

---

**meta***policysearchDocumentation*

**Dennis Lee, Ignasi Clavera, Jonas Rothfuss**

**Sep 06, 2019**



---

## Contents:

---

<b>1</b>	<b>Meta-Policy Search</b>	<b>3</b>
1.1	Baselines	3
1.1.1	Baseline (Interface)	3
1.1.2	Linear Feature Baseline	3
1.1.3	LinearTimeBaseline	4
1.2	Environments	5
1.2.1	MetaEnv (Interface)	5
1.3	Meta-Algorithms	6
1.3.1	MAML-Algorithm (Interface)	6
1.3.2	ProMP-Algorithm	7
1.3.3	TRPO-MAML-Algorithm	8
1.3.4	VPD-MAML-Algorithm	8
1.4	Optimizers	10
1.4.1	Conjugate Gradient Optimizer	10
1.4.2	MAML First Order Optimizer	12
1.5	Policies	13
1.5.1	Policy Interfaces	13
1.5.2	Gaussian-Policies	17
1.6	Samplers	22
1.6.1	Sampler	22
1.6.2	Sample Processor	23
1.6.3	Vectorized Environment Executor	25
1.7	Meta-Trainer	27
<b>2</b>	<b>Indices and tables</b>	<b>29</b>
	<b>Python Module Index</b>	<b>31</b>
	<b>Index</b>	<b>33</b>



Despite recent progress, deep reinforcement learning (RL) still relies heavily on hand-crafted features and reward functions as well as engineered problem specific inductive bias. Meta-RL aims to forego such reliance by acquiring inductive bias in a data-driven manner. A particular instance of meta learning that has proven successful in RL is gradient-based meta-learning.

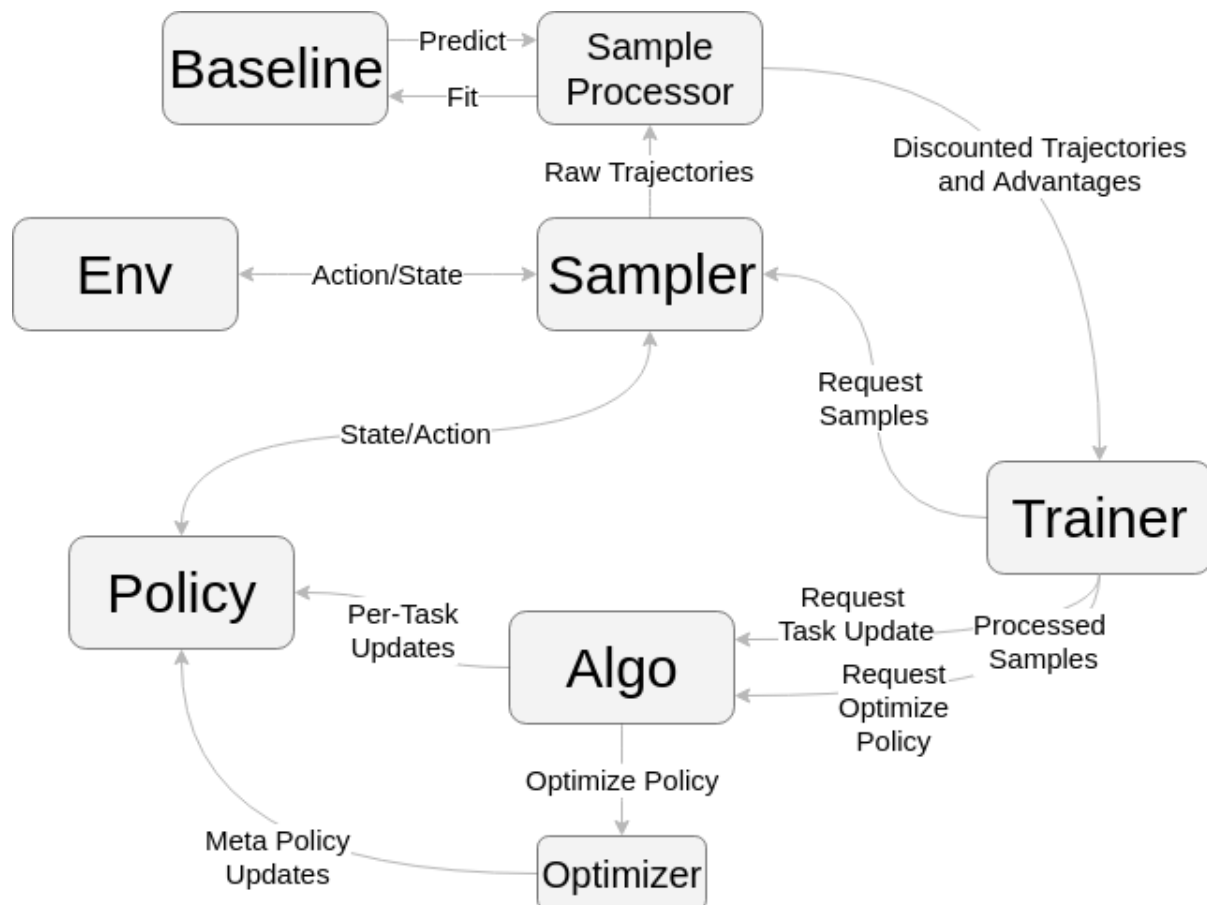
The code repository provides implementations of various gradient-based Meta-RL methods including

- ProMP: Proximal Meta-Policy Search (Rothfuss et al., 2018)
- MAML: Model Agnostic Meta-Learning (Finn et al., 2017)
- E-MAML: Exploration MAML (Al-Shedivat et al., 2018, Stadie et al., 2018)

The code was written as part of ProMP. Further information and experimental results can be found on our [website](#). This documentation specifies the API and interaction of the algorithm's components. Overall, on iteration of gradient-based Meta-RL consists of the followings steps:

1. Sample trajectories with pre update policy
2. Perform gradient step for each task to obtain updated/adapted policy
3. Sample trajectories with the updated/adapted policy
4. Perform a meta-policy optimization step, changing the pre-updates policy parameters

This high level structure of the algorithm is implemented in the Meta-Trainer class. The overall structure and interaction of the code components is depicted in the following figure:





## 1.1 Baselines

### 1.1.1 Baseline (Interface)

**class** `meta_policy_search.baselines.Baseline`

Reward baseline interface

**fit** (*paths*)

Fits the baseline model with the provided paths

**Parameters** *paths* – list of paths

**get\_param\_values** ()

Returns the parameter values of the baseline object

**log\_diagnostics** (*paths*, *prefix*)

Log extra information per iteration based on the collected paths

**predict** (*path*)

Predicts the reward baselines for a provided trajectory / path

**Parameters** *path* – dict of lists/numpy array containing trajectory / path information such as “observations”, “rewards”, ...

Returns: numpy array of the same length as `paths[“observations”]` specifying the reward baseline

**set\_params** (*value*)

Sets the parameter values of the baseline object

**Parameters** *value* – parameter value to be set

### 1.1.2 Linear Feature Baseline

**class** `meta_policy_search.baselines.LinearFeatureBaseline` (*reg\_coeff=1e-05*)

Linear (polynomial) time-state dependent return baseline model (see. Duan et al. 2016, “Benchmarking Deep

Reinforcement Learning for Continuous Control”, ICML)

Fits the following linear model

$$\text{reward} = b_0 + b_1 * \text{obs} + b_2 * \text{obs}^2 + b_3 * t + b_4 * t^2 + b_5 * t^3$$

**Parameters** `reg_coeff` – list of paths

**fit** (*paths*, *target\_key*=‘returns’)

Fits the linear baseline model with the provided paths via damped least squares

**Parameters**

- **paths** (*list*) – list of paths
- **target\_key** (*str*) – path dictionary key of the target that shall be fitted (e.g. “returns”)

**get\_param\_values** (*\*\*tags*)

Returns the parameter values of the baseline object

**Returns** numpy array of linear\_regression coefficients

**log\_diagnostics** (*paths*, *prefix*)

Log extra information per iteration based on the collected paths

**predict** (*path*)

Abstract Class for the LinearFeatureBaseline and the LinearTimeBaseline Predicts the linear reward baselines estimates for a provided trajectory / path. If the baseline is not fitted - returns zero baseline

**Parameters** **path** (*dict*) – dict of lists/numpy array containing trajectory / path information such as “observations”, “rewards”, ...

