

# Project: Train a Quadcopter How to Fly

Design an agent to fly a quadcopter, and then train it using a reinforcement learning algorithm of your choice!

Try to apply the techniques you have learnt, but also feel free to come up with innovative ideas and test them.

## Instructions

Take a look at the files in the directory to better understand the structure of the project.

- `task.py` : Define your task (environment) in this file.
- `agents/` : Folder containing reinforcement learning agents.
  - `policy_search.py` : A sample agent has been provided here.
  - `agent.py` : Develop your agent here.
- `physics_sim.py` : This file contains the simulator for the quadcopter. **DO NOT MODIFY THIS FILE.**

For this project, you will define your own task in `task.py`. Although we have provided a example task to get you started, you are encouraged to change it. Later in this notebook, you will learn more about how to amend this file.

You will also design a reinforcement learning agent in `agent.py` to complete your chosen task.

You are welcome to create any additional files to help you to organize your code. For instance, you may find it useful to define a `model.py` file defining any needed neural network architectures.

## Controlling the Quadcopter

We provide a sample agent in the code cell below to show you how to use the sim to control the quadcopter. This agent is even simpler than the sample agent that you'll examine (in `agents/policy_search.py`) later in this notebook!

The agent controls the quadcopter by setting the revolutions per second on each of its four rotors. The provided agent in the `Basic_Agent` class below always selects a random action for each of the four rotors. These four speeds are returned by the `act` method as a list of four floating-point numbers.

For this project, the agent that you will implement in `agents/agent.py` will have a far more intelligent method for selecting actions!

In [1]:

```
import random

class Basic_Agent():
    def __init__(self, task):
        self.task = task

    def act(self):
        new_thrust = random.gauss(450., 25.)
        return [new_thrust + random.gauss(0., 1.) for x in range(4)]
```

Run the code cell below to have the agent select actions to control the quadcopter.

Feel free to change the provided values of `runtime`, `init_pose`, `init_velocities`, and `init_angle_velocities` below to change the starting conditions of the quadcopter.

The `labels` list below annotates statistics that are saved while running the simulation. All of this information is saved in a text file `data.txt` and stored in the dictionary `results`.

In [2]:

```
%load_ext autoreload
%autoreload 2

import csv
import numpy as np
from task import Task

# Modify the values below to give the quadcopter a different starting position.
runtime = 5. # time limit of the episode
init_pose = np.array([0., 0., 10., 0., 0., 0.]) # initial pose
init_velocities = np.array([0., 0., 0.]) # initial velocities
init_angle_velocities = np.array([0., 0., 0.]) # initial angle velocities
file_output = 'data.txt' # file name for saved results

# Setup
task = Task(init_pose, init_velocities, init_angle_velocities, runtime)
agent = Basic_Agent(task)
done = False
labels = ['time', 'x', 'y', 'z', 'phi', 'theta', 'psi', 'x_velocity',
          'y_velocity', 'z_velocity', 'phi_velocity', 'theta_velocity',
          'psi_velocity', 'rotor_speed1', 'rotor_speed2', 'rotor_speed3', 'rotor_sp
results = {x : [] for x in labels}

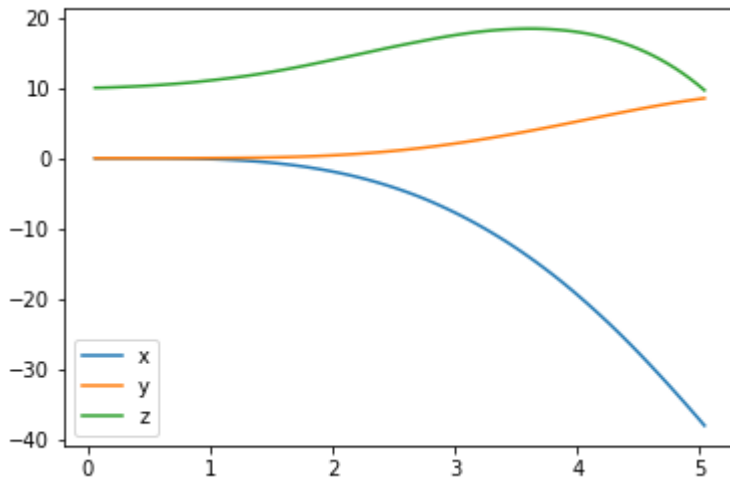
# Run the simulation, and save the results.
with open(file_output, 'w') as csvfile:
    writer = csv.writer(csvfile)
    writer.writerow(labels)
    while True:
        rotor_speeds = agent.act()
        _, _, done = task.step(rotor_speeds)
        to_write = [task.sim.time] + list(task.sim.pose) + list(task.sim.v) + list(
            for ii in range(len(labels)):
                results[labels[ii]].append(to_write[ii])
        writer.writerow(to_write)
        if done:
            break
```

