± 0.011	1.336 ± 0.050	3.722 ± 0.231
MC $q = 10^2$	0.026 ± 0.001	0.025 ± 0.001	0.055 ± 0.001	0.093 ± 0.001	0.099 ± 0.001	0.025 ± 0.001	0.055 ± 0.001	0.388 ± 0.015	1.261 ± 0.069	3.492 ± 0.291
MC $q = 10^3$	0.026 ± 0.002	0.025 ± 0.002	0.055 ± 0.001	0.093 ± 0.002	0.100 ± 0.001	0.025 ± 0.002	0.055 ± 0.001	0.390 ± 0.013	1.268 ± 0.059	3.488 ± 0.242
PS $q = 1$	0.037 ± 0.001	0.035 ± 0.001	0.094 ± 0.003	0.163 ± 0.005	0.174 ± 0.007	0.035 ± 0.001	0.101 ± 0.005	0.781 ± 0.027	2.334 ± 0.064	8.881 ± 0.616
PS $q = 10$	0.034 ± 0.001	0.031 ± 0.001	0.075 ± 0.001	0.126 ± 0.002	0.135 ± 0.003	0.031 ± 0.001	0.075 ± 0.001	0.467 ± 0.014	1.307 ± 0.029	5.012 ± 0.353
PS $q = 10^2$	0.030 ± 0.003	0.028 ± 0.002	0.060 ± 0.006	0.099 ± 0.010	0.107 ± 0.011	0.028 ± 0.002	0.058 ± 0.006	0.304 ± 0.007	0.816 ± 0.027	2.613 ± 0.193
PS $q = 10^3$	0.024 ± 0.001	0.024 ± 0.001	0.047 ± 0.001	0.077 ± 0.001	0.083 ± 0.001	0.024 ± 0.001	0.045 ± 0.001	0.239 ± 0.008	0.684 ± 0.037	1.646 ± 0.095
PGD-50	0.032 ± 0.005	0.039 ± 0.006	0.051 ± 0.007	0.076 ± 0.010	0.081 ± 0.011	0.039 ± 0.006	0.048 ± 0.007	0.101 ± 0.016	0.187 ± 0.040	0.270 ± 0.033

Table 5: Robust accuracy (PGD-100) of experiments on MNIST for perturbations in the ℓ_∞ ball of radius $\epsilon = 0.3$.

Training method	Robust accuracy
Standard	4.69%
MC $q = 1$	25.21%
MC $q = 10$	43.54%
MC $q = 10^2$	47.27%
MC $q = 10^3$	45.45%
PS $q = 1$	22.27%
PS $q = 10$	43.57%
PS $q = 10^2$	65.11%
PS $q = 10^3$	69.45%
PGD-50	91.55%

using Gaussian random walk Metropolis Hastings to sample from the unnormalized loss distribution rather than Hamiltonian Monte Carlo due to the non-differentiable perturbation distribution. The evaluation adversarial loss was approximated by averaging the maximum loss value encountered for each example during the Metropolis Hastings sampling process. As with the case of the ℓ_∞ -norm ball perturbation distribution on CIFAR-10, a larger number of samples are needed to get good estimates of the desired integral, making training using these estimates challenging, as shown by the lack of improvement in robust performance for larger values of q for models trained according to MC $q = 10^2$ vs. MC $q = 10$.

A.2 Additional figures

We include additional plots showing the convergence of Monte Carlo and path sampling estimates on a single test batch given an increasing number of samples. In Figure 2 we plot intermediate- q robustness estimates for $q = 10$ and $q = 100$ on a standard trained model over the ℓ_∞ -norm ball

Table 6: Mean and standard deviation over 3 training runs with different random seeds for experiments on CIFAR-10 for the ℓ_∞ -norm ball with radius 0.03. The results show that the intermediate- q robustness metrics, \hat{Z}_{MC} and \hat{Z}_{PS+HMC} , interpolate between loss over random samples (when $q = 1$) and adversarial loss (when $q = \infty$), with the path sampling estimator consistently resulting in higher (better) estimates for larger values of q . How to train for better intermediate- q robust performance on CIFAR-10, while still being computationally reasonable, remains an open question.