**Returns** numpy array of the same length as paths[“observations”] specifying the reward baseline

**Return type** (np.ndarray)

**set\_params** (*value*, *\*\*tags*)

Sets the parameter values of the baseline object

**Parameters** **value** – numpy array of linear\_regression coefficients

### 1.1.3 LinearTimeBaseline

**class** meta\_policy\_search.baselines.LinearTimeBaseline (*reg\_coeff*=1e-05)

Linear (polynomial) time-dependent reward baseline model

Fits the following linear model

$$\text{reward} = b_0 + b_3 * t + b_4 * t^2 + b_5 * t^3$$

**Parameters** `reg_coeff` – list of paths

**fit** (*paths*, *target\_key*=‘returns’)

Fits the linear baseline model with the provided paths via damped least squares

**Parameters**

- **paths** (*list*) – list of paths
- **target\_key** (*str*) – path dictionary key of the target that shall be fitted (e.g. “returns”)

**get\_param\_values** (*\*\*tags*)

Returns the parameter values of the baseline object

**Returns** numpy array of linear\_regression coefficients



**log\_diagnostics** (*paths, prefix*)

Log extra information per iteration based on the collected paths

**predict** (*path*)

Abstract Class for the LinearFeatureBaseline and the LinearTimeBaseline Predicts the linear reward baseline estimates for a provided trajectory / path. If the baseline is not fitted - returns zero baseline

**Parameters** **path** (*dict*) – dict of lists/numpy array containing trajectory / path information such as “observations”, “rewards”, ...

**Returns** numpy array of the same length as paths[“observations”] specifying the reward baseline

**Return type** (np.ndarray)

**set\_params** (*value, \*\*tags*)

Sets the parameter values of the baseline object

**Parameters** **value** – numpy array of linear\_regression coefficients

## 1.2 Environments

### 1.2.1 MetaEnv (Interface)

**class** meta\_policy\_search.envs.base.**MetaEnv** (*\*args, \*\*kwargs*)

Wrapper around OpenAI gym environments, interface for meta learning

**get\_task** ()

Gets the task that the agent is performing in the current environment

**Returns** task of the meta-learning environment

**Return type** task

**log\_diagnostics** (*paths, prefix*)

Logs env-specific diagnostic information

**Parameters**

- **paths** (*list*) – list of all paths collected with this env during this iteration
- **prefix** (*str*) – prefix for logger

**sample\_tasks** (*n\_tasks*)

Samples task of the meta-environment

**Parameters** **n\_tasks** (*int*) – number of different meta-tasks needed

**Returns** an (n\_tasks) length list of tasks

**Return type** tasks (list)

**set\_task** (*task*)

Sets the specified task to the current environment

**Parameters** **task** – task of the meta-learning environment

## 1.3 Meta-Algorithms

### 1.3.1 MAML-Algorithm (Interface)

```
class meta_policy_search.meta_algos.MAMLAlgo (policy, inner_lr=0.1, meta_batch_size=20,
                                             num_inner_grad_steps=1, trainable_inner_step_size=False)
```

Bases: meta\_policy\_search.meta\_algos.base.MetaAlgo

Provides some implementations shared between all MAML algorithms

#### Parameters

- **policy** (*Policy*) – policy object
- **inner\_lr** (*float*) – gradient step size used for inner step
- **meta\_batch\_size** (*int*) – number of meta-learning tasks
- **num\_inner\_grad\_steps** (*int*) – number of gradient updates taken per maml iteration
- **trainable\_inner\_step\_size** (*boolean*) – whether make the inner step size a trainable variable

#### **build\_graph** ()

Creates meta-learning computation graph

Pseudocode:

```
for task in meta_batch_size:
    make_vars
    init_dist_info_sym
for step in num_grad_steps:
    for task in meta_batch_size:
        make_vars
        update_dist_info_sym
set objectives for optimizer
```

#### **make\_vars** (prefix=“")

**Parameters** **prefix** (*str*) – a string to prepend to the name of each variable

**Returns** a tuple containing lists of placeholders for each input type and meta task

**Return type** (tuple)

#### **optimize\_policy** (all\_samples\_data, log=True)

Performs MAML outer step for each task

#### Parameters

- **all\_samples\_data** (*list*) – list of lists of lists of samples (each is a dict) split by gradient update and meta task
- **log** (*bool*) – whether to log statistics

**Returns** None

### 1.3.2 ProMP-Algorithm

```
class meta_policy_search.meta_algos.ProMP (*args,          name='ppo_maml',      learn-
                                         ing_rate=0.001,        num_ppo_steps=5,
                                         num_minibatches=1,    clip_eps=0.2,
                                         target_inner_step=0.01,
                                         init_inner_kl_penalty=0.01,      adap-
                                         tive_inner_kl_penalty=True,      an-
                                         neal_factor=1.0, **kwargs)
```

Bases: meta\_policy\_search.meta\_algos.base.MAMLAlgo

ProMP Algorithm

#### Parameters

- **policy** (*Policy*) – policy object
- **name** (*str*) – tf variable scope
- **learning\_rate** (*float*) – learning rate for optimizing the meta-objective
- **num\_ppo\_steps** (*int*) – number of ProMP steps (without re-sampling)
- **num\_minibatches** (*int*) – number of minibatches for computing the ppo gradient steps
- **clip\_eps** (*float*) – PPO clip range
- **target\_inner\_step** (*float*) – target inner kl divergence, used only when adaptive\_inner\_kl\_penalty is true
- **init\_inner\_kl\_penalty** (*float*) – initial penalty for inner kl
- **adaptive\_inner\_kl\_penalty** (*bool*) – whether to used a fixed or adaptive kl penalty on inner gradient update
- **anneal\_factor** (*float*) – multiplicative factor for annealing clip\_eps. If anneal\_factor < 1, clip\_eps <- anneal\_factor \* clip\_eps at each iteration
- **inner\_lr** (*float*) – gradient step size used for inner step
- **meta\_batch\_size** (*int*) – number of meta-learning tasks
- **num\_inner\_grad\_steps** (*int*) – number of gradient updates taken per maml iteration
- **trainable\_inner\_step\_size** (*boolean*) – whether make the inner step size a trainable variable

**build\_graph()**

Creates the computation graph

**make\_vars** (*prefix=""*)

**Parameters** **prefix** (*str*) – a string to prepend to the name of each variable

**Returns** a tuple containing lists of placeholders for each input type and meta task

**Return type** (tuple)

**optimize\_policy** (*all\_samples\_data, log=True*)

Performs MAML outer step

#### Parameters

- **all\_samples\_data** (*list*) – list of lists of lists of samples (each is a dict) split by gradient update and meta task
- **log** (*bool*) – whether to log statistics

**Returns** None

### 1.3.3 TRPO-MAML-Algorithm

```
class meta_policy_search.meta_algos.TRPOMAML (*args, name='trpo_maml', step_size=0.01,
                                             inner_type='likelihood_ratio', exploration=
                                             False, **kwargs)
```