Run the code cell below to visualize how the position of the quadcopter evolved during the simulation.

In [3]:

```
import matplotlib.pyplot as plt
%matplotlib inline

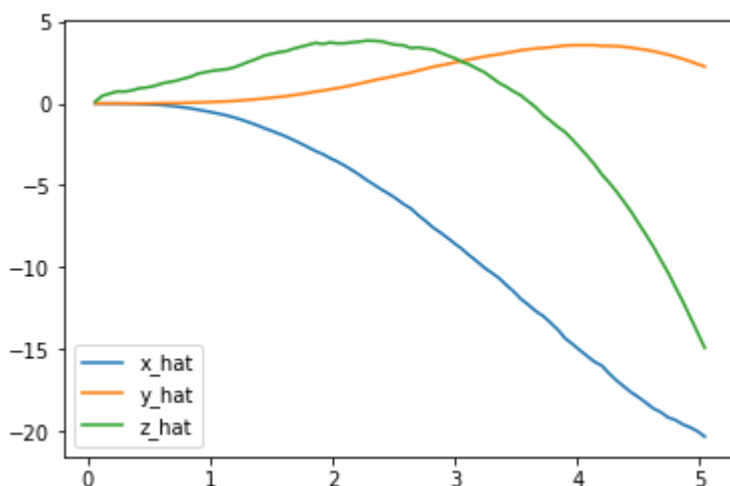
plt.plot(results['time'], results['x'], label='x')
plt.plot(results['time'], results['y'], label='y')
plt.plot(results['time'], results['z'], label='z')
plt.legend()
_ = plt.ylim()
```



The next code cell visualizes the velocity of the quadcopter.

In [4]:

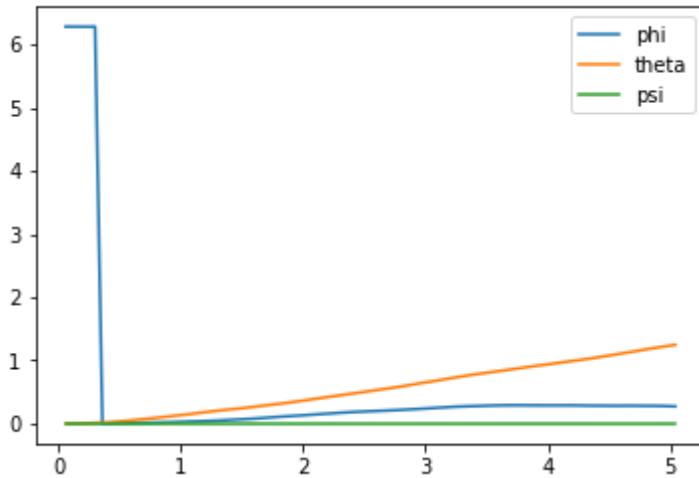
```
plt.plot(results['time'], results['x_velocity'], label='x_hat')
plt.plot(results['time'], results['y_velocity'], label='y_hat')
plt.plot(results['time'], results['z_velocity'], label='z_hat')
plt.legend()
_ = plt.ylim()
```



Next, you can plot the Euler angles (the rotation of the quadcopter over the  $x$ -,  $y$ -, and  $z$ -axes),

In [5]:

