

# Week 4 Briefing

Making robots learn by trial and error

Jack Cashman

July 19, 2024

# Plan for presentation

- 1 Overview and Motivation
- 2 Attempting to Improve Beta Policy Performance

# Overview and Motivation

- 1 Last week I found that the beta distribution is a promising alternative as a prior for the policy in continuous and bounded action spaces.
- 2 This week, I've been experimenting in order to boost the performance of PPO-clip when using a beta policy.
- 3 Today I'll recap some of these experiments, as well as some issues and how I fixed them.

# Overview and Motivation

- 1 Last week I found that the beta distribution is a promising alternative as a prior for the policy in continuous and bounded action spaces.
- 2 This week, I've been experimenting in order to boost the performance of PPO-clip when using a beta policy.
- 3 Today I'll recap some of these experiments, as well as some issues and how I fixed them.

# Overview and Motivation

- ① Last week I found that the beta distribution is a promising alternative as a prior for the policy in continuous and bounded action spaces.
- ② This week, I've been experimenting in order to boost the performance of PPO-clip when using a beta policy.
- ③ Today I'll recap some of these experiments, as well as some issues and how I fixed them.

# A Note On Numerical Instability

Given  $d$  independent, beta-distributed random variables  $X_1, \dots, X_d$ , each with parameters  $\alpha_i, \beta_i$ , the joint distribution,  $f$ , is:

$$f(x_1, \dots, x_d) \propto \prod_{i=1}^d x_i^{\alpha_i-1} (1 - x_i)^{\beta_i-1}$$

If either  $\alpha_i < 1$  or  $\beta_i < 1$ , for some  $1 \leq i \leq d$ , then  $f$  is unbounded. This is problematic, as in PPO I need to call a method that returns  $\log(f(\mathbf{x}))$  for a point  $\mathbf{x} \in [0, 1]^d$ . The distribution is bounded if we enforce  $\alpha_i, \beta_i \geq 1$  for all  $i$ . I've used this restriction in all of the following implementations.

# Restrictions on $\alpha$ and $\beta$

- For numerical stability I require  $\alpha_i, \beta_i > 1$ .
- In computation, the  $\alpha_i$  and  $\beta_i$  that is output by the NN are relatively close to this lower bound.
- This distribution is *very* stochastic. For larger  $\alpha_i$  and  $\beta_i$ , the policy becomes less stochastic.
- Perhaps introducing larger lower bounds on the  $\alpha_i$  and  $\beta_i$  can improve performance?
- I've tested lower bounds of 1, 3, 5, 25 on 5 seeds for 3 environments.

# Restrictions on $\alpha$ and $\beta$

- For numerical stability I require  $\alpha_i, \beta_i > 1$ .
- In computation, the  $\alpha_i$  and  $\beta_i$  that is output by the NN are relatively close to this lower bound.
- This distribution is *very* stochastic. For larger  $\alpha_i$  and  $\beta_i$ , the policy becomes less stochastic.
- Perhaps introducing larger lower bounds on the  $\alpha_i$  and  $\beta_i$  can improve performance?
- I've tested lower bounds of 1, 3, 5, 25 on 5 seeds for 3 environments.



# Restrictions on $\alpha$ and $\beta$

- For numerical stability I require  $\alpha_i, \beta_i > 1$ .
- In computation, the  $\alpha_i$  and  $\beta_i$  that is output by the NN are relatively close to this lower bound.
- This distribution is *very* stochastic. For larger  $\alpha_i$  and  $\beta_i$ , the policy becomes less stochastic.
- Perhaps introducing larger lower bounds on the  $\alpha_i$  and  $\beta_i$  can improve performance?
- I've tested lower bounds of 1, 3, 5, 25 on 5 seeds for 3 environments.

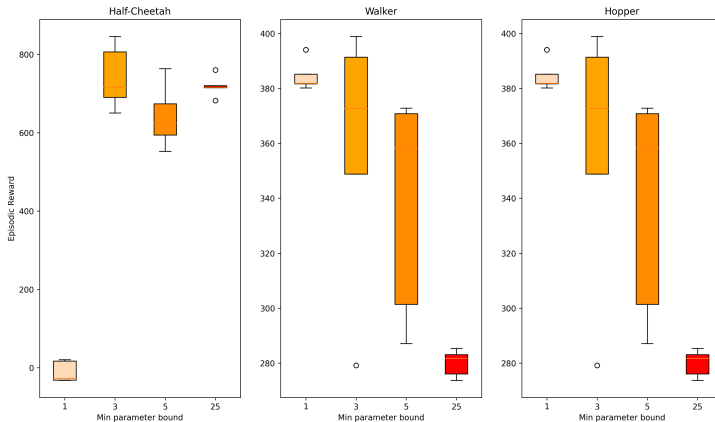
# Restrictions on $\alpha$ and $\beta$

- For numerical stability I require  $\alpha_i, \beta_i > 1$ .
- In computation, the  $\alpha_i$  and  $\beta_i$  that is output by the NN are relatively close to this lower bound.
- This distribution is *very* stochastic. For larger  $\alpha_i$  and  $\beta_i$ , the policy becomes less stochastic.
- Perhaps introducing larger lower bounds on the  $\alpha_i$  and  $\beta_i$  can improve performance?
- I've tested lower bounds of 1, 3, 5, 25 on 5 seeds for 3 environments.

# Restrictions on $\alpha$ and $\beta$

- For numerical stability I require  $\alpha_i, \beta_i > 1$ .
- In computation, the  $\alpha_i$  and  $\beta_i$  that is output by the NN are relatively close to this lower bound.
- This distribution is *very* stochastic. For larger  $\alpha_i$  and  $\beta_i$ , the policy becomes less stochastic.
- Perhaps introducing larger lower bounds on the  $\alpha_i$  and  $\beta_i$  can improve performance?
- I've tested lower bounds of 1, 3, 5, 25 on 5 seeds for 3 environments.

# Restrictions on $\alpha$ and $\beta$ - Results



It appears that requiring  $\alpha_i, \beta_i \geq 3$  yields the best performance.

# Different NN models

- I've been using 2 NNs to predict the vectors  
 $\alpha = (\alpha_1 \dots \alpha_d)^\top$  and  $\beta = (\beta_1 \dots \beta_d)^\top$
- To ensure that both vectors satisfy  $\alpha_i, \beta_i > 3$ , I push the output through the ReLU function, and then add 3
- Nan suggested that we may be able to instead use a single neural network. So, output  $(\alpha_1 \dots \alpha_d \beta_1 \dots \beta_d)^\top$
- To test this, I've compared over 5 seeds for 3 separate mujoco environments; Half-Cheetah, Hopper, and Walker-2d

# Different NN models

- I've been using 2 NNs to predict the vectors  $\alpha = (\alpha_1 \dots \alpha_d)^\top$  and  $\beta = (\beta_1 \dots \beta_d)^\top$
- To ensure that both vectors satisfy  $\alpha_i, \beta_i > 3$ , I push the output through the ReLU function, and then add 3
- Nan suggested that we may be able to instead use a single neural network. So, output  $(\alpha_1 \dots \alpha_d \beta_1 \dots \beta_d)^\top$
- To test this, I've compared over 5 seeds for 3 separate mujoco environments; Half-Cheetah, Hopper, and Walker-2d

# Different NN models

- I've been using 2 NNs to predict the vectors  $\alpha = (\alpha_1 \dots \alpha_d)^\top$  and  $\beta = (\beta_1 \dots \beta_d)^\top$
- To ensure that both vectors satisfy  $\alpha_i, \beta_i > 3$ , I push the output through the ReLU function, and then add 3
- Nan suggested that we may be able to instead use a single neural network. So, output  $(\alpha_1 \dots \alpha_d \beta_1 \dots \beta_d)^\top$
- To test this, I've compared over 5 seeds for 3 separate mujoco environments; Half-Cheetah, Hopper, and Walker-2d

# Different NN models

- I've been using 2 NNs to predict the vectors  $\alpha = (\alpha_1 \dots \alpha_d)^\top$  and  $\beta = (\beta_1 \dots \beta_d)^\top$
- To ensure that both vectors satisfy  $\alpha_i, \beta_i > 3$ , I push the output through the ReLU function, and then add 3
- Nan suggested that we may be able to instead use a single neural network. So, output  $(\alpha_1 \dots \alpha_d \beta_1 \dots \beta_d)^\top$
- To test this, I've compared over 5 seeds for 3 separate mujoco environments; Half-Cheetah, Hopper, and Walker-2d



# Different NN models - Results

Overall, using a single network to predict both the vector of  $\alpha_i$ 's and  $\beta_i$ 's appears to be more effective.

	Cheetah	Hopper	Walker
Double Net	284 $\pm$ 46.6	358.2 $\pm$ 43.2	347 $\pm$ 20.3
Single Net	631 $\pm$ 48.4	383.4 $\pm$ 15.6	346.4 $\pm$ 10.8

Combining this with the aforementioned restriction that  $\alpha_i, \beta_i > 3$ , the algorithm performs similar to, or better than the `cleanrl` implementation that assume a diagonal Gaussian prior for the policy.

# Conclusion

The performance when a Beta prior is assumed vs the performance when a diagonal Gaussian prior is assumed across 5 random seeds:

	Cheetah	Hopper	Walker
Gaussian Prior	$32.5 \pm 139.4$	$765.6 \pm 237.5$	$365.7 \pm 13.6$
Beta Prior	$742 \pm 72.8$	$415.5 \pm 59.6$	$384.6 \pm 5$

The beta distribution serves as a promising alternative to the traditional diagonal Gaussian in certain continuous and bounded action spaces.