Bases: meta\_policy\_search.meta\_algos.base.MAMLAlgo

Algorithm for TRPO MAML

#### Parameters

- **policy** (*Policy*) – policy object
- **name** (*str*) – tf variable scope
- **step\_size** (*int*) – trust region size for the meta policy optimization through TPRO
- **inner\_type** (*str*) – One of 'log\_likelihood', 'likelihood\_ratio', 'dice', choose which inner update to use
- **exploration** (*bool*) – whether to use E-MAML or MAML
- **inner\_lr** (*float*) – gradient step size used for inner step
- **meta\_batch\_size** (*int*) – number of meta-learning tasks
- **num\_inner\_grad\_steps** (*int*) – number of gradient updates taken per maml iteration
- **trainable\_inner\_step\_size** (*boolean*) – whether make the inner step size a trainable variable

**build\_graph** ()

Creates the computation graph

**make\_vars** (*prefix=""*)

**Parameters** **prefix** (*str*) – a string to prepend to the name of each variable

**Returns** a tuple containing lists of placeholders for each input type and meta task

**Return type** (tuple)

**optimize\_policy** (*all\_samples\_data*, *log=True*)

Performs MAML outer step

#### Parameters

- **all\_samples\_data** (*list*) – list of lists of lists of samples (each is a dict) split by gradient update and meta task
- **log** (*bool*) – whether to log statistics

**Returns** None

### 1.3.4 VPG-MAML-Algorithm

```
class meta_policy_search.meta_algos.VPGMAML (*args, name='vpg_maml',
                                             learning_rate=0.001, inner_type='likelihood_ratio',
                                             exploration=False, **kwargs)
```

Bases: meta\_policy\_search.meta\_algos.base.MAMLAlgo

Algorithm for PPO MAML

### Parameters

- **policy** (*Policy*) – policy object
- **name** (*str*) – tf variable scope
- **learning\_rate** (*float*) – learning rate for the meta-objective
- **exploration** (*bool*) – use exploration / pre-update sampling term / E-MAML term
- **inner\_type** (*str*) – inner optimization objective - either `log_likelihood` or `likelihood_ratio`
- **inner\_lr** (*float*) – gradient step size used for inner step
- **meta\_batch\_size** (*int*) – number of meta-learning tasks
- **num\_inner\_grad\_steps** (*int*) – number of gradient updates taken per maml iteration
- **trainable\_inner\_step\_size** (*boolean*) – whether make the inner step size a trainable variable

**build\_graph** ()

Creates the computation graph

**make\_vars** (*prefix=""*)

**Parameters** **prefix** (*str*) – a string to prepend to the name of each variable

**Returns** a tuple containing lists of placeholders for each input type and meta task

**Return type** (tuple)

**optimize\_policy** (*all\_samples\_data, log=True*)

Performs MAML outer step

### Parameters

- **all\_samples\_data** (*list*) – list of lists of lists of samples (each is a dict) split by gradient update and meta task
- **log** (*bool*) – whether to log statistics

**Returns** None

## 1.4 Optimizers

### 1.4.1 Conjugate Gradient Optimizer

```
class meta_policy_search.optimizers.ConjugateGradientOptimizer (cg_iters=10,
                                                             reg_coeff=0,
                                                             subsam-
                                                             ple_factor=1.0,
                                                             back-
                                                             track_ratio=0.8,
                                                             max_backtracks=15,
                                                             de-
                                                             bug_nan=False,
                                                             ac-
                                                             cept_violation=False,
                                                             hvp_approach=<meta_policy_search.op-
                                                             object>)
```

Bases: meta\_policy\_search.optimizers.base.Optimizer

Performs constrained optimization via line search. The search direction is computed using a conjugate gradient algorithm, which gives  $x = A^{-1}g$ , where  $A$  is a second order approximation of the constraint and  $g$  is the gradient of the loss function.

#### Parameters

- **cg\_iters** (*int*) – The number of conjugate gradients iterations used to calculate  $A^{-1}g$
- **reg\_coeff** (*float*) – A small value so that  $A \rightarrow A + \text{reg} * I$
- **subsample\_factor** (*float*) – Subsampling factor to reduce samples when using “conjugate gradient. Since the computation time for the descent direction dominates, this can greatly reduce the overall computation time.
- **backtrack\_ratio** (*float*) – ratio for decreasing the step size for the line search
- **max\_backtracks** (*int*) – maximum number of backtracking iterations for the line search
- **debug\_nan** (*bool*) – if set to True, NanGuard will be added to the compilation, and ipdb will be invoked when nan is detected
- **accept\_violation** (*bool*) – whether to accept the descent step if it violates the line search condition after exhausting all backtracking budgets
- **hvp\_approach** (*obj*) – Hessian vector product approach

**build\_graph** (*loss, target, input\_ph\_dict, leq\_constraint*)

Sets the objective function and target weights for the optimize function

#### Parameters

- **loss** (*tf\_op*) – minimization objective
- **target** (*Policy*) – Policy whose values we are optimizing over
- **inputs** (*list*) – tuple of tf.placeholders for input data which may be subsampled. The first dimension corresponds to the number of data points
- **extra\_inputs** (*list*) – tuple of tf.placeholders for hyperparameters (e.g. learning rate, if annealed)

- **leq\_constraint** (*tuple*) – A constraint provided as a tuple (f, epsilon), of the form  $f(*inputs) \leq \epsilon$ .

**constraint\_val** (*input\_val\_dict*)

Computes the value of the KL-divergence between pre-update policies for given inputs

**Parameters**

- **inputs** (*list*) – inputs needed to compute the inner KL
- **extra\_inputs** (*list*) – additional inputs needed to compute the inner KL

**Returns** value of the loss

**Return type** (float)

**gradient** (*input\_val\_dict*)

Computes the gradient of the loss function

**Parameters**

- **inputs** (*list*) – inputs needed to compute the gradient
- **extra\_inputs** (*list*) – additional inputs needed to compute the loss function

**Returns** flattened gradient

**Return type** (np.ndarray)

**loss** (*input\_val\_dict*)

Computes the value of the loss for given inputs

**Parameters**

- **inputs** (*list*) – inputs needed to compute the loss function
- **extra\_inputs** (*list*) – additional inputs needed to compute the loss function

**Returns** value of the loss

**Return type** (float)

**optimize** (*input\_val\_dict*)

Carries out the optimization step

**Parameters**

- **inputs** (*list*) – inputs for the optimization
- **extra\_inputs** (*list*) – extra inputs for the optimization
- **subsample\_grouped\_inputs** (*None or list*) – subsample data from each element of the list

## 1.4.2 MAML First Order Optimizer

```
class meta_policy_search.optimizers.MAMLFIRSTORDEROPTIMIZER (tf_optimizer_cls=<class
    'tensorflow.python.training.adam.AdamOptimizer'>,
    tf_optimizer_args=None,
    learning_rate=0.001,
    max_epochs=1,
    tolerance=1e-06,
    num_minibatches=1,
    verbose=False)
```

Bases: meta\_policy\_search.optimizers.base.Optimizer

Optimizer for first order methods (SGD, Adam)

### Parameters

- **tf\_optimizer\_cls** (*tf.train.optimizer*) – desired tensorflow optimizer for training
- **tf\_optimizer\_args** (*dict or None*) – arguments for the optimizer
- **learning\_rate** (*float*) – learning rate
- **max\_epochs** – number of maximum epochs for training
- **tolerance** (*float*) – tolerance for early stopping. If the loss function decreases less than the specified tolerance
- **an epoch, then the training stops.** (*after*) –
- **num\_minibatches** (*int*) – number of mini-batches for performing the gradient step. The mini-batch size is
- **size//num\_minibatches.** (*batch*) –
- **verbose** (*bool*) – Whether to log or not the optimization process

**build\_graph** (*loss, target, input\_ph\_dict*)

Sets the objective function and target weights for the optimize function

### Parameters

- **loss** (*tf.Op*) – minimization objective
- **target** (*Policy*) – Policy whose values we are optimizing over
- **input\_ph\_dict** (*dict*) – dict containing the placeholders of the computation graph corresponding to loss