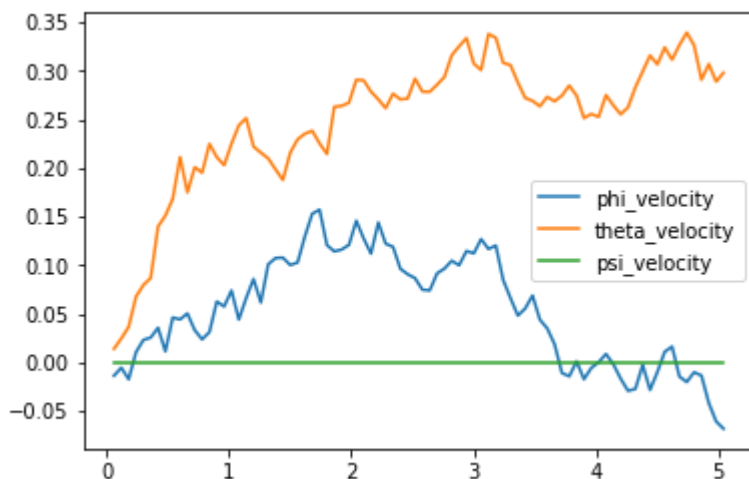
```
plt.plot(results['time'], results['phi'], label='phi')
plt.plot(results['time'], results['theta'], label='theta')
plt.plot(results['time'], results['psi'], label='psi')
plt.legend()
_ = plt.ylim()
```



before plotting the velocities (in radians per second) corresponding to each of the Euler angles.

In [6]:

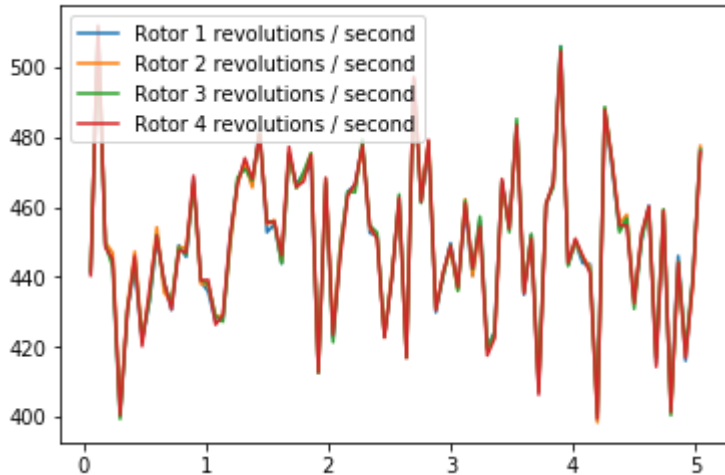
```
plt.plot(results['time'], results['phi_velocity'], label='phi_velocity')
plt.plot(results['time'], results['theta_velocity'], label='theta_velocity')
plt.plot(results['time'], results['psi_velocity'], label='psi_velocity')
plt.legend()
_ = plt.ylim()
```



Finally, you can use the code cell below to print the agent's choice of actions.

In [7]:

```
plt.plot(results['time'], results['rotor_speed1'], label='Rotor 1 revolutions / second')
plt.plot(results['time'], results['rotor_speed2'], label='Rotor 2 revolutions / second')
plt.plot(results['time'], results['rotor_speed3'], label='Rotor 3 revolutions / second')
plt.plot(results['time'], results['rotor_speed4'], label='Rotor 4 revolutions / second')
plt.legend()
_ = plt.ylim()
```



When specifying a task, you will derive the environment state from the simulator. Run the code cell below to print the values of the following variables at the end of the simulation:

- `task.sim.pose` (the position of the quadcopter in  $(x, y, z)$  dimensions and the Euler angles),
- `task.sim.v` (the velocity of the quadcopter in  $(x, y, z)$  dimensions), and
- `task.sim.angular_v` (radians/second for each of the three Euler angles).