Train method	Standard	\hat{Z}_{MC}				\hat{Z}_{PS+HMC}				Adv. loss
		$q = 1$	$q = 10$	$q = 10^2$	$q = 10^3$	$q = 1$	$q = 10$	$q = 10^2$	$q = 10^3$	
Standard	0.382 ± 0.002	0.453 ± 0.006	0.787 ± 0.023	1.153 ± 0.037	1.216 ± 0.039	0.453 ± 0.006	0.841 ± 0.026	2.718 ± 0.090	4.991 ± 0.117	18.142 ± 0.448
MC $q = 1$	0.400 ± 0.005	0.405 ± 0.004	0.532 ± 0.006	0.717 ± 0.009	0.756 ± 0.010	0.405 ± 0.004	0.546 ± 0.006	1.490 ± 0.015	3.140 ± 0.017	14.240 ± 0.011
MC $q = 10$	0.393 ± 0.002	0.398 ± 0.003	0.468 ± 0.004	0.598 ± 0.005	0.630 ± 0.005	0.398 ± 0.003	0.471 ± 0.004	1.037 ± 0.013	2.365 ± 0.019	12.051 ± 0.036
MC $q = 10^2$	0.399 ± 0.003	0.402 ± 0.003	0.466 ± 0.003	0.589 ± 0.003	0.620 ± 0.004	0.402 ± 0.003	0.468 ± 0.003	0.980 ± 0.006	2.269 ± 0.020	12.084 ± 0.135
MC $q = 10^3$	0.399 ± 0.003	0.405 ± 0.003	0.469 ± 0.002	0.593 ± 0.003	0.625 ± 0.004	0.405 ± 0.003	0.471 ± 0.002	0.993 ± 0.005	2.302 ± 0.009	12.173 ± 0.128
PGD-10	0.731 ± 0.005	0.733 ± 0.005	0.734 ± 0.005	0.743 ± 0.005	0.761 ± 0.005	0.733 ± 0.005	0.734 ± 0.005	0.743 ± 0.005	0.796 ± 0.003	1.411 ± 0.009

Table 7: Mean and standard deviation of robustness estimates over 3 training runs with different random seeds for experiments on CIFAR-10 with non-differentiable parameterizations of flips, rotation, translation, and scaling. The results show that the intermediate- q robustness metrics, \hat{Z}_{MC} and \hat{Z}_{PS} , interpolate between loss over random samples (when $q = 1$) and adversarial loss (when $q = \infty$), with the path sampling estimator consistently resulting in higher (better) estimates for larger values of q . However, training for intermediate- q robust performance on CIFAR-10 on this perturbation distribution is more challenging due to the computational complexity.

Train method	Standard	\hat{Z}_{MC}				\hat{Z}_{PS}				Adv. loss
		$q = 1$	$q = 10$	$q = 10^2$	$q = 10^3$	$q = 1$	$q = 10$	$q = 10^2$	$q = 10^3$	
Standard	0.186 ± 0.007	0.450 ± 0.017	2.268 ± 0.067	3.687 ± 0.095	3.865 ± 0.114	0.444 ± 0.017	2.450 ± 0.068	4.636 ± 0.113	4.889 ± 0.111	5.625 ± 0.134
MC $q = 1$	0.154 ± 0.011	0.191 ± 0.002	0.800 ± 0.026	1.246 ± 0.024	1.300 ± 0.023	0.186 ± 0.004	0.879 ± 0.018	1.615 ± 0.039	1.711 ± 0.018	2.021 ± 0.036
MC $q = 10$	0.963 ± 0.014	1.019 ± 0.012	1.086 ± 0.010	1.348 ± 0.027	1.423 ± 0.029	1.015 ± 0.012	1.078 ± 0.010	1.416 ± 0.033	1.502 ± 0.078	1.596 ± 0.038
MC $q = 10^2$	2.014 ± 0.001	2.131 ± 0.002	2.131 ± 0.001	2.137 ± 0.002	2.164 ± 0.001	2.130 ± 0.002	2.131 ± 0.002	2.136 ± 0.002	2.171 ± 0.001	2.190 ± 0.006

perturbation distribution on CIFAR-10 ($\epsilon = 0.3$). Again, we see that HMC-based path sampling provides a much better estimate of the integral than random sampling, especially for higher q , for the same number of iterations (where iterations for the MC estimate is equal to m samples, and iterations for the PS+HMC estimate is equal to m samples times L leapfrog steps). Additionally, convergence plots for robustness estimates over the spatial transformations on CIFAR-10 (described in Section 5.2) are shown in Figure 3 specifically for $q = 1$ and $q = 100$ on a standard trained model over the spatial transformations. We see that path sampling and Monte Carlo sampling converge to nearly the same estimate for $q = 1$, but quickly diverge for $q = 100$, showing that path sampling works much better as an estimator even when the perturbation set is non-differentiable and we do not have the advantage of the Hamiltonian Monte Carlo sampler.

A.3 Experimental details

On both the ℓ_∞ experiments and the spatial transform CIFAR-10 experiments, we train a PreAct ResNet18 architecture using the SGD optimizer for 200 epochs with a starting learning rate of 0.1, Nesterov momentum of 0.9, and weight decay 0.0005. On MNIST we train using the Adam optimizer for 10 epochs with a starting learning rate of 0.001. We use a convolutional ReLU architecture with two convolutional layers with 32 and 64 channels and kernel sizes of 4×4 , which are followed by a

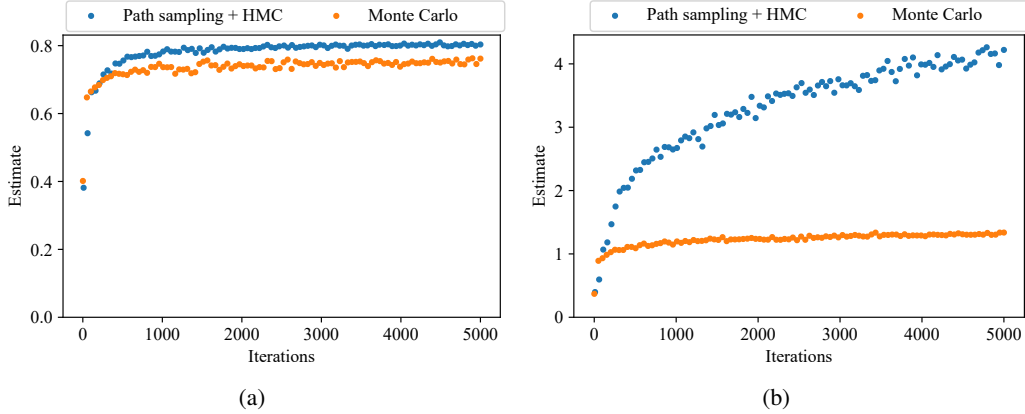


Figure 2: Comparison of the convergence of the path sampling with HMC estimate ($\hat{Z}_{\text{PS+HMC}}$) and Monte Carlo sampling estimate (\hat{Z}_{MC}) of the functional q -norm of the loss over the ℓ_∞ -norm ball perturbation distribution for (a) $q = 10$, and (b) $q = 1000$ with increasing iterations on a standard trained model on CIFAR-10. Iterations corresponds m for the Monte Carlo estimator, and to $m \times L$ for the path sampling + HMC estimator, where we fix $L = 10$. For a smaller $q = 10$ (a), path sampling with HMC is slightly better than Monte Carlo sampling. For a larger $q = 1000$ (b), using path sampling with HMC allows for a much better estimate of robustness.

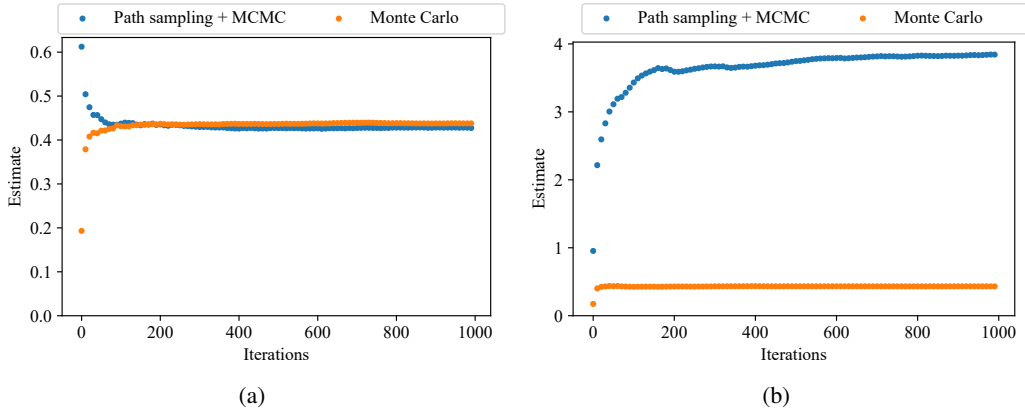


Figure 3: Comparison of the convergence of the path sampling (\hat{Z}_{PS}) and Monte Carlo sampling estimate (\hat{Z}_{MC}) of the functional q -norm of the loss over the spatial transforms for (a) $q = 1$, and (b) $q = 100$ with increasing iterations (corresponding to number of samples m) on a standard trained model on CIFAR-10. For a smaller $q = 1$ (a), Monte Carlo sampling and path sampling converge to the nearly the same estimate, whereas for a larger $q = 100$ (b), the estimates diverge quickly, with path sampling providing a much better estimate given the same number of samples.

fully connected layer with 1024 units. Each training and evaluation run is performed using a single Quadro RTX 8000 GPU.

A.3.1 ℓ_∞ -norm ball

We use a learning rate schedule that divides the starting learning rate by 10 halfway and two thirds of the way through training. On CIFAR-10 training, we train for 200 epochs and we also do *not* use random flip/crop data augmentation that is typically used for training CIFAR-10. For adversarial training and evaluation, we set the step size such that $\alpha = 2.5 \cdot \epsilon/m$, where m is the number of PGD steps, and we use early stopping for adversarial training based on the adversarial validation loss. For training and evaluation runs using HMC-based path sampling, we use $\sigma = 0.1$, and set the step size

such that $\alpha = \rho \cdot \sigma^2 / L$. For $q = 1$ and $q = 10$, we set $\rho = 0.6$, for $q = 100$ we set $\rho = 0.4$, and for $q = 1000$, we set $\rho = 0.2$.

A.3.2 Non-differentiable spatial transforms

We perform standard training for 50 epochs, and training using the Monte Carlo sampling estimator (for $q = 1$, $q = 10$, and $q = 100$) for 200 epochs. We use a learning rate schedule that linearly increases from 0 to the maximum value of 0.1 for the first two fifths of training epochs, and then linearly decreases to 0. When training using the MC estimator, we do not perform random flip/crop data augmentation.