**loss** (*input\_val\_dict*)

Computes the value of the loss for given inputs

**Parameters** **input\_val\_dict** (*dict*) – dict containing the values to be fed into the computation graph

**Returns** value of the loss

**Return type** (float)

**optimize** (*input\_val\_dict*)

Carries out the optimization step



**Parameters** `input_val_dict` (*dict*) – dict containing the values to be fed into the computation graph

**Returns** (float) loss before optimization

## 1.5 Policies

### 1.5.1 Policy Interfaces

```
class meta_policy_search.policies.Policy(obs_dim, action_dim, name='policy',
                                         hidden_sizes=(32, 32), learn_std=True,
                                         hidden_nonlinearity=<function tanh>,
                                         output_nonlinearity=None, **kwargs)
```

Bases: meta\_policy\_search.utils.serializable.Serializable

A container for storing the current pre and post update policies Also provides functions for executing and updating policy parameters

---

**Note:** the preupdate policy is stored as tf.Variables, while the postupdate policy is stored in numpy arrays and executed through tf.placeholders

---

#### Parameters

- **obs\_dim** (*int*) – dimensionality of the observation space -> specifies the input size of the policy
- **action\_dim** (*int*) – dimensionality of the action space -> specifies the output size of the policy
- **name** (*str*) – Name used for scoping variables in policy
- **hidden\_sizes** (*tuple*) – size of hidden layers of network
- **learn\_std** (*bool*) – whether to learn variance of network output
- **hidden\_nonlinearity** (*Operation*) – nonlinearity used between hidden layers of network
- **output\_nonlinearity** (*Operation*) – nonlinearity used after the final layer of network

**build\_graph** ()

Builds computational graph for policy

**distribution**

Returns this policy's distribution

**Returns** this policy's distribution

**Return type** (Distribution)

**distribution\_info\_keys** (*obs, state\_infos*)

#### Parameters

- **obs** (*placeholder*) – symbolic variable for observations
- **state\_infos** (*dict*) – a dictionary of placeholders that contains information about the

- of the policy at the time it received the observation (*state*)

–

**Returns** a dictionary of tf placeholders for the policy output distribution

**Return type** (dict)

**distribution\_info\_sym** (*obs\_var*, *params=None*)

Return the symbolic distribution information about the actions.

**Parameters**

- **obs\_var** (*placeholder*) – symbolic variable for observations
- **params** (*None or dict*) – a dictionary of placeholders that contains information about the
- of the policy at the time it received the observation (*state*)

–

**Returns** a dictionary of tf placeholders for the policy output distribution

**Return type** (dict)

**get\_action** (*observation*)

Runs a single observation through the specified policy

**Parameters** **observation** (*array*) – single observation

**Returns** array of arrays of actions for each env

**Return type** (array)

**get\_actions** (*observations*)

Runs each set of observations through each task specific policy

**Parameters** **observations** (*array*) – array of arrays of observations generated by each task and env

**Returns**

**array of arrays of actions for each env** (*meta\_batch\_size*) x (*batch\_size*) x (*action\_dim*)  
and array of arrays of agent\_info dicts

**Return type** (tuple)

**get\_param\_values** ()

Gets a list of all the current weights in the network (in original code it is flattened, why?)

**Returns** list of values for parameters

**Return type** (list)

**get\_params** ()

Get the tf.Variables representing the trainable weights of the network (symbolic)

**Returns** a dict of all trainable Variables

**Return type** (dict)

**likelihood\_ratio\_sym** (*obs*, *action*, *dist\_info\_old*, *policy\_params*)

Computes the likelihood p\_new(obs|act)/p\_old ratio between

**Parameters**

- **obs** (*tf.Tensor*) – symbolic variable for observations
- **action** (*tf.Tensor*) – symbolic variable for actions

- **dist\_info\_old** (*dict*) – dictionary of tf.placeholders with old policy information
- **policy\_params** (*dict*) – dictionary of the policy parameters (each value is a tf.Tensor)

**Returns** likelihood ratio

**Return type** (tf.Tensor)

**log\_diagnostics** (*paths*)

Log extra information per iteration based on the collected paths

**log\_likelihood\_sym** (*obs, action, policy\_params*)

Computes the log likelihood p(obs|act)

**Parameters**

- **obs** (*tf.Tensor*) – symbolic variable for observations
- **action** (*tf.Tensor*) – symbolic variable for actions
- **policy\_params** (*dict*) – dictionary of the policy parameters (each value is a tf.Tensor)

**Returns** log likelihood

**Return type** (tf.Tensor)

**set\_params** (*policy\_params*)

Sets the parameters for the graph

**Parameters** **policy\_params** (*dict*) – of variable names and corresponding parameter values

**class** meta\_policy\_search.policies.**MetaPolicy** (*\*args, \*\*kwargs*)

Bases: meta\_policy\_search.policies.base.Policy

**build\_graph** ()

Also should create lists of variables and corresponding assign ops

**distribution**

Returns this policy's distribution

**Returns** this policy's distribution

**Return type** (Distribution)

**distribution\_info\_keys** (*obs, state\_infos*)

**Parameters**

- **obs** (*placeholder*) – symbolic variable for observations
- **state\_infos** (*dict*) – a dictionary of placeholders that contains information about the
- **of the policy at the time it received the observation** (*state*) –

**Returns** a dictionary of tf placeholders for the policy output distribution

**Return type** (dict)

**distribution\_info\_sym** (*obs\_var, params=None*)

Return the symbolic distribution information about the actions.

**Parameters**

- **obs\_var** (*placeholder*) – symbolic variable for observations
- **params** (*None or dict*) – a dictionary of placeholders that contains information about the
- **of the policy at the time it received the observation** (*state*)

**Returns** a dictionary of tf placeholders for the policy output distribution

**Return type** (dict)

**get\_action** (*observation*)

Runs a single observation through the specified policy

**Parameters** **observation** (*array*) – single observation

**Returns** array of arrays of actions for each env

**Return type** (array)

**get\_actions** (*observations*)

Runs each set of observations through each task specific policy

**Parameters** **observations** (*array*) – array of arrays of observations generated by each task and env

**Returns**

**array of arrays of actions for each env** (*meta\_batch\_size*) x (*batch\_size*) x (*action\_dim*)  
and array of arrays of agent\_info dicts

**Return type** (tuple)

**get\_param\_values** ()

Gets a list of all the current weights in the network (in original code it is flattened, why?)

**Returns** list of values for parameters

**Return type** (list)

**get\_params** ()

Get the tf.Variables representing the trainable weights of the network (symbolic)

**Returns** a dict of all trainable Variables

**Return type** (dict)

**likelihood\_ratio\_sym** (*obs, action, dist\_info\_old, policy\_params*)

Computes the likelihood  $p_{\text{new}}(\text{obs}|\text{act})/p_{\text{old}}$  ratio between

**Parameters**

- **obs** (*tf.Tensor*) – symbolic variable for observations
- **action** (*tf.Tensor*) – symbolic variable for actions
- **dist\_info\_old** (*dict*) – dictionary of tf.placeholders with old policy information
- **policy\_params** (*dict*) – dictionary of the policy parameters (each value is a tf.Tensor)

**Returns** likelihood ratio

**Return type** (tf.Tensor)

**log\_diagnostics** (*paths*)

Log extra information per iteration based on the collected paths

**log\_likelihood\_sym** (*obs, action, policy\_params*)

Computes the log likelihood  $p(\text{obs}|\text{act})$

**Parameters**

- **obs** (*tf.Tensor*) – symbolic variable for observations
- **action** (*tf.Tensor*) – symbolic variable for actions
- **policy\_params** (*dict*) – dictionary of the policy parameters (each value is a *tf.Tensor*)

**Returns** log likelihood

**Return type** (*tf.Tensor*)

**policies\_params\_feed\_dict**

returns fully prepared feed dict for feeding the currently saved policy parameter values into the lightweight policy graph