In [8]:

```
# the pose, velocity, and angular velocity of the quadcopter at the end of the episode
print(task.sim.pose)
print(task.sim.v)
print(task.sim.angular_v)
```

```
[ -37.97685458   8.51616276   9.66000427   0.27242426   1.24666878
    0.          ]
[ -20.34653995   2.25242229 -14.92625181]
[ -0.06812903   0.29787148   0.          ]
```

In the sample task in `task.py`, we use the 6-dimensional pose of the quadcopter to construct the state of the environment at each timestep. However, when amending the task for your purposes, you are welcome to expand the size of the state vector by including the velocity information. You can use any combination of the pose, velocity, and angular velocity - feel free to tinker here, and construct the state to suit your task.

## The Task

A sample task has been provided for you in `task.py`. Open this file in a new window now.

The `__init__()` method is used to initialize several variables that are needed to specify the task.

- The simulator is initialized as an instance of the `PhysicsSim` class (from `physics_sim.py`).

- Inspired by the methodology in the original DDPG paper, we make use of action repeats. For each timestep of the agent, we step the simulation `action_repeats` timesteps. If you are not familiar with action repeats, please read the **Results** section in [the DDPG paper \(https://arxiv.org/abs/1509.02971\)](https://arxiv.org/abs/1509.02971).
- We set the number of elements in the state vector. For the sample task, we only work with the 6-dimensional pose information. To set the size of the state ( `state_size` ), we must take action repeats into account.
- The environment will always have a 4-dimensional action space, with one entry for each rotor ( `action_size=4` ). You can set the minimum ( `action_low` ) and maximum ( `action_high` ) values of each entry here.
- The sample task in this provided file is for the agent to reach a target position. We specify that target position as a variable.

The `reset()` method resets the simulator. The agent should call this method every time the episode ends. You can see an example of this in the code cell below.

The `step()` method is perhaps the most important. It accepts the agent's choice of action `rotor_speeds` , which is used to prepare the next state to pass on to the agent. Then, the reward is computed from `get_reward()` . The episode is considered done if the time limit has been exceeded, or the quadcopter has travelled outside of the bounds of the simulation.

In the next section, you will learn how to test the performance of an agent on this task.

## The Agent

The sample agent given in `agents/policy_search.py` uses a very simplistic linear policy to directly compute the action vector as a dot product of the state vector and a matrix of weights. Then, it randomly perturbs the parameters by adding some Gaussian noise, to produce a different policy. Based on the average reward obtained in each episode ( `score` ), it keeps track of the best set of parameters found so far, how the score is changing, and accordingly tweaks a scaling factor to widen or tighten the noise.

Run the code cell below to see how the agent performs on the sample task.

In [9]:

```

import sys
import pandas as pd
from agents.policy_search import PolicySearch_Agent
from task import Task

num_episodes = 1000
target_pos = np.array([0., 0., 10.])
task = Task(target_pos=target_pos)
agent = PolicySearch_Agent(task)

for i_episode in range(1, num_episodes+1):
    state = agent.reset_episode() # start a new episode
    while True:
        action = agent.act(state)
        next_state, reward, done = task.step(action)
        agent.step(reward, done)
        state = next_state
        if done:
            print("\rEpisode = {:4d}, score = {:.3f} (best = {:.3f}), noise_scale = {:.3f}
                  i_episode, agent.score, agent.best_score, agent.noise_scale), end="
            break
    sys.stdout.flush()

```

Episode = 1000, score = 1.213 (best = 1.539), noise\_scale = 3.25

This agent should perform very poorly on this task. And that's where you come in!

## Define the Task, Design the Agent, and Train Your Agent!

Amend `task.py` to specify a task of your choosing. If you're unsure what kind of task to specify, you may like to teach your quadcopter to takeoff, hover in place, land softly, or reach a target pose.

After specifying your task, use the sample agent in `agents/policy_search.py` as a template to define your own agent in `agents/agent.py`. You can borrow whatever you need from the sample agent, including ideas on how you might modularize your code (using helper methods like `act()`, `learn()`, `reset_episode()`, etc.).

