# Efficient Baseline-free Sampling in Parameter Exploring Policy Gradients: Super Symmetric PGPE

Frank Sehnke

Zentrum für Sonnenenergie- und Wasserstoff-Forschung,
Industriestr. 6, Stuttgart, BW 70565 Germany

**Abstract.** Policy Gradient methods that explore directly in parameter space are among the most effective and robust direct policy search methods and have drawn a lot of attention lately. The basic method from this field, Policy Gradients with Parameter-based Exploration, uses two samples that are symmetric around the current hypothesis to circumvent misleading reward in *asymmetrical* reward distributed problems gathered with the usual baseline approach. The exploration parameters are still updated by a baseline approach - leaving the exploration prone to asymmetric reward distributions. In this paper we will show how the exploration parameters can be sampled quasi symmetric despite having limited instead of free parameters for exploration. We give a transformation approximation to get quasi symmetric samples with respect to the exploration without changing the overall sampling distribution. Finally we will demonstrate that sampling symmetrically also for the exploration parameters is superior in needs of samples and robustness than the original sampling approach.

## 1 Introduction

Policy Gradient (PG) methods that explore directly in parameter space have some major advantages over standard PG methods, like described in [1,2,3,4,5,6] and [7] and have therefore drawn a lot of attention in the last years. The basic method from the field of Parameter Exploring Policy Gradients (PEPG) [8], Policy Gradients with Parameter-based Exploration (PGPE) [1], uses two samples that are symmetric around the current hypothesis to circumvent misleading reward in *asymmetrical* reward distributed problems, gathered with the usual baseline approach. [4] shows that Symmetric Sampling (SyS) is superior even to the optimal baseline. The exploration parameters, however, are still updated by a baseline approach - leaving the exploration prone to asymmetric reward distributions. While the optimal baseline improved this issue substantially, like shown again by [4], it is likely that removing the baseline altogether by a SyS wrt. the exploration parameters will be again superior. Because the exploration parameters are standard deviations that are bounded between zero and infinity, there exist no correct symmetric samples wrt. the exploration parameters.

We will, however, show how the exploration parameters can be sampled quasi symmetric. We give therefore a transformation approximation to get quasi symmetric samples without changing the overall sampling distribution significantly, so that the PGPE assumptions based on normal distributed samples still hold. Finally we will demonstrate via experiments that sampling symmetrically also for the exploration parameters is superior in needs of samples and robustness compared to the original sampling approach, if confronted with search spaces with significant amounts of local optima.

## 2    Method

In this section we derive the super-symmetric sampling (SupSyS) method. We show how the method relates to SyS and sampling with a baseline, thereby summarizing the derivation from [1] for SyS and baseline sampling PGPE.

### 2.1    Parameter-Based Exploration

To stay conform with the nomenclature of [1] and [4], we assume a Markovian environment that produces a cumulative reward $r$ for a fixed length *episode*, *history*, *trajectory* or *roll-out*. In this setting, the goal of reinforcement learning is to find the optimal policy parameters $\boldsymbol{\theta}$ that maximize the agent's expected reward

$$J(\boldsymbol{\theta}) = \int_H p(h|\boldsymbol{\theta})r(h)dh. \tag{1}$$

An obvious way to maximize $J(\boldsymbol{\theta})$ is to estimate $\nabla_{\boldsymbol{\theta}} J$ and use it to carry out gradient ascent optimization. The probabilistic policy used in standard PG is replaced with a probability distribution over the parameters $\boldsymbol{\theta}$ for PGPE. The advantage of this approach is that the actions are deterministic, and an entire history can therefore be generated from a single parameter sample. This reduction in samples-per-history is what reduces the variance in the gradient estimate (see [1] for details).

We name the distribution over parameters in accordance with [1] $\boldsymbol{\rho}$. The expected reward with a given $\boldsymbol{\rho}$ is

$$J(\boldsymbol{\rho}) = \int_{\boldsymbol{\Theta}} \int_H p(h, \boldsymbol{\theta}|\boldsymbol{\rho})r(h)dhd\boldsymbol{\theta}. \tag{2}$$

Differentiating this form of the expected return with respect to $\boldsymbol{\rho}$ and applying sampling methods (first choosing $\boldsymbol{\theta}$ from $p(\boldsymbol{\theta}|\boldsymbol{\rho})$, then running the agent to generate $h$ from $p(h|\boldsymbol{\theta})$) yields the following gradient estimator:

$$\nabla_{\boldsymbol{\rho}} J(\boldsymbol{\rho}) \approx \frac{1}{N} \sum_{n=1}^{N} \nabla_{\boldsymbol{\rho}} \log p(\boldsymbol{\theta}|\boldsymbol{\rho})r(h^n). \tag{3}$$

Assuming that $\boldsymbol{\rho}$ consists of a set of means $\{\mu_i\}$ and standard deviations $\{\sigma_i\}$ that determine an independent normal distribution for each parameter $\theta_i$ in

$\boldsymbol{\theta}$ gives the following forms for the derivative of the characteristic eligibility $\log p(\boldsymbol{\theta}|\boldsymbol{\rho})$ with respect to $\mu_i$ and $\sigma_i$

$$\nabla_{\mu_i} \log p(\boldsymbol{\theta}|\boldsymbol{\rho}) = \frac{(\theta_i - \mu_i)}{\sigma_i^2}, \qquad \nabla_{\sigma_i} \log p(\boldsymbol{\theta}|\boldsymbol{\rho}) = \frac{(\theta_i - \mu_i)^2 - \sigma_i^2}{\sigma_i^3}, \qquad (4)$$

which can be substituted into Eq. (3) to approximate the $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ gradients.

## 2.2   Sampling with a Baseline

Given enough samples, Eq. (3) will determine the reward gradient to arbitrary accuracy. However each sample requires rolling out an entire state-action history, which is expensive. Following [9], we obtain a cheaper gradient estimate by drawing a single sample $\boldsymbol{\theta}$ and comparing its reward $r$ to a baseline reward $b$ given e.g. by a moving average over previous samples. Intuitively, if $r > b$ we adjust $\boldsymbol{\rho}$ so as to increase the probability of $\boldsymbol{\theta}$, and $r < b$ we do the opposite. If, as in [9], we use a step size $\alpha_i = \alpha\sigma_i^2$ in the direction of positive gradient (where $\alpha$ is a constant) we get the following parameter update equations:

$$\Delta\mu_i = \alpha(r - b)(\theta_i - \mu_i), \qquad \Delta\sigma_i = \alpha(r - b)\frac{(\theta_i - \mu_i)^2 - \sigma_i^2}{\sigma_i}. \qquad (5)$$

Usually the baseline is realized as decaying or moving average baseline of the form:

$$b(n) = \gamma r(h^{n-1}) + (1 - \gamma)b(n - 1) \qquad \text{or} \qquad b(n) = \sum_{n=N-m}^{N} r(h^n)/m \qquad (6)$$

[4] showed recently that an optimal baseline can be achieved for PGPE and the algorithm converges significantly faster with an optimal baseline of the form:

$$b^* = \frac{\mathbb{E}[r(h)||\nabla_{\boldsymbol{\rho}} \log p(\boldsymbol{\theta}|\boldsymbol{\rho})||^2]}{\mathbb{E}[||\nabla_{\boldsymbol{\rho}} \log p(\boldsymbol{\theta}|\boldsymbol{\rho})||^2]}. \qquad (7)$$

## 2.3   Symmetric Sampling

While sampling with a baseline is efficient and reasonably accurate for most scenarios, it has several drawbacks. In particular, if the reward distribution is strongly skewed then the comparison between the sample reward and the baseline reward is misleading. A more robust gradient approximation can be found by measuring the difference in reward between two symmetric samples on either side of the current mean. That is, we pick a perturbation $\boldsymbol{\epsilon}$ from the distribution $\mathcal{N}(0, \boldsymbol{\sigma})$, then create symmetric parameter samples $\boldsymbol{\theta}^+ = \boldsymbol{\mu} + \boldsymbol{\epsilon}$ and $\boldsymbol{\theta}^- = \boldsymbol{\mu} - \boldsymbol{\epsilon}$. Defining $r^+$ as the reward given by $\boldsymbol{\theta}^+$ and $r^-$ as the reward given by $\boldsymbol{\theta}^-$. We can insert the two samples into Eq. (3) and make use of Eq. (4) to obtain

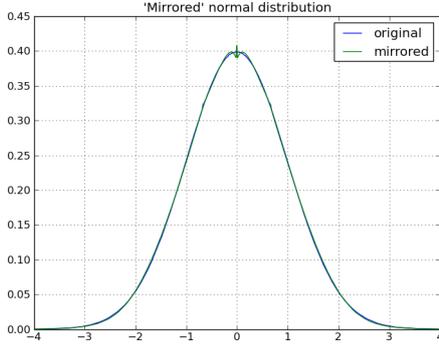$$\nabla_{\mu_i} J(\boldsymbol{\rho}) \approx \frac{\epsilon_i(r^+ - r^-)}{2\sigma_i^2}, \qquad (8)$$

**Fig. 1.** Normal distribution and the final approximation of the 'mirrored' distribution.
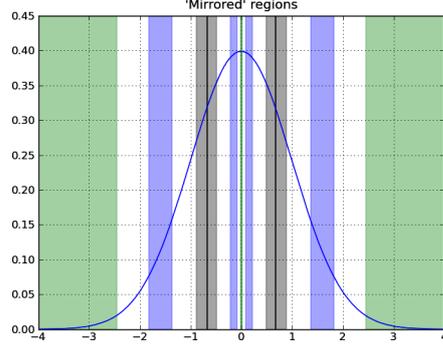
**Fig. 2.** Normal distribution and the regions that are transfered into each other by 'reflecting' the samples on the other side of the median deviation.

which resembles the *central difference* approximation used in finite difference methods. Using the same step sizes as before gives the following update equation for the $\boldsymbol{\mu}$ terms

$$\Delta\mu_i = \frac{\alpha\epsilon_i(r^+ - r^-)}{2}.$$  (9)

The updates for the standard deviations are more involved. As $\boldsymbol{\theta}^+$ and $\boldsymbol{\theta}^-$ are by construction equally probable under a given $\boldsymbol{\sigma}$, the difference between them cannot be used to estimate the $\boldsymbol{\sigma}$ gradient. Instead we take the mean $\frac{r^+ + r^-}{2}$ of the two rewards and compare it to the baseline reward $b$. This approach yields

$$\Delta\sigma_i = \alpha\left(\frac{r^+ + r^-}{2} - b\right)\left(\frac{\epsilon_i^2 - \sigma_i^2}{\sigma_i}\right)$$  (10)

SyS removes the problem of misleading baselines, and therefore improves the $\boldsymbol{\mu}$ gradient estimates. It also improves the $\boldsymbol{\sigma}$ gradient estimates, since both samples are equally probable under the current distribution, and therefore reinforce each other as predictors of the benefits of altering $\boldsymbol{\sigma}$. Even though symmetric sampling requires twice as many histories per update, [1] and [4] have shown that it gives a considerable improvement in convergence quality and time.

## 2.4   Super-Symmetric Sampling

While SyS removes the misleading baseline problem for the $\boldsymbol{\mu}$ gradient estimate, the $\boldsymbol{\sigma}$ gradient still uses a baseline and is prone to this problem. On the other hand there is no correct *symmetric* sample with respect to the standard deviation, because the standard deviation is bounded on the one *side* to 0 and is unbounded on the positive *side*. Another problem is that $\frac{2}{3}$ of the samples are on one *side* of the standard deviation and only $\frac{1}{3}$ on the other - *mirroring* the

samples to the opposite side of the standard deviation in some way, would therefore deform the normal distribution so much, that it would no longer be a close enough approximation to fulfill the assumptions that lead to the PGPE update rules.

We therefore chose to define the normal distribution via the mean and the median deviation $\phi$. The median deviation is due to the nice properties of the normal distribution simply defined by: $\phi = 0.67449 \cdot \sigma$. We can therefore draw samples from the new defined normal distribution: $\epsilon \sim \mathcal{N}_m(0, \phi)$.

The median deviation has by construction an equal amount of samples on either side and solves therefore the symmetry problem of *mirroring* samples. The update rule Eq. (9) stays unchanged while Eq. (10) is only scaled by $\frac{1}{0.67449}$ (the factor that transforms $\phi$ in $\sigma$) that can be substituted in $\alpha_\sigma$.

While the update rules stay the same for normal distributed sampling using the median deviation (despite a larger $\alpha_\sigma$), the median deviation is still also bounded on one side. Because the *mirroring* cannot be solved in closed form we resort to approximation via a polynomial that can be transfered to an infinite series. We found a good approximation for *mirroring* samples by:

$$a_i = \frac{\phi_i - \mid \epsilon_i \mid}{\phi_i}, \qquad \epsilon_i^* = sign(\epsilon_i) \cdot \phi_i \cdot \begin{cases} e^{c_1 \frac{|a_i|^3 - |a_i|}{\log(|a_i|)} + c_2|a_i|} & \text{if } a_i \leq 0 \\ e^{a_i}/(1. - a_i^3)^{c_3 a_i} & \text{if } a_i > 0, \end{cases} \tag{11}$$

with the following constants: $c_1 = -0.06655, c_2 = -0.9706, c_3 = 0.124$. This *mirrored* distribution has a standard deviation of 1.002 times the original standard deviation and looks like depicted in Fig. 1. Fig. 2 shows the regions of samples that are transfered into each other while generating the quasi symmetric samples.

Additional to the symmetric sample with respect to the mean hypothesis, now we also can generate two quasi symmetric samples with respect to the median deviation. We named this set of four samples super symmetric samples (SupSyS-samples). They allow for completely baseline free update rules, not only for the $\mu$ update but also for the $\sigma$ updates.

Therefore the two symmetric sample pairs are used to update $\mu$ according to Eq. (9). $\sigma$ is updated in a similar way by using the mean reward of each symmetric sample pair, there $r^{++}$ is the mean reward of the original symmetric sample pair and $r^{--}$ is the mean reward of the *mirrored* sample pair. The SupSyS update rule for the $\sigma$ update is given by:

$$\Delta\sigma_i = \frac{\alpha \frac{\epsilon_i^2 - \sigma_i^2}{\sigma_i}(r^{++} - r^{--})}{2}. \tag{12}$$

## 3   Experiments and Results

We use the square function as search space instance with no local optima and the Rastrigin function (see Fig. 8) as search space with exponentially many local optima, to test the different behavior of SupSyS- and SyS-PGPE. The two meta-parameters connected with SyS-PGPE as well as with SupSyS-PGPE, namely
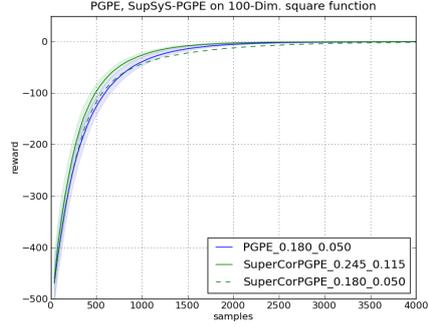
**Fig. 3.** Convergence plots of PGPE and SupSyS-PGPE on the 100 dimensional square function. The mean and standard deviation of 200 independent runs are shown.
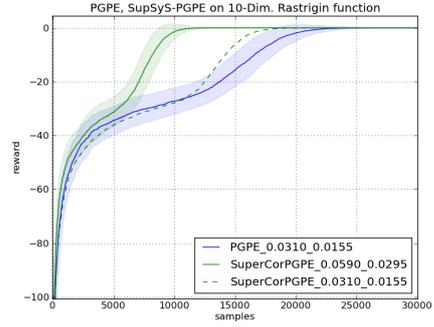
**Fig. 4.** Convergence plots of PGPE and SupSyS-PGPE on the 10 dimensional Rastrigin function. The mean and standard deviation of 200 independent runs are shown.

the step sizes for the $\mu$ and $\sigma$ updates, were optimized for every experiment via a grid search. The Figures 3 to 6 show the means and standard deviations of 200 independent runs each. It can be seen in Fig. 3 that for a search space with no local optima SupSyS-PGPE shows no advantage over standard SyS-PGPE. However, despite using 4 samples per update the performance is also not reduced by using SupSyS-PGPE — the two methods become merely equivalent. The situation changes drastically if the Rastrigin function is used as test function. Not only needs SupSyS-PGPE about half the samples compared to PGPE, the effect seems also to become stronger the higher dimensional the search space
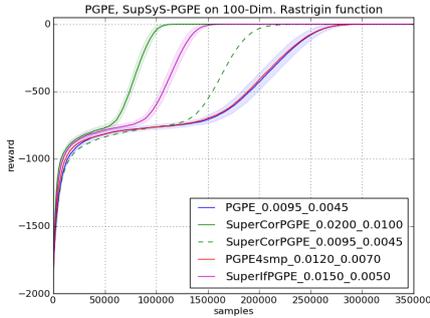


**Fig. 5.** Convergence plots of PGPE, PGPE with 4 samples (PGPE4smp), conditional SupSyS-PGPE (SupIf-PGPE) and SupSyS-PGPE on the 100 dimensional Rastrigin function. The mean and standard deviation of 200 independent runs are shown.
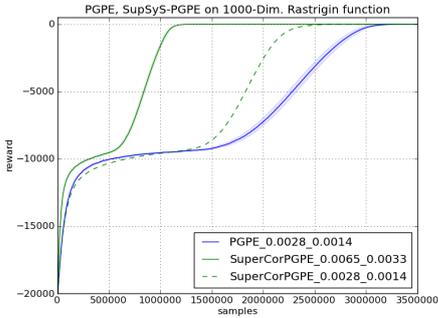
**Fig. 6.** Convergence plots of PGPE and SupSyS-PGPE on the 1000 dimensional Rastrigin function. The mean and standard deviation of 200 independent runs are shown.
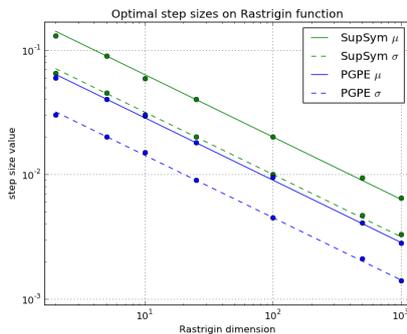
**Fig. 7.** Optimal meta-parameters for the multi-dimensional Rastrigin function for PGPE and SupSyS-PGPE.
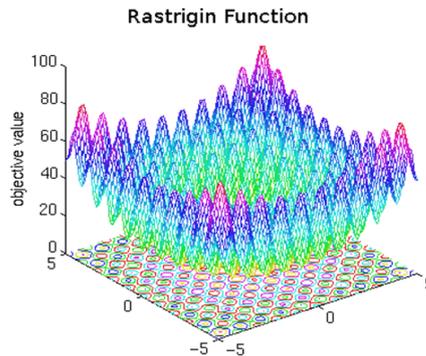
**Fig. 8.** Visualization of the 2D Rastrigin function.

gets (see Fig. 4 to Fig. 6). We also added SupSyS-PGPE plots with the for SyS-PGPE optimal (less greedy) meta parameters to show that the effect is not only due to the more aggressive meta parameters. This runs were also more efficient than for PGPE, the effect was however not so distinct.

In Fig. 5 we also show a standard PGPE experiment with 4 samples (2 SyS samples — *PGPE4smp*) instead of 2 to show that the improved performance is not due to the different samples per update. Fig. 5 additionally shows an experiment (SupIf-PGPE) there symmetric samples are only drawn if the first sample(s) result in worse reward than a decaying average baseline. The intuitive idea behind symmetric samples was initially that changing the parameters *away* from the current sample if the sample resulted in lower than average reward may move the mean hypothesis still in a worse region of the parameter space. Search spaces like the one given in the Rastrigin function can visualize this problem. For SupIf-PGPE one Sample is drawn. If the reward is larger than the baseline then an update is done immediately. If not, a symmetric sample is drawn. Is the mean reward connected with both samples better than the baseline an SyS-PGPE update is done. If also this mean reward is worse than the baseline, a full SupSyS-PGPE update with 2 additional SyS samples is performed. As can be seen in Fig. 5 the performance is worse by some degree — the difference is however small enough that maybe the optimal baseline approach would improve this method enough to be challenging to SupSyS-PGPE (see also Sec. 4).

The optimal meta-parameters are an exponential function of the search space dimension, like to expect, so that we observe a line in the *loglog*-plot of Fig. 7. For SupSyS-PGPE the meta-parameters are about 2 times larger than for SyS-PGPE. This is partly because SupSyS-PGPE uses four samples per update instead of two. But the optimal meta-parameters are also larger than for the PGPE4smp experiment so that the symmetric nature of the four SupSyS samples obviously brings additional stability in the gradient estimate than a pure averaging over 4 samples would.

## 4   Conclusions and Future Work

We introduced SupSyS-PGPE, a completely baseline free PGPE that uses quasi-symmetric samples wrt. the exploration parameters. We showed that on the Rastrigin function, as example for a test function with exponentially many local optima, this novel method is clearly superior to standard SyS-PGPE and that both methods become equivalent in performance if the search space lack *distracting* local optima.

For future work we want to highlight that SupSyS-PGPE can be easily combined with other extensions of PGPE. Multi-modal PGPE [10] can be equipped straight forward with SupSyS sampling. Also the natural gradient used for PGPE in [3] can be defined over the SupSyS gradient instead over the vanilla gradient. If the full 4 super symmetric sample set is only used if the first samples are worse than a baseline (like described as SupIf-PGPE in Sec. 3) a combination with the optimal baseline (described for PGPE in [4]) can yield a superior method to both SupSyS-PGPE and optimal baseline PGPE. Also importance mixing introduced for PGPE by [5] is applicable to SupSyS-PGPE.

Finally a big open point for future work is the validation of the mere theoretical findings on real world problems, e.g. robotic tasks, for SupSyS-PGPE and its combination with other PGPE extensions.

## References

1. Sehnke, F., Osendorfer, C., Rückstieß, T., Graves, A., Peters, J., Schmidhuber, J.: Parameter-exploring policy gradients. Neural Networks **23**(4) (2010) 551–559
2. Rückstieß, T., Sehnke, F., Schaul, T., Wierstra, D., Sun, Y., Schmidhuber, J.: Exploring parameter space in reinforcement learning. Paladyn. Journal of Behavioral Robotics **1**(1) (2010) 14–24
3. Miyamae, A., Nagata, Y., Ono, I.: Natural Policy Gradient Methods with Parameter-based Exploration for Control Tasks. In: NIPS. (2010) 1–9
4. Zhao, T., Hachiya, H., Niu, G., Sugiyama, M.: Analysis and improvement of policy gradient estimation. Neural networks : the official journal of the International Neural Network Society (October 2011) 1–30
5. Zhao, T., Hachiya, H., Tangkaratt, V., Morimoto, J., Sugiyama, M.: Efficient sample reuse in policy gradients with parameter-based exploration. arXiv preprint arXiv:1301.3966 (2013)
6. Stulp, F., Sigaud, O.: Path integral policy improvement with covariance matrix adaptation. arXiv preprint arXiv:1206.4621 (2012)
7. Wierstra, D., Schaul, T., Peters, J., Schmidhuber, J.: Natural evolution strategies. In: Evolutionary Computation, 2008. CEC 2008., IEEE (2008) 3381–3387
8. Sehnke, F.: Parameter exploring policy gradients and their implications
9. Williams, R.J.: Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning **8** (1992) 229–256
10. Sehnke, F., Graves, A., Osendorfer, C., Schmidhuber, J.: Multimodal parameter-exploring policy gradients. In: Machine Learning and Applications (ICMLA), 2010 Ninth International Conference on, IEEE (2010) 113–118