**set\_params** (*policy\_params*)

Sets the parameters for the graph

**Parameters** **policy\_params** (*dict*) – of variable names and corresponding parameter values

**switch\_to\_pre\_update** ()

Switches `get_action` to pre-update policy

**update\_task\_parameters** (*updated\_policies\_parameters*)

**Parameters**

- **updated\_policies\_parameters** (*list*) – List of size meta-batch size. Each contains a dict with the policies
- **as numpy arrays** (*parameters*) –

## 1.5.2 Gaussian-Policies

**class** `meta_policy_search.policies.GaussianMLPPolicy` (*\*args, init\_std=1.0, min\_std=1e-06, \*\*kwargs*)

Bases: `meta_policy_search.policies.base.Policy`

Gaussian multi-layer perceptron policy (diagonal covariance matrix) Provides functions for executing and updating policy parameters A container for storing the current pre and post update policies

**Parameters**

- **obs\_dim** (*int*) – dimensionality of the observation space -> specifies the input size of the policy
- **action\_dim** (*int*) – dimensionality of the action space -> specifies the output size of the policy
- **name** (*str*) – name of the policy used as `tf` variable scope
- **hidden\_sizes** (*tuple*) – tuple of integers specifying the hidden layer sizes of the MLP
- **hidden\_nonlinearity** (*tf.op*) – nonlinearity function of the hidden layers
- **output\_nonlinearity** (*tf.op or None*) – nonlinearity function of the output layer

- **learn\_std** (*boolean*) – whether the standard\_dev / variance is a trainable or fixed variable
- **init\_std** (*float*) – initial policy standard deviation
- **min\_std** (*float*) – minimal policy standard deviation

**build\_graph** ()

Builds computational graph for policy

**distribution**

Returns this policy's distribution

**Returns** this policy's distribution

**Return type** (Distribution)

**distribution\_info\_keys** (*obs, state\_infos*)**Parameters**

- **obs** (*placeholder*) – symbolic variable for observations
- **state\_infos** (*dict*) – a dictionary of placeholders that contains information about the
- **of the policy at the time it received the observation** (*state*) –

**Returns** a dictionary of tf placeholders for the policy output distribution

**Return type** (dict)

**distribution\_info\_sym** (*obs\_var, params=None*)

Return the symbolic distribution information about the actions.

**Parameters**

- **obs\_var** (*placeholder*) – symbolic variable for observations
- **params** (*dict*) – a dictionary of placeholders or vars with the parameters of the MLP

**Returns** a dictionary of tf placeholders for the policy output distribution

**Return type** (dict)

**get\_action** (*observation*)

Runs a single observation through the specified policy and samples an action

**Parameters** **observation** (*ndarray*) – single observation - shape: (obs\_dim,)

**Returns** single action - shape: (action\_dim,)

**Return type** (ndarray)

**get\_actions** (*observations*)

Runs each set of observations through each task specific policy

**Parameters** **observations** (*ndarray*) – array of observations - shape: (batch\_size, obs\_dim)

**Returns** array of sampled actions - shape: (batch\_size, action\_dim)

**Return type** (ndarray)

**get\_param\_values** ()

Gets a list of all the current weights in the network (in original code it is flattened, why?)

**Returns** list of values for parameters

**Return type** (list)

**get\_params** ()

Get the tf.Variables representing the trainable weights of the network (symbolic)

**Returns** a dict of all trainable Variables

**Return type** (dict)

**likelihood\_ratio\_sym** (*obs, action, dist\_info\_old, policy\_params*)

Computes the likelihood  $p_{\text{new}}(\text{obs})/p_{\text{old}}$  ratio between

**Parameters**

- **obs** (*tf.Tensor*) – symbolic variable for observations
- **action** (*tf.Tensor*) – symbolic variable for actions
- **dist\_info\_old** (*dict*) – dictionary of tf.placeholders with old policy information
- **policy\_params** (*dict*) – dictionary of the policy parameters (each value is a *tf.Tensor*)

**Returns** likelihood ratio

**Return type** (*tf.Tensor*)

**load\_params** (*policy\_params*)

**Parameters** **policy\_params** (*ndarray*) – array of policy parameters for each task

**log\_diagnostics** (*paths, prefix=""*)

Log extra information per iteration based on the collected paths

**log\_likelihood\_sym** (*obs, action, policy\_params*)

Computes the log likelihood  $p(\text{obs})$

**Parameters**

- **obs** (*tf.Tensor*) – symbolic variable for observations
- **action** (*tf.Tensor*) – symbolic variable for actions
- **policy\_params** (*dict*) – dictionary of the policy parameters (each value is a *tf.Tensor*)

**Returns** log likelihood

**Return type** (*tf.Tensor*)

**set\_params** (*policy\_params*)

Sets the parameters for the graph

**Parameters** **policy\_params** (*dict*) – of variable names and corresponding parameter values

**class** meta\_policy\_search.policies.**MetaGaussianMLPPolicy** (*meta\_batch\_size, \*args, \*\*kwargs*)

Bases: meta\_policy\_search.policies.gaussian\_mlp\_policy.GaussianMLPPolicy, meta\_policy\_search.policies.base.MetaPolicy

**build\_graph** ()

Builds computational graph for policy

**distribution**

Returns this policy's distribution

**Returns** this policy's distribution

**Return type** (Distribution)

**distribution\_info\_keys** (*obs, state\_infos*)

**Parameters**

- **obs** (*placeholder*) – symbolic variable for observations
- **state\_infos** (*dict*) – a dictionary of placeholders that contains information about the
- **of the policy at the time it received the observation** (*state*)

**Returns** a dictionary of tf placeholders for the policy output distribution

**Return type** (dict)

**distribution\_info\_sym** (*obs\_var, params=None*)

Return the symbolic distribution information about the actions.

**Parameters**

- **obs\_var** (*placeholder*) – symbolic variable for observations
- **params** (*dict*) – a dictionary of placeholders or vars with the parameters of the MLP

**Returns** a dictionary of tf placeholders for the policy output distribution

**Return type** (dict)

**get\_action** (*observation, task=0*)

Runs a single observation through the specified policy and samples an action

**Parameters** **observation** (*ndarray*) – single observation - shape: (obs\_dim,)

**Returns** single action - shape: (action\_dim,)

**Return type** (ndarray)

**get\_actions** (*observations*)

**Parameters** **observations** (*list*) – List of numpy arrays of shape (meta\_batch\_size, batch\_size, obs\_dim)

**Returns** A tuple containing a list of numpy arrays of action, and a list of list of dicts of agent infos

**Return type** (tuple)

**get\_param\_values** ()

Gets a list of all the current weights in the network (in original code it is flattened, why?)

**Returns** list of values for parameters

**Return type** (list)

**get\_params** ()

Get the tf.Variables representing the trainable weights of the network (symbolic)

**Returns** a dict of all trainable Variables

**Return type** (dict)

**likelihood\_ratio\_sym** (*obs, action, dist\_info\_old, policy\_params*)

Computes the likelihood  $p_{\text{new}}(\text{obs})/p_{\text{old}}$  ratio between



**Parameters**

- **obs** (*tf.Tensor*) – symbolic variable for observations
- **action** (*tf.Tensor*) – symbolic variable for actions
- **dist\_info\_old** (*dict*) – dictionary of tf.placeholders with old policy information
- **policy\_params** (*dict*) – dictionary of the policy parameters (each value is a *tf.Tensor*)

**Returns** likelihood ratio

**Return type** (*tf.Tensor*)

**load\_params** (*policy\_params*)

**Parameters** **policy\_params** (*ndarray*) – array of policy parameters for each task

**log\_diagnostics** (*paths, prefix=""*)

Log extra information per iteration based on the collected paths

**log\_likelihood\_sym** (*obs, action, policy\_params*)

Computes the log likelihood  $p(\text{obs}|\text{act})$

**Parameters**

- **obs** (*tf.Tensor*) – symbolic variable for observations
- **action** (*tf.Tensor*) – symbolic variable for actions
- **policy\_params** (*dict*) – dictionary of the policy parameters (each value is a *tf.Tensor*)