Note that it is **highly unlikely** that the first agent and task that you specify will learn well. You will likely have to tweak various hyperparameters and the reward function for your task until you arrive at reasonably good behavior.

As you develop your agent, it's important to keep an eye on how it's performing. Use the code above as inspiration to build in a mechanism to log/save the total rewards obtained in each episode to file. If the episode rewards are gradually increasing, this is an indication that your agent is learning.

In [10]:

```

import sys
import shutil
import pandas as pd
import csv
import numpy as np
from agents.DDPG import DDPG
from task import Task

import matplotlib.pyplot as plt
%matplotlib notebook

# Task : Hover at (0, 0, 10.)

# Modify the values below to give the quadcopter a different starting position.
runtime = 5. # time limit of the episode
init_pose = np.array([0., 0., 10., 0., 0., 0.]) # initial pose
init_velocities = np.array([0., 0., 0.]) # initial velocities
init_angle_velocities = np.array([0., 0., 0.]) # initial angle velocities
file_output = 'data.txt' # file name for saved results
best_file_output = "best_data.txt" # file name for saved best results

num_episodes = 450
target_pos = np.array([0., 0., 10.])

# Setup task and agent(environment)
task = Task(init_pose=init_pose, init_velocities=init_velocities, init_angle_veloci
agent = DDPG(task)

labels = ['time', # time
          'x', 'y', 'z', 'phi', 'theta', 'psi', # pose
          'x_velocity', 'y_velocity', 'z_velocity', # v
          'phi_velocity', 'theta_velocity', 'psi_velocity'] # angular_v
results = {x : [] for x in labels}

# Function for track history in an episode
fig, ax = plt.subplots(1,2, figsize=(10,4))
def show_hist(t,x,y,z,r, ax=ax, fig=fig):
    x_h,=ax[0].plot(t,x,'r',label='x')
    y_h,=ax[0].plot(t,y,'g',label='y')
    z_h,=ax[0].plot(t,z,'b',label='z')
    ax[0].set_title("x,y,z history")
    ax[0].legend(handles=[x_h, y_h, z_h])
    ax[0].set_xlabel("time")

    reward_h,=ax[1].plot(t,r,'b',label='reward')
    ax[1].legend(handles=[reward_h])
    ax[1].set_title("reward history")
    ax[1].set_xlabel("time")
    fig.canvas.draw()

# DDPG training
total_rewards = []
best_total_reward = 0
for i_episode in range(1, num_episodes+1):
    state = agent.reset_episode()
    total_reward = 0
    reward_record = []

    # track status of current episode

```



```
t, x, y, z, r = [],[],[], [],[]
```

```
with open(file_output, 'w') as csvfile:
```

```
    writer = csv.writer(csvfile)
```

```
    writer.writerow(labels)
```

```
    while True:
```

```
        action = agent.act(state)
```

```
        next_state, reward, done = task.step(action)
```

```
        reward_record.append(reward)
```

```
        agent.step(action, reward, next_state, done)
```

```
        state = next_state
```

```
        # write quadcopter status in file
```

```
        to_write = [task.sim.time] + list(task.sim.pose) + list(task.sim.v) + l
```

```
        writer.writerow(to_write)
```

```
        # save quadcopter status
```

```
        t.append(task.sim.time)
```

```
        x.append(task.sim.pose[0])
```

```
        y.append(task.sim.pose[1])
```

```
        z.append(task.sim.pose[2])
```

```
        total_reward = sum(reward_record)
```

```
        r.append(total_reward)
```

```
    if done:
```

```
        # flush cache and close file of current episode
```

```
        csvfile.close()
```

```
        # save total_reward for every episode in total_rewards [for tracking]
```

```
        total_reward = sum(reward_record)
```

```
        total_rewards.append(total_reward)
```

```
        # if current total_reward is the best, copy current file to the best
```

```
        if total_reward > best_total_reward:
```

```
            best_total_reward = total_reward
```

```
            shutil.copyfile(file_output, best_file_output)
```

```
        print("\rEpisode = {:4d}, total_reward = {:.3f} (best = {:.3f})".format(
            i_episode, total_reward, best_total_reward), end="")
```

```
        # show quadcopter status history every 10 episodes [for monitoring]
```

```
        if (i_episode % 10 == 0):
```

```
            show_hist(t,x,y,z,r)
```

```
    break
```

Using TensorFlow backend.

```
Episode = 450, total_reward = 37.983 (best = 43.778)
```

## Plot the Rewards

Once you are satisfied with your performance, plot the episode rewards, either from a single run, or averaged over multiple runs.

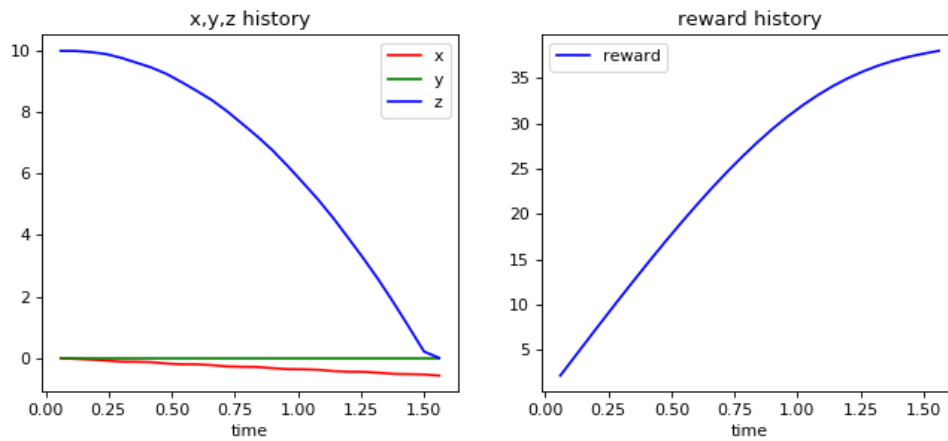
In [14]:

```

# Plot the last episode history
import matplotlib.pyplot as plt
%matplotlib notebook
fig, ax = plt.subplots(1,2, figsize=(10,4))
x_h,=ax[0].plot(t,x,'r',label='x')
y_h,=ax[0].plot(t,y,'g',label='y')
z_h,=ax[0].plot(t,z,'b',label='z')
ax[0].set_title("x,y,z history")
ax[0].legend(handles=[x_h, y_h, z_h])
ax[0].set_xlabel("time")

reward_h,=ax[1].plot(t,r,'b',label='reward')
ax[1].legend(handles=[reward_h])
ax[1].set_title("reward history")
ax[1].set_xlabel("time")
fig.canvas.draw()

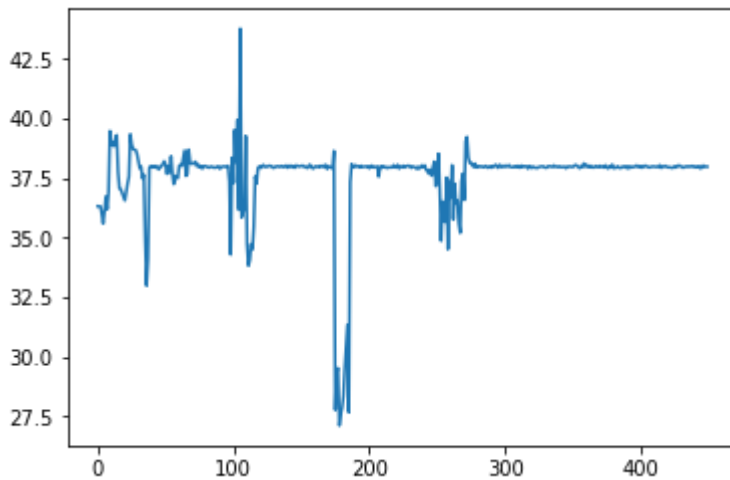
```



In [15]:

```
## TODO: Plot the rewards.  
import matplotlib.pyplot as plt  
%matplotlib inline  
  
plt.plot(range(len(total_rewards)), total_rewards)  
_ = plt.ylim()  
  
print("Last 10 episode average reward :", np.mean(total_rewards[:10]))
```