**Returns** log likelihood

**Return type** (*tf.Tensor*)

**policies\_params\_feed\_dict**

returns fully prepared feed dict for feeding the currently saved policy parameter values into the lightweight policy graph

**set\_params** (*policy\_params*)

Sets the parameters for the graph

**Parameters** **policy\_params** (*dict*) – of variable names and corresponding parameter values

**switch\_to\_pre\_update** ()

Switches `get_action` to pre-update policy

**update\_task\_parameters** (*updated\_policies\_parameters*)

**Parameters**

- **updated\_policies\_parameters** (*list*) – List of size meta-batch size. Each contains a dict with the policies
- **as numpy arrays** (*parameters*) –

## 1.6 Samplers

### 1.6.1 Sampler

**class** meta\_policy\_search.samplers.Sampler(*env, policy, batch\_size, max\_path\_length*)

Bases: object

Sampler interface

#### Parameters

- **env** (*gym.Env*) – environment object
- **policy** (*meta\_policy\_search.policies.policy*) – policy object
- **batch\_size** (*int*) – number of trajectories per task
- **max\_path\_length** (*int*) – max number of steps per trajectory

**obtain\_samples** ()

Collect *batch\_size* trajectories

**Returns** A list of paths.

**Return type** (list)

**class** meta\_policy\_search.samplers.MetaSampler(*env, policy, rollouts\_per\_meta\_task, meta\_batch\_size, max\_path\_length, envs\_per\_task=None, parallel=False*)

Bases: meta\_policy\_search.samplers.base.Sampler

Sampler for Meta-RL

#### Parameters

- **env** (*meta\_policy\_search.envs.base.MetaEnv*) – environment object
- **policy** (*meta\_policy\_search.policies.base.Policy*) – policy object
- **batch\_size** (*int*) – number of trajectories per task
- **meta\_batch\_size** (*int*) – number of meta tasks
- **max\_path\_length** (*int*) – max number of steps per trajectory
- **envs\_per\_task** (*int*) – number of envs to run vectorized for each task (influences the memory usage)

**obtain\_samples** (*log=False, log\_prefix=""*)

Collect *batch\_size* trajectories from each task

#### Parameters

- **log** (*boolean*) – whether to log sampling times
- **log\_prefix** (*str*) – prefix for logger

**Returns** A dict of paths of size [*meta\_batch\_size*] x (*batch\_size*) x [5] x (*max\_path\_length*)

**Return type** (dict)

**update\_tasks** ()

Samples a new goal for each meta task

## 1.6.2 Sample Processor

```
class meta_policy_search.samplers.SampleProcessor (baseline,           discount=0.99,
                                                gae_lambda=1,           normalize_
                                                adv=False,             adv=True,
                                                positive_adv=False,    positive_adv=True,
                                                return_baseline=None)
```

Bases: object

### Sample processor interface

- fits a reward baseline (use zero baseline to skip this step)
- performs Generalized Advantage Estimation to provide advantages (see Schulman et al. 2015 - <https://arxiv.org/abs/1506.02438>)

### Parameters

- **baseline** (*Baseline*) – a reward baseline object
- **discount** (*float*) – reward discount factor
- **gae\_lambda** (*float*) – Generalized Advantage Estimation lambda
- **normalize\_adv** (*bool*) – indicates whether to normalize the estimated advantages (zero mean and unit std)
- **positive\_adv** (*bool*) – indicates whether to shift the (normalized) advantages so that they are all positive

```
process_samples (paths, log=False, log_prefix="")
```

### Processes sampled paths. This involves:

- computing discounted rewards (returns)
- fitting baseline estimator using the path returns and predicting the return baselines
- estimating the advantages using GAE (+ advantage normalization if desired)
- stacking the path data
- logging statistics of the paths

### Parameters

- **paths** (*list*) – A list of paths of size (batch\_size) x [5] x (max\_path\_length)
- **log** (*boolean*) – indicates whether to log
- **log\_prefix** (*str*) – prefix for the logging keys

**Returns** Processed sample data of size [7] x (batch\_size x max\_path\_length)

**Return type** (dict)

```
class meta_policy_search.samplers.DiceSampleProcessor (baseline,    max_path_length,
                                                discount=0.99,
                                                gae_lambda=1,           normalize_
                                                adv=True,             adv=True,
                                                positive_adv=False,    positive_adv=True,
                                                return_baseline=None)
```

Bases: meta\_policy\_search.samplers.base.SampleProcessor

### Sample processor for DICE implementations

- fits a reward baseline (use zero baseline to skip this step)
- computes adjusted rewards (reward - baseline)
- normalize adjusted rewards if desired
- zero-pads paths to `max_path_length`
- stacks the padded path data

### Parameters

- **baseline** (*Baseline*) – a time dependent reward baseline object
- **max\_path\_length** (*int*) – maximum path length
- **discount** (*float*) – reward discount factor
- **normalize\_adv** (*bool*) – indicates whether to normalize the estimated advantages (zero mean and unit std)
- **positive\_adv** (*bool*) – indicates whether to shift the (normalized) advantages so that they are all positive
- **return\_baseline** (*Baseline*) – (optional) a state(-time) dependent baseline - if provided it is also fitted and used to calculate GAE advantage estimates

**process\_samples** (*paths*, *log=False*, *log\_prefix=""*)

### Processes sampled paths, This involves:

- computing discounted rewards
- fitting a reward baseline
- computing adjusted rewards (reward - baseline)
- normalizing adjusted rewards if desired
- stacking the padded path data
- creating a mask which indicates padded values by zero and original values by one
- logging statistics of the paths

### Parameters

- **paths** (*list*) – A list of paths of size (batch\_size) x [5] x (max\_path\_length)
- **log** (*boolean*) – indicates whether to log
- **log\_prefix** (*str*) – prefix for the logging keys

### Returns

**Processed sample data. A dict containing the following items with respective shapes:**

- **mask**: (batch\_size, max\_path\_length)
- **observations**: (batch\_size, max\_path\_length, ndim\_act)
- **actions**: (batch\_size, max\_path\_length, ndim\_obs)
- **rewards**: (batch\_size, max\_path\_length)
- **adjusted\_rewards**: (batch\_size, max\_path\_length)
- **env\_infos**: dict of ndarrays of shape (batch\_size, max\_path\_length, ?)

- `agent_infos`: dict of ndarrays of shape (batch\_size, max\_path\_length, ?)

**Return type** (dict)

```
class meta_policy_search.samplers.MetaSampleProcessor (baseline, discount=0.99,
                                                    gae_lambda=1, normalize_adv=False, positive_adv=False)
```

Bases: meta\_policy\_search.samplers.base.SampleProcessor

```
process_samples (paths_meta_batch, log=False, log_prefix="")
```

**Processes sampled paths. This involves:**

- computing discounted rewards (returns)
- fitting baseline estimator using the path returns and predicting the return baselines
- estimating the advantages using GAE (+ advantage normalization if desired)
- stacking the path data
- logging statistics of the paths

**Parameters**

- **paths\_meta\_batch** (*dict*) – A list of dict of lists, size: [meta\_batch\_size] x (batch\_size) x [5] x (max\_path\_length)
- **log** (*boolean*) – indicates whether to log
- **log\_prefix** (*str*) – prefix for the logging keys

**Returns** Processed sample data among the meta-batch; size: [meta\_batch\_size] x [7] x (batch\_size x max\_path\_length)

**Return type** (list of dicts)

### 1.6.3 Vectorized Environment Executor

```
class meta_policy_search.samplers.vectorized_env_executor.MetaIterativeEnvExecutor (env,
                                                                              meta_batch_size,
                                                                              envs_per_task,
                                                                              max_path_length)
```

Bases: object

Wraps multiple environments of the same kind and provides functionality to reset / step the environments in a vectorized manner. Internally, the environments are executed iteratively.

**Parameters**

- **env** (meta\_policy\_search.envs.base.MetaEnv) – meta environment object
- **meta\_batch\_size** (*int*) – number of meta tasks
- **envs\_per\_task** (*int*) – number of environments per meta task
- **max\_path\_length** (*int*) – maximum length of sampled environment paths - if the max\_path\_length is reached, the respective environment is reset

**num\_envs**

Number of environments

**Returns** number of environments

**Return type** (int)

**reset** ()

Resets the environments

**Returns** list of (np.ndarray) with the new initial observations.

**Return type** (list)

**set\_tasks** (*tasks*)

Sets a list of tasks to each environment

**Parameters** *tasks* (*list*) – list of the tasks for each environment

**step** (*actions*)

Steps the wrapped environments with the provided actions

**Parameters** *actions* (*list*) – lists of actions, of length meta\_batch\_size x envs\_per\_task

**Returns**

**(tuple): a length 4 tuple of lists, containing obs (np.array), rewards (float), dones (bool), env\_infos (dict).** Each list is of length meta\_batch\_size x envs\_per\_task (assumes that every task has same number of envs)

**class** meta\_policy\_search.samplers.vectorized\_env\_executor.**MetaParallelEnvExecutor** (*env*, *meta\_batch\_size*, *envs\_per\_task*, *max\_path\_length*)

Bases: object

Wraps multiple environments of the same kind and provides functionality to reset / step the environments in a vectorized manner. Thereby the environments are distributed among meta\_batch\_size processes and executed in parallel.

**Parameters**

- **env** (*meta\_policy\_search.envs.base.MetaEnv*) – meta environment object
- **meta\_batch\_size** (*int*) – number of meta tasks
- **envs\_per\_task** (*int*) – number of environments per meta task
- **max\_path\_length** (*int*) – maximum length of sampled environment paths - if the max\_path\_length is reached, the respective environment is reset

**num\_envs**

Number of environments

**Returns** number of environments

**Return type** (int)

**reset** ()

Resets the environments of each worker

**Returns** list of (np.ndarray) with the new initial observations.

**Return type** (list)

**set\_tasks** (*tasks=None*)

Sets a list of tasks to each worker

**Parameters** *tasks* (*list*) – list of the tasks for each worker

**step** (*actions*)

Executes actions on each env

**Parameters** **actions** (*list*) – lists of actions, of length meta\_batch\_size x envs\_per\_task

**Returns**

**(tuple): a length 4 tuple of lists, containing obs (np.array), rewards (float), dones (bool), env\_infos (dict)**  
each list is of length meta\_batch\_size x envs\_per\_task (assumes that every task has same number of envs)

## 1.7 Meta-Trainer

**class** meta\_policy\_search.meta\_trainer.Trainer (*algo, env, sampler, sample\_processor, policy, n\_itr, start\_itr=0, num\_inner\_grad\_steps=1, sess=None*)

Bases: object

Performs steps of meta-policy search.

Pseudocode:

```

for iter in n_iter:
    sample tasks
    for task in tasks:
        for adapt_step in num_inner_grad_steps
            sample trajectories with policy
            perform update/adaptation step
            sample trajectories with post-update policy
        perform meta-policy gradient step(s)

```

**Parameters**

- **algo** (*Algo*) –
- **env** (*Env*) –
- **sampler** (*Sampler*) –
- **sample\_processor** (*SampleProcessor*) –
- **baseline** (*Baseline*) –
- **policy** (*Policy*) –
- **n\_itr** (*int*) – Number of iterations to train for
- **start\_itr** (*int*) – Number of iterations policy has already trained for, if reloading
- **num\_inner\_grad\_steps** (*int*) – Number of inner steps per maml iteration
- **sess** (*tf.Session*) – current tf session (if we loaded policy, for example)

**get\_itr\_snapshot** (*itr*)

Gets the current policy and env for storage

**train** ()

Trains policy on env using algo

Pseudocode:

```
for itr in n_itr:
    for step in num_inner_grad_steps:
        sampler.sample()
        algo.compute_updated_dists()
    algo.optimize_policy()
    sampler.update_goals()
```



## CHAPTER 2

---

### Indices and tables

---

- genindex
- modindex



## m

`meta_policy_search.baselines`, 3  
`meta_policy_search.envs.base`, 5  
`meta_policy_search.meta_algos`, 6  
`meta_policy_search.meta_trainer`, 27  
`meta_policy_search.optimizers`, 10  
`meta_policy_search.policies`, 13  
`meta_policy_search.samplers`, 22  
`meta_policy_search.samplers.vectorized_env_executor`,  
25



**B**

Baseline (class in `meta_policy_search.baselines`), 3  
 build\_graph() (`meta_policy_search.meta_algos.MAML` `method`), 6  
 build\_graph() (`meta_policy_search.meta_algos.ProMP` `method`), 7  
 build\_graph() (`meta_policy_search.meta_algos.TRPO` `method`), 8  
 build\_graph() (`meta_policy_search.meta_algos.VPG` `method`), 9  
 build\_graph() (`meta_policy_search.optimizers.ConjugateGradientOptimizer` `method`), 10  
 build\_graph() (`meta_policy_search.optimizers.MAMLFirstOrderOptimizer` `method`), 12  
 build\_graph() (`meta_policy_search.policies.GaussianMLPPolicy` `method`), 18  
 build\_graph() (`meta_policy_search.policies.MetaGaussianMLPPolicy` `method`), 19  
 build\_graph() (`meta_policy_search.policies.MetaPolicy` `method`), 15  
 build\_graph() (`meta_policy_search.policies.Policy` `method`), 13

**C**

ConjugateGradientOptimizer (class in `meta_policy_search.optimizers`), 10  
 constraint\_val() (`meta_policy_search.optimizers.ConjugateGradientOptimizer` `method`), 11

**D**

DiceSampleProcessor (class in `meta_policy_search.samplers`), 23  
 distribution(`meta_policy_search.policies.GaussianMLPPolicy` `attribute`), 18  
 distribution(`meta_policy_search.policies.MetaGaussianMLPPolicy` `attribute`), 19  
 distribution(`meta_policy_search.policies.MetaPolicy` `attribute`), 15  
 fit() (`meta_policy_search.baselines.Baseline` `method`), 3  
 fit() (`meta_policy_search.baselines.LinearFeatureBaseline` `method`), 4  
 fit() (`meta_policy_search.baselines.LinearTimeBaseline` `method`), 4  
 GaussianMLPPolicy (class in `meta_policy_search.policies`), 17  
 get\_action() (`meta_policy_search.policies.GaussianMLPPolicy` `method`), 18

`get_action()` (*meta\_policy\_search.policies.MetaGaussianMLPPolicy* method), 19  
`get_action()` (*meta\_policy\_search.policies.MetaPolicy* method), 16  
`get_action()` (*meta\_policy\_search.policies.Policy* method), 14  
`get_actions()` (*meta\_policy\_search.policies.GaussianMLPPolicy* method), 16  
`get_actions()` (*meta\_policy\_search.policies.MetaGaussianMLPPolicy* method), 20  
`get_actions()` (*meta\_policy\_search.policies.MetaPolicy* method), 16  
`get_actions()` (*meta\_policy\_search.policies.Policy* method), 14  
`get_itr_snapshot()` (*meta\_policy\_search.meta\_trainer.Trainer* method), 27  
`get_param_values()` (*meta\_policy\_search.baselines.Baseline* method), 3  
`get_param_values()` (*meta\_policy\_search.baselines.LinearFeatureBaseline* method), 4  
`get_param_values()` (*meta\_policy\_search.baselines.LinearTimeBaseline* method), 4  
`get_param_values()` (*meta\_policy\_search.policies.GaussianMLPPolicy* method), 18  
`get_param_values()` (*meta\_policy\_search.policies.MetaGaussianMLPPolicy* method), 20  
`get_param_values()` (*meta\_policy\_search.policies.MetaPolicy* method), 16  
`get_param_values()` (*meta\_policy\_search.policies.Policy* method), 14  
`get_params()` (*meta\_policy\_search.policies.GaussianMLPPolicy* method), 19  
`get_params()` (*meta\_policy\_search.policies.MetaGaussianMLPPolicy* method), 20  
`get_params()` (*meta\_policy\_search.policies.MetaPolicy* method), 16  
`get_params()` (*meta\_policy\_search.policies.Policy* method), 14  
`get_task()` (*meta\_policy\_search.envs.base.MetaEnv* method), 5  
`gradient()` (*meta\_policy\_search.optimizers.ConjugateGradientOptimizer* method), 11  
**L**  
`likelihood_ratio_sym()` (*meta\_policy\_search.policies.Policy* method),  
(*meta\_policy\_search.policies.GaussianMLPPolicy* method), 15

loss () (meta\_policy\_search.optimizers.ConjugateGradientOptimizer (method), 11)

loss () (meta\_policy\_search.optimizers.MAMLFirstOrderOptimizer (method), 12)

## M

make\_vars () (meta\_policy\_search.meta\_algos.MAMLAlgo (method), 6)

make\_vars () (meta\_policy\_search.meta\_algos.ProMP (method), 7)

make\_vars () (meta\_policy\_search.meta\_algos.TRPOMAML (method), 8)

make\_vars () (meta\_policy\_search.meta\_algos.VPGMAML (method), 9)

MAMLAlgo (class in meta\_policy\_search.meta\_algos), 6

MAMLFirstOrderOptimizer (class in meta\_policy\_search.optimizers), 12

meta\_policy\_search.baselines (module), 3

meta\_policy\_search.envs.base (module), 5

meta\_policy\_search.meta\_algos (module), 6

meta\_policy\_search.meta\_trainer (module), 27

meta\_policy\_search.optimizers (module), 10

meta\_policy\_search.policies (module), 13

meta\_policy\_search.samplers (module), 22

meta\_policy\_search.samplers.vectorized\_env\_executor (module), 25

MetaEnv (class in meta\_policy\_search.envs.base), 5

MetaGaussianMLPPolicy (class in meta\_policy\_search.policies), 19

MetaIterativeEnvExecutor (class in meta\_policy\_search.samplers.vectorized\_env\_executor), 25

MetaParallelEnvExecutor (class in meta\_policy\_search.samplers.vectorized\_env\_executor), 26

MetaPolicy (class in meta\_policy\_search.policies), 15

MetaSampleProcessor (class in meta\_policy\_search.samplers), 25

MetaSampler (class in meta\_policy\_search.samplers), 22

## N

num\_envs (meta\_policy\_search.samplers.vectorized\_env\_executor (attribute), 25)

num\_envs (meta\_policy\_search.samplers.vectorized\_env\_executor (attribute), 26)

## O

obtain\_samples () (meta\_policy\_search.samplers.MetaSampler (method), 22)

obtain\_samples () (meta\_policy\_search.samplers.SampleProcessor (method), 22)

optimize () (meta\_policy\_search.optimizers.ConjugateGradientOptimizer (method), 11)

optimize () (meta\_policy\_search.optimizers.MAMLFirstOrderOptimizer (method), 12)

optimize\_policy () (meta\_policy\_search.meta\_algos.MAMLAlgo (method), 6)

optimize\_policy () (meta\_policy\_search.meta\_algos.ProMP (method), 7)

optimize\_policy () (meta\_policy\_search.meta\_algos.TRPOMAML (method), 8)

optimize\_policy () (meta\_policy\_search.meta\_algos.VPGMAML (method), 9)

optimize\_policy () (meta\_policy\_search.meta\_algos.VPGMAML (method), 9)

optimize\_policy () (meta\_policy\_search.meta\_algos.VPGMAML (method), 9)

## P

policies\_params\_feed\_dict (meta\_policy\_search.policies.MetaGaussianMLPPolicy (attribute), 21)

policies\_params\_feed\_dict (meta\_policy\_search.policies.MetaPolicy (attribute), 17)

Policy (class in meta\_policy\_search.policies), 13

predict () (meta\_policy\_search.baselines.Baseline (method), 3)

predict () (meta\_policy\_search.baselines.LinearFeatureBaseline (method), 4)

predict () (meta\_policy\_search.baselines.LinearTimeBaseline (method), 5)

process\_samples () (meta\_policy\_search.samplers.DiceSampleProcessor (method), 24)

process\_samples () (meta\_policy\_search.samplers.MetaSampleProcessor (method), 25)

process\_samples () (meta\_policy\_search.samplers.SampleProcessor (method), 23)

process\_samples () (meta\_policy\_search.samplers.SampleProcessor (method), 23)

process\_samples () (meta\_policy\_search.samplers.SampleProcessor (method), 23)

ProMP (class in meta\_policy\_search.meta\_algos), 7

## R

render () (meta\_policy\_search.samplers.vectorized\_env\_executor (method), 26)

render () (meta\_policy\_search.samplers.vectorized\_env\_executor (method), 26)

## S

sample\_tasks () (meta\_policy\_search.envs.base.MetaEnv (method), 5)

sample\_tasks () (meta\_policy\_search.samplers (class in meta\_policy\_search.samplers), 23)

Sampler (class in meta\_policy\_search.samplers), 22

*set\_params()* (*meta\_policy\_search.baselines.Baseline*  
*method*), 3  
*set\_params()* (*meta\_policy\_search.baselines.LinearFeatureBaseline*  
*method*), 4  
*set\_params()* (*meta\_policy\_search.baselines.LinearTimeBaseline*  
*method*), 5  
*set\_params()* (*meta\_policy\_search.policies.GaussianMLPPolicy*  
*method*), 19  
*set\_params()* (*meta\_policy\_search.policies.MetaGaussianMLPPolicy*  
*method*), 21  
*set\_params()* (*meta\_policy\_search.policies.MetaPolicy*  
*method*), 17  
*set\_params()* (*meta\_policy\_search.policies.Policy*  
*method*), 15  
*set\_task()* (*meta\_policy\_search.envs.base.MetaEnv*  
*method*), 5  
*set\_tasks()* (*meta\_policy\_search.samplers.vectorized\_env\_executor.MetaIterativeEnvExecutor*  
*method*), 26  
*set\_tasks()* (*meta\_policy\_search.samplers.vectorized\_env\_executor.MetaParallelEnvExecutor*  
*method*), 26  
*step()* (*meta\_policy\_search.samplers.vectorized\_env\_executor.MetaIterativeEnvExecutor*  
*method*), 26  
*step()* (*meta\_policy\_search.samplers.vectorized\_env\_executor.MetaParallelEnvExecutor*  
*method*), 26  
*switch\_to\_pre\_update()*  
(*meta\_policy\_search.policies.MetaGaussianMLPPolicy*  
*method*), 21  
*switch\_to\_pre\_update()*  
(*meta\_policy\_search.policies.MetaPolicy*  
*method*), 17

## T

*train()* (*meta\_policy\_search.meta\_trainer.Trainer*  
*method*), 27  
Trainer (*class in meta\_policy\_search.meta\_trainer*),  
27  
TRPOMAML (*class in meta\_policy\_search.meta\_algos*), 8

## U

*update\_task\_parameters()*  
(*meta\_policy\_search.policies.MetaGaussianMLPPolicy*  
*method*), 21  
*update\_task\_parameters()*  
(*meta\_policy\_search.policies.MetaPolicy*  
*method*), 17  
*update\_tasks()* (*meta\_policy\_search.samplers.MetaSampler*  
*method*), 22

## V

VPGMAML (*class in meta\_policy\_search.meta\_algos*), 8