Last 10 episode average reward : 36.603818983939945



In [16]:

```
## Plot the final result.
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline

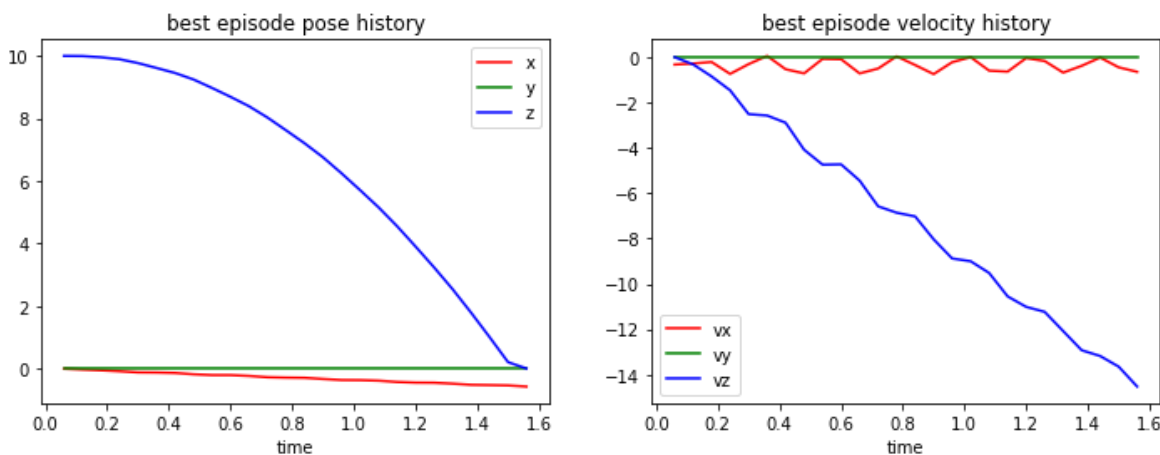
filename = "best_data.txt"
filename = "data.txt"

best_result = pd.read_csv(filename)

fig, ax = plt.subplots(1,2, figsize=(12,4))
x,=ax[0].plot(best_result['time'],best_result['x'],'r',label='x')
y,=ax[0].plot(best_result['time'],best_result['y'],'g',label='y')
z,=ax[0].plot(best_result['time'],best_result['z'],'b',label='z')
ax[0].set_title("best episode pose history")
ax[0].legend(handles=[x, y, z])
ax[0].set_xlabel("time")

vx,=ax[1].plot(best_result['time'],best_result['x_velocity'],'r',label='vx')
vy,=ax[1].plot(best_result['time'],best_result['y_velocity'],'g',label='vy')
vz,=ax[1].plot(best_result['time'],best_result['z_velocity'],'b',label='vz')

ax[1].legend(handles=[vx, vy, vz])
ax[1].set_title("best episode velocity history")
ax[1].set_xlabel("time")
fig.canvas.draw()
```



## Reflections

**Question 1:** Describe the task that you specified in `task.py`. How did you design the reward function?

**Answer:**

- I specified a hover task that quadcopter try to keep at  $(x,y,z) = (0, 0, 10)$
- The reward function is defined as  $\text{sigmoid}(1 - 0.3 * |\text{pose} - \text{target\_pose}|)$ . That is, current pose getting far away from `target_pose` should be penalized. Also, sigmoid function is used to try to keep reward between +1 and -1.

**Question 2:** Discuss your agent briefly, using the following questions as a guide:

- What learning algorithm(s) did you try? What worked best for you?

- What was your final choice of hyperparameters (such as  $\alpha$ ,  $\gamma$ ,  $\epsilon$ , etc.)?
- What neural network architecture did you use (if any)? Specify layers, sizes, activation functions, etc.

**Answer:**

**1. I used the recommended DDPG algorithm. After several attempts, it can learn to hover for a short period, but still fail to keep hovering all the time.**

**2. Here's hyperparameters I chose :**

- Algorithm parameters
  - Discount factor
    - $\gamma = 0.99$
  - Soft-update
    - $\tau = 0.01$
- Replay buffer
  - size = 100000, batch\_size = 64
- OU Noise
  - $\mu = 0$ ,  $\theta = 0.15$ ,  $\sigma = 0.2$
- Learning rate (for all network)
  - $\alpha = 0.001$
- Dropout probability
  - 0.5
- L2 regularization
  - 0.01

**3. Here's my neural network architecture :**

- Actor Network
  - Input : state
  - Output : action
  - Layers :
    - Dense : size = 32, L2\_regularizer=0.01
    - BatchNormalization
    - ReLU
    - Dropout : prob = 0.5
    - Dense Layer : size = 64, L2\_regularizer=0.01
    - BatchNormalization
    - ReLU
    - Dropout : prob = 0.5
    - Dense Layer : size = 128, L2\_regularizer=0.01
    - BatchNormalization
    - ReLU
    - Dropout : prob = 0.5
    - Dense Layer : size = 32, L2\_regularizer=0.01
    - BatchNormalization
    - ReLU
    - Dropout : prob = 0.5
    - Dense : size = action\_size
- Critic Network
  - Input : state, action
  - Output : Q value
  - State Network :
    - Input : state

- Layers :
  - Dense : size = 32, L2\_regularizer=0.01
  - BatchNormalization
  - ReLU
  - Dropout : prob = 0.5
  - Dense : size = 64, L2\_regularizer=0.01
  - BatchNormalization
  - ReLU
  - Dropout : prob = 0.5
  - Dense : size = 128, L2\_regularizer=0.01
  - BatchNormalization
  - ReLU
  - Dropout : prob = 0.5
  - Dense Layer : size = 32, L2\_regularizer=0.01
  - BatchNormalization
  - ReLU
  - Dropout : prob = 0.5
  - Dense : size = 1
- Action Network
  - Input : action
  - Layers :
    - Dense : size = 32, L2\_regularizer=0.01
    - BatchNormalization
    - ReLU
    - Dropout : prob = 0.5
    - Dense Layer : size = 64, L2\_regularizer=0.01
    - BatchNormalization
    - ReLU
    - Dropout : prob = 0.5
    - Dense Layer : size = 128, L2\_regularizer=0.01
    - BatchNormalization
    - ReLU
    - Dropout : prob = 0.5
    - Dense : size = 32, L2\_regularizer=0.01
    - BatchNormalization
    - ReLU
    - Dropout : prob = 0.5
- Combination of State and Action network
  - Layers :
    - Add : add output of state network and action network
    - ReLU

**Question 3:** Using the episode rewards plot, discuss how the agent learned over time.

- Was it an easy task to learn or hard?
- Was there a gradual learning curve, or an aha moment?
- How good was the final performance of the agent? (e.g. mean rewards over the last 10 episodes)

**Answer:**

- Though it seems to be a simple task, I found it's not that easy as I take all physical factors into consideration. Intuitively, too many factors seem to make the learning algorithm hard to find a promising policy to follow because there are too many choices. Finally, I only chose pose as the main target.
- It seems gradually learning at the beginning, and enter to certain stable phase. However, it occasionally trap into poor total reward phase, but after several episodes, it exits to that local minimum region.
- The final mean rewards over the last 10 episodes is 36.6

**Question 4:** Briefly summarize your experience working on this project. You can use the following prompts for ideas.

- What was the hardest part of the project? (e.g. getting started, plotting, specifying the task, etc.)
- Did you find anything interesting in how the quadcopter or your agent behaved?

**Answer:**

- It's very hard to get started and realize how the DDPG sample code works. I had some troubles in linking the concepts learned from nanodegree videos to the codes. After figuring out how to design task, I found the most challenging part is to design reward and tweak hyperparameter to make DDPG learn a good policy. There seems to be lots of know-hows to tweak DDPG learning models.
- Though my agent fails to hover all the time (descends in z direction at the end), it's pose values in x direction and y direction keep near 0 that I want.

In [ ]: