

# TSSL Lab 3 - Nonlinear state space models and Sequential Monte Carlo

In this lab we will make use of a non-linear state space model for analyzing the dynamics of SARS-CoV-2, the virus causing covid-19. We will use an epidemiological model referred to as a Susceptible-Exposed-Infectious-Recovered (SEIR) model. It is a stochastic adaptation of the model used by the The Public Health Agency of Sweden for predicting the spread of covid-19 in the Stockholm region early in the pandemic, see [Estimates of the peak-day and the number of infected individuals during the covid-19 outbreak in the Stockholm region, Sweden February – April 2020](https://www.folkhalsomyndigheten.se/publicerat-material/publikationsarkiv/e/estimates-of-the-peak-day-and-the-number-of-infected-individuals-during-the-covid-19-outbreak-in-the-stockholm-region-sweden-february--april-2020/) (<https://www.folkhalsomyndigheten.se/publicerat-material/publikationsarkiv/e/estimates-of-the-peak-day-and-the-number-of-infected-individuals-during-the-covid-19-outbreak-in-the-stockholm-region-sweden-february--april-2020/>).

The background and details of the SEIR model that we will use are available in the document *TSSL Lab 3 Predicting Covid-19 Description of the SEIR model* on LISAM. Please read through the model description before starting on the lab assignments to get a feeling for what type of model that we will work with.

---

## DISCLAIMER

Even though we will use a type of model that is common in epidemiological studies and analyze real covid-19 data, you should *NOT* read too much into the results of the lab. The model is intentionally simplified to fit the scope of the lab, it is not validated, and it involves several model parameters that are set somewhat arbitrarily. The lab is intended to be an illustration of how we can work with nonlinear state space models and Sequential Monte Carlo methods to solve a problem of practical interest, but the actual predictions made by the final model should be taken with a big grain of salt.

---

We load a few packages that are useful for solving this lab assignment.

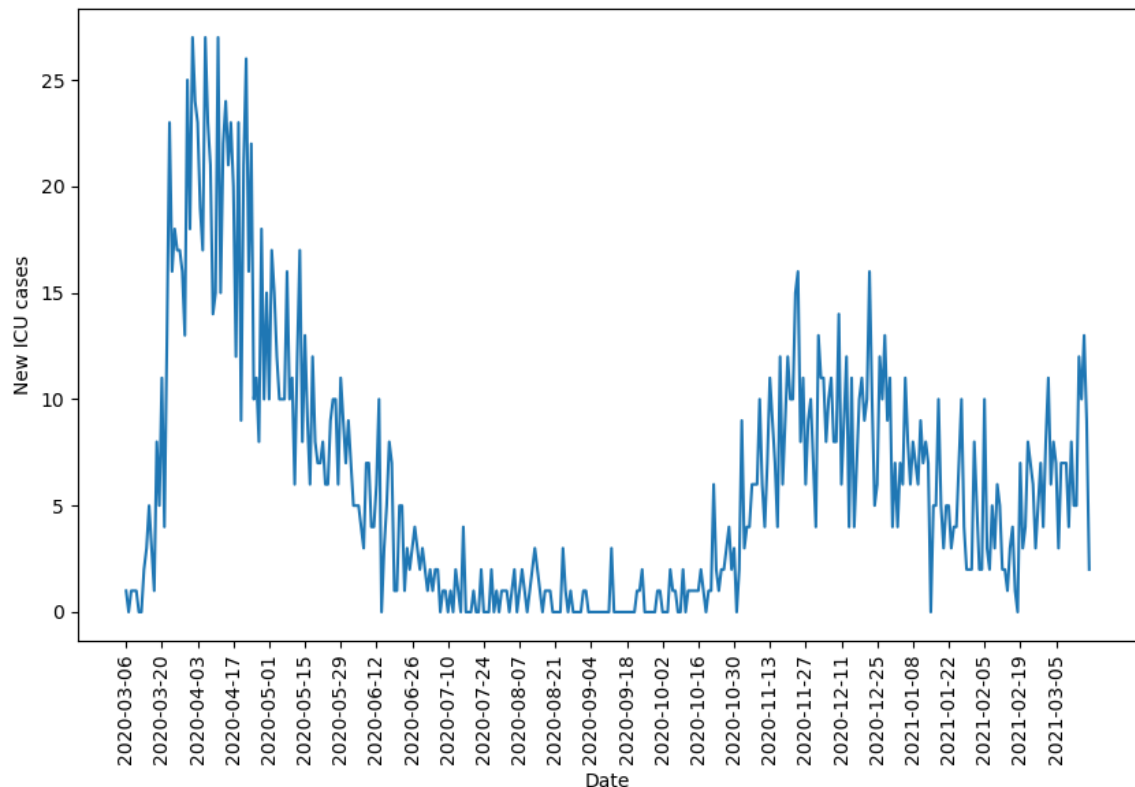
```
In [1]: import pandas # Loading data / handling data frames
import numpy as np
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (10,6) # Increase default size of plots
```

## 3.1 A first glance at the data

The data that we will use in this lab is a time series consisting of daily covid-19-related intensive care cases in Stockholm from March 2020 to March 2021. As always, we start by loading and plotting the data.

```
In [2]: data=pandas.read_csv('SIR_Stockholm.csv',header=0)
```

```
In [2]: data=pandas.read_csv('SIR_Stockholm.csv',header=0)
y_sthlm = data['ICU'].values
u_sthlm = data['Date'].values
ndata = len(y_sthlm)
plt.plot(u_sthlm,y_sthlm)
plt.xticks(range(0, ndata, 14), u_sthlm[::14], rotation = 90) # Show only one
plt.xlabel('Date')
plt.ylabel('New ICU cases')
plt.show()
```



**Q0:** What type of values can the observations  $y_t$  take? Is a Gaussian likelihood model a good choice if we want to respect the properties of the data?

**A:** Discrete data. As can be seen from the plot, the number of ICU cases is more likely to fit a complex distribution than a Gaussian distribution. Using Gaussian distribution is not a good option.

## 3.2 Setting up and simulating the SEIR model

In this section we will set up a SEIR model and use this to simulate a synthetic data set. You should keep these simulated trajectories, we will use them in the following sections.

```
In [3]: from tssltools_lab3 import Param, SEIR
```

```
In [3]: from tssltools_lab3 import Param, SEIR

        """For Stockholm the population is probably roughly 2.5 million."""
        population_size = 2500000

        """ Binomial probabilities (p_se, p_ei, p_ir, and p_ic) and the transmission
        pse = 0          # This controls the rate of spontaneous s->e transitions. It is
        pei = 1 / 5.1    # Based on FHM report
        pir = 1 / 5      # Based on FHM report
        pic = 1 / 1000   # Quite arbitrary!
        rho = 0.3        # Quite arbitrary!

        """ The instantaneous contact rate b[t] is modeled as
        b[t] = exp(z[t])
        z[t] = z[t-1] + epsilon[t], epsilon[t] ~ N(0,sigma_epsilon^2)
        """
        sigma_epsilon = .1

        """ For setting the initial state of the simulation"""
        i0 = 1000 # Mean number of infectious individuals at initial time point
        e0 = 5000 # Mean number of exposed...
        r0 = 0    # Mean number of recovered
        s0 = population_size - i0 - e0 - r0 # Mean number of susceptible
        init_mean = np.array([s0, e0, i0, 0.], dtype=np.float64) # The last 0. is the

        """All the above parameters are stored in params."""
        params = Param(pse, pei, pir, pic, rho, sigma_epsilon, init_mean, population_s

        """ Create a model instance"""
        model = SEIR(params)
```

**Q1:** Generate 10 different trajectories of length 200 from the model and plot them in one figure. Does the trajectories look reasonable? Could the data have been generated using this model?

For reproducibility, we set the seed of the random number generator to 0 before simulating the trajectories using `np.random.seed(0)`

Save these 10 generated trajectories for future use.

(hint: The SEIR class has a simulate method)

**A1:** In the real data, multiple peaks can be seen, but the model can only generate trajectories with only one peak. Can't generate the data using this model.

```
In [4]: np.random.seed(0)
```

```
In [4]: np.random.seed(0)
        help(model.simulate)
```

Help on method simulate in module tssltools\_lab3:

simulate(T, N=1) method of tssltools\_lab3.SEIR instance

Simulates the SEIR model for a given number of time steps. Multiple trajectories can be simulated simultaneously.

Parameters

-----

T : int

Number of time steps to simulate the model for.

N : int, optional

Number of independent trajectories to simulate. The default is 1.

Returns

-----

alpha : ndarray

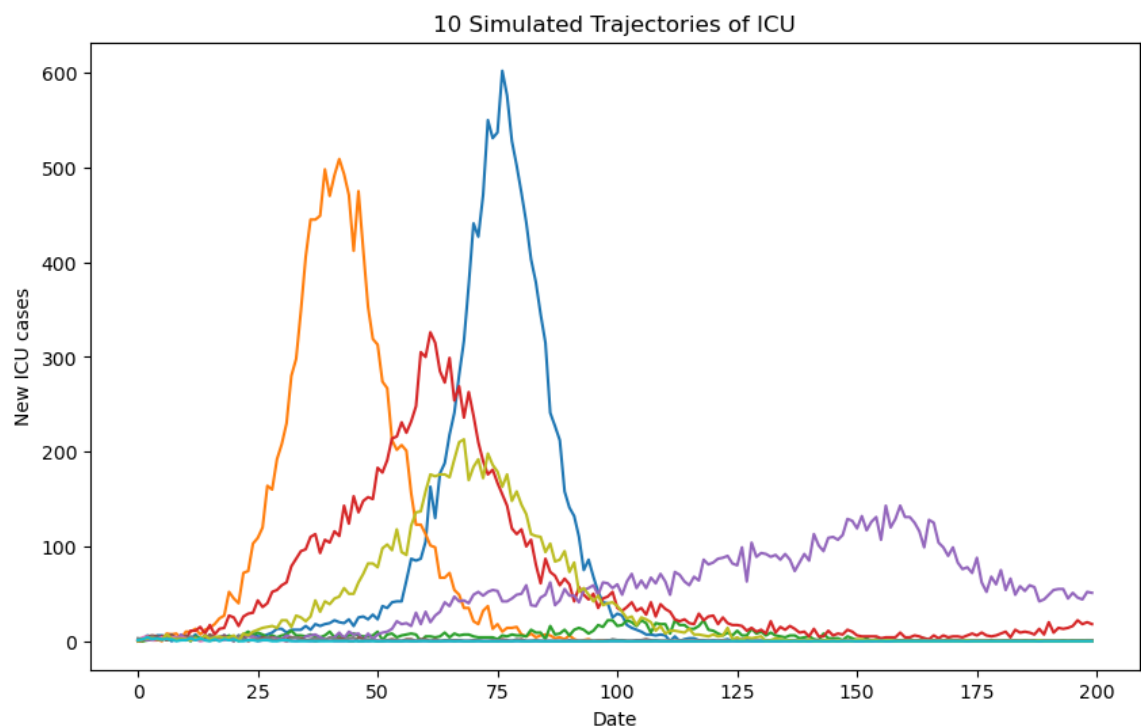
Array of size (d,N,T) with state trajectories. alpha[:,i,:] is the i:th trajectory.

y : ndarray

Array of size (1,N,T) with observations.

```
In [5]: alpha,y=model.simulate(200,N=10)

        for i in range(10):
            plt.plot(y[0,i,:])
        plt.title("10 Simulated Trajectories of ICU")
        plt.xlabel('Date')
        plt.ylabel('New ICU cases')
        plt.show()
```



### 3.3 Sequential Importance Sampling

Next, we pick out one trajectory that we will use for filtering. We use simulated data to start with, since we then know the true underlying SEIR states and can compare the filter results with the ground truth.

**Q2:** Implement the **Sequential Importance Sampling** algorithm by filling in the following functions.

The **exp\_norm** function should return the normalized weights and the log average of the unnormalized weights. For numerical reasons, when calculating the weights we should "normalize" the log-weights first by removing the maximal value.

Let  $\bar{\omega}_t = \max(\log \omega_t^i)$  and take the exponential of  $\log \tilde{\omega}_t^i = \log \omega_t^i - \bar{\omega}_t$ . Normalizing  $\tilde{\omega}_t^i$  will yield the normalized weights!

For the log average of the unnormalized weights, care has to be taken to get the correct output,  $\log(1/N \sum_{i=1}^N \tilde{\omega}_t^i) = \log(1/N \sum_{i=1}^N \omega_t^i) - \bar{\omega}_t$ . We are going to need this in the future, so best to implement it right away.

*(hint: look at the SEIR model class, it contains all necessary functions for propagation and weighting)*

$$\frac{\omega_t^i}{\Omega_t} = \frac{\tilde{\omega}_t^i}{\tilde{\Omega}_t}$$

and,

$$\log(1/N \sum_{i=1}^N \tilde{\omega}_t^i) = \log(1/N \sum_{i=1}^N \omega_t^i) - \bar{\omega}_t$$

(From PP Exercise 3)

```
In [6]: from tssltools_lab3 import smc_res
```

```

In [6]: from tssltools_lab3 import smc_res

def exp_norm(logwgt):
    """
    Exponentiates and normalizes the log-weights.

    Parameters
    -----
    logwgt : ndarray
        Array of size (N,) with log-weights.

    Returns
    -----
    wgt : ndarray
        Array of size (N,) with normalized weights, wgt[i] = exp(logwgt[i])/sum
        but computed in a /numerically robust way/!
    logZ : float
        log of the normalizing constant, logZ = log(sum(exp(logwgt))),
        but computed in a /numerically robust way/!
    """
    max_logwgt = max(logwgt)
    nor_logwgt = logwgt - max_logwgt
    wgt = np.exp(nor_logwgt) / sum(np.exp(nor_logwgt))

    logZ = np.log(np.mean(np.exp(nor_logwgt))) + max_logwgt

    return wgt, logZ

def ESS(wgt):
    """
    Computes the effective sample size.

    Parameters
    -----
    wgt : ndarray
        Array of size (N,) with normalized importance weights.

    Returns
    -----
    ess : float
        Effective sample size.
    """
    ess = (sum(wgt)**2) / sum(wgt**2)

    return ess

def sis_filter(model, y, N):
    d = model.d
    n = len(y)

    # Allocate memory
    particles = np.zeros((d, N, n), dtype = float) # All generated particles
    logW = np.zeros((1, N, n)) # Unnormalized Log-weight
    W = np.zeros((1, N, n)) # Normalized weight
    alpha_filt = np.zeros((d, 1, n)) # Store filter mean
    N_eff = np.zeros(n) # Efficient number of particles
    logZ = 0. # Log-Likelihood estimate

    # Filter Loop
    for t in range(n):
        # Sample from "bootstrap proposal"
        if t == 0:
            particles[:, :, 0] = model.sample_state(N=N) # Initialize from p(c
            logW[0, :, 0] = model.log_lik(v[0], particles[:, :, 0]) # Compute

```

```

if t==0:
    particles[:, :, 0] = model.sample_state(N=N) # Initialize from p(c
    logW[0, :, 0] = model.log_lik(y[t], particles[:, :, 0]) # Compute
else:
    particles[:, :, t] = model.sample_state(alpha0=particles[:, :, t-1]
    logW[0, :, t] = model.log_lik(y[t], particles[:, :, t])+logW[0, :,

# Normalize the importance weights and compute N_eff
W[0, :, t], _ = exp_norm(logW[0, :, t])
N_eff[t] = ESS(W[0,:,t])

# Compute filter estimates
alpha_filt[:, 0, t] = np.dot(W[0, :, t],particles[:, :, t].T)

return smc_res(alpha_filt, particles, W, logW=logW, N_eff=N_eff)

```

**Q3:** Choose one of the simulated trajectories and run the SIS algorithm using  $N = 100$  particles. Show plots comparing the filter means from the SIS algorithm with the underlying truth of the Infected, Exposed and Recovered.

Also show a plot of how the ESS behaves over the run.

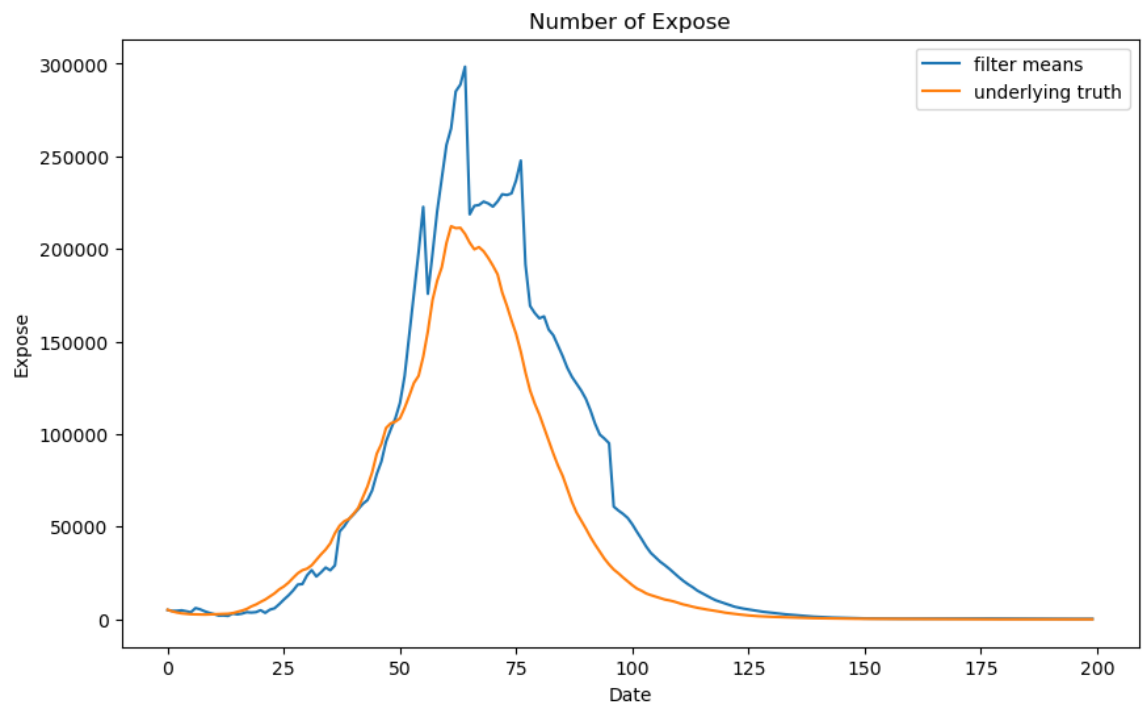
(hint: In the model we use the  $S, E, I$  as states, but  $S$  will be much larger than the others. To calculate  $R$ , note that  $S + E + I + R = \text{Population}$ )

In [7]: chosen\_trajectorie=8



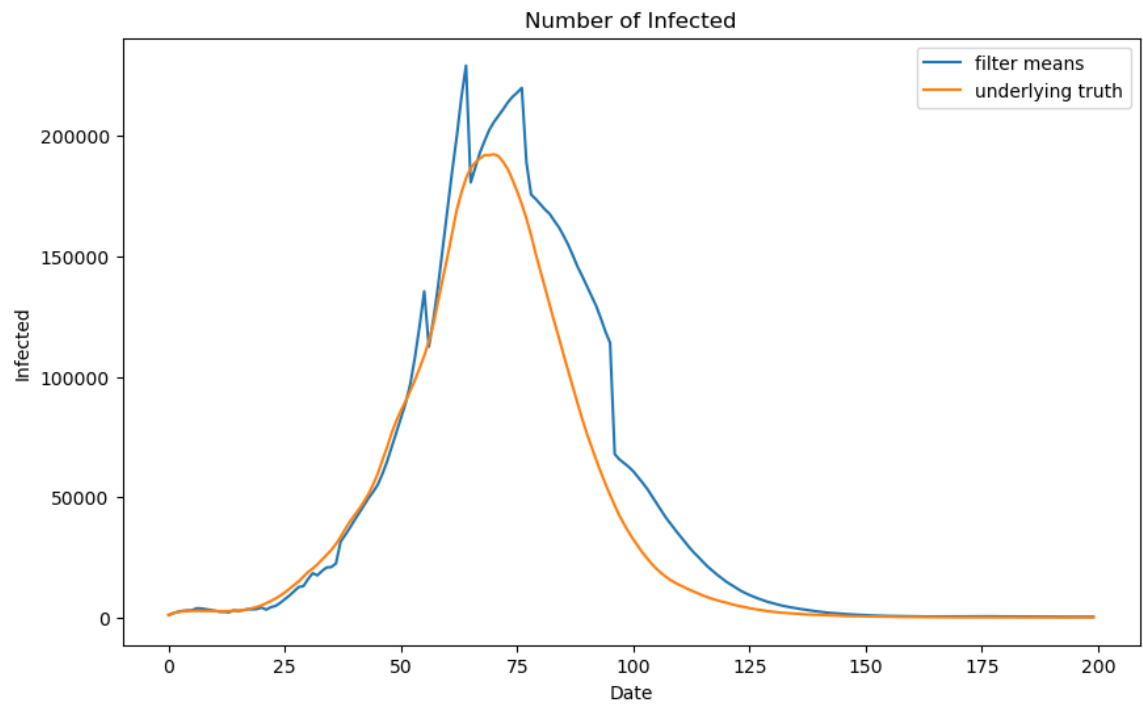
```
In [7]: chosen_trajectorie=8
y1=y[0,chosen_trajectorie,:]
alpha1=alpha[:,chosen_trajectorie,:]
smc_res1=sis_filter(model, y=y1, N=100)

# Plot Expose
index=1
plt.plot(smc_res1.alpha_filt[index,0,:],label='filter means')
plt.plot(alpha1[index,:],label='underlying truth')
plt.title("Number of Expose")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Expose')
plt.show()
```



```
In [8]: # Plot Infected
```

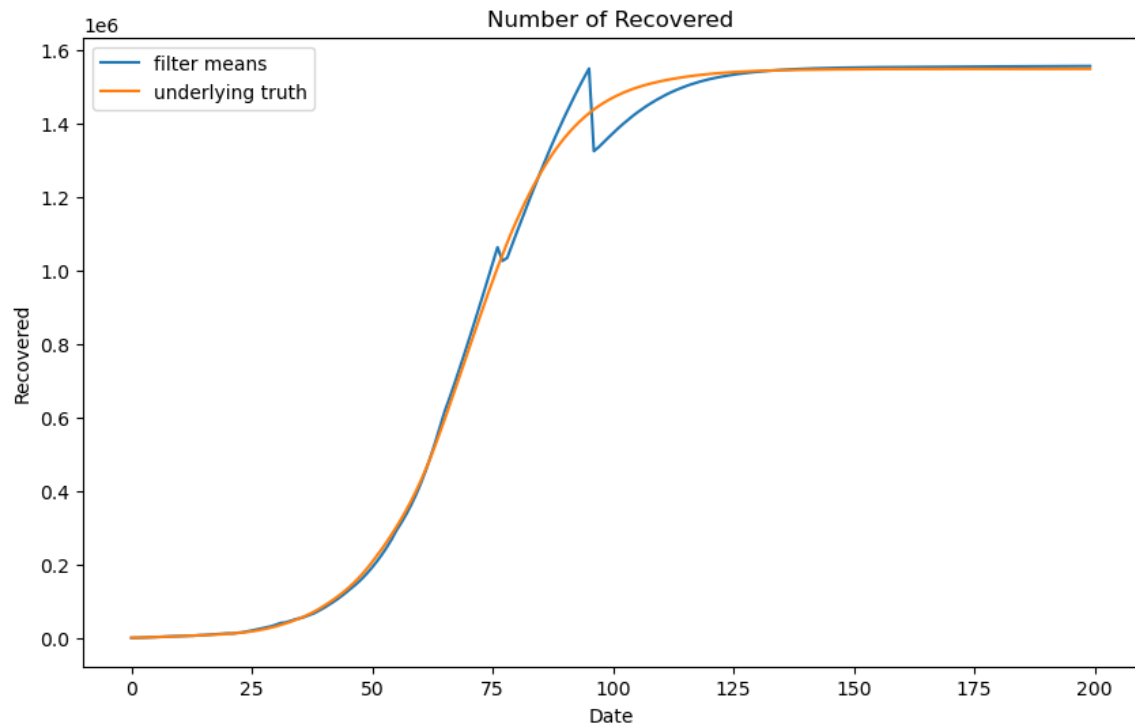
```
In [8]: # Plot Infected
index=2
plt.plot(smc_res1.alpha_filt[index,0,:],label='filter means')
plt.plot(alpha1[index,:],label='underlying truth')
plt.title("Number of Infected")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Infected')
plt.show()
```



```
In [9]: # Plot Recovered
```

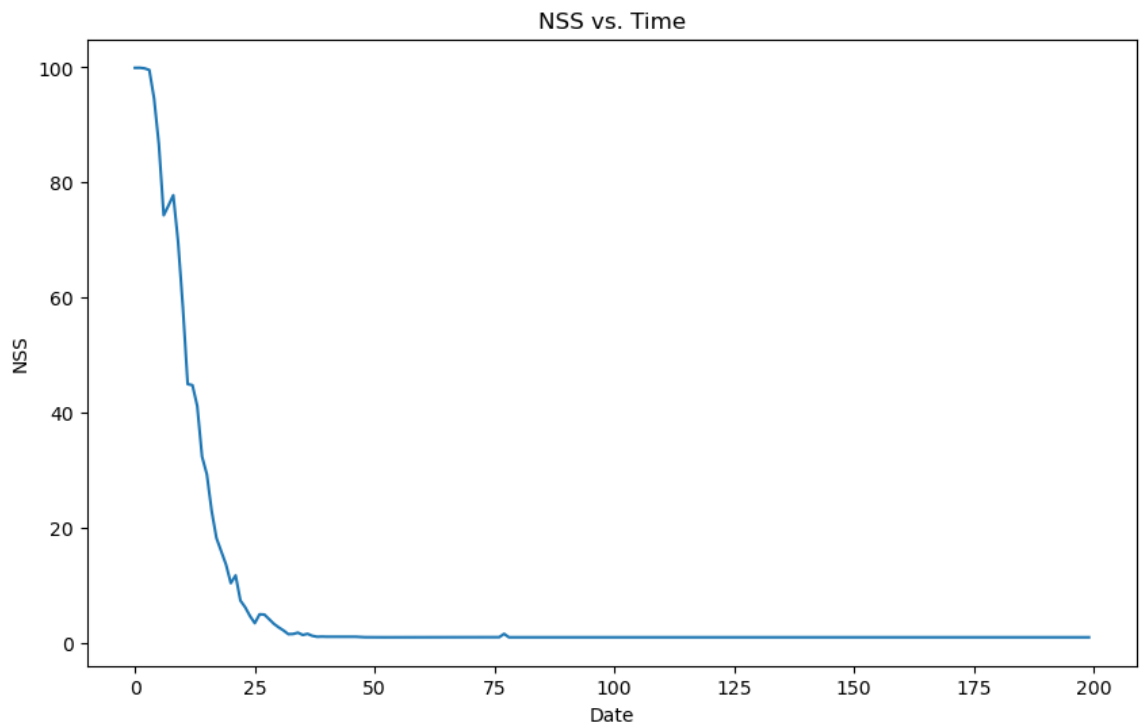
In [9]: # Plot Recovered

```
plt.plot(2500000-np.sum(smc_res1.alpha_filt[0:3,0,:],axis=0),label='filter means')
plt.plot(2500000-np.sum(alpha1[0:3,:],axis=0),label='underlying truth')
plt.title("Number of Recovered")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Recovered')
plt.show()
```



In [10]: plt.plot(smc\_res1.N\_eff)

```
In [10]: plt.plot(smc_res1.N_eff)
plt.title("NSS vs. Time")
#plt.legend()
plt.xlabel('Date')
plt.ylabel('NSS')
plt.show()
```



### 3.4 Sequential Importance Sampling with Resampling

Pick the same simulated trajectory as for the previous section.

**Q4:** Implement the **Sequential Importance Sampling with Resampling** or **Bootstrap Particle Filter** by completing the code below.

$$l(y_{1:n}) = \sum_{t=1}^n \log\left(\frac{1}{N} \sum_{i=1}^N \omega_t^i\right)$$

(From PP Exercise 3)

```
In [11]: def bpf(model, y, numParticles):
```

```

In [11]: def bpf(model, y, numParticles):
    d = model.d
    n = len(y)
    N = numParticles

    # Allocate memory
    particles = np.zeros((d, N, n), dtype = float) # All generated particles
    logW = np.zeros((1, N, n)) # Unnormalized log-weight
    W = np.zeros((1, N, n)) # Normalized weight
    alpha_filt = np.zeros((d, 1, n)) # Store filter mean
    N_eff = np.zeros(n) # Efficient number of particles
    logZ = 0. # Log-Likelihood estimate

    # Filter Loop
    for t in range(n):
        # Sample from "bootstrap proposal"
        if t == 0: # Initialize from prior
            particles[:, :, 0] = model.sample_state(N=N)
        else: # Resample and propagate according to dynamics
            ind = np.random.choice(N, N, replace=True, p=W[0, :, t-1])
            resampled_particles = particles[:, ind, t-1]
            particles[:, :, t] = model.sample_state(alpha0=resampled_particles

        # Compute weights
        logW[0, :, t] = model.log_lik(y[t], particles[:, :, t]) #In bpf don't
        W[0, :, t], logZ_now = exp_norm(logW[0, :, t])
        logZ += logZ_now # Update Log-Likelihood estimate
        N_eff[t] = ESS(W[0, :, t])

        # Compute filter estimates
        alpha_filt[:, 0, t] = np.dot(W[0, :, t], particles[:, :, t].T)

    return smc_res(alpha_filt, particles, W, N_eff = N_eff, logZ = logZ)

```

**Q5:** Use the same simulated trajectory as above and run the BPF algorithm using  $N = 100$  particles. Show plots comparing the filter means from the Bootstrap Particle Filter algorithm with the underlying truth of the Infected, Exposed and Recovered. Also show a plot of how the ESS behaves over the run. Compare this with the results from the SIS algorithm.

**A:** Compare to the results from the SIS algorithm, the ESS in the BPF will not converge to zero as the algorithm resample each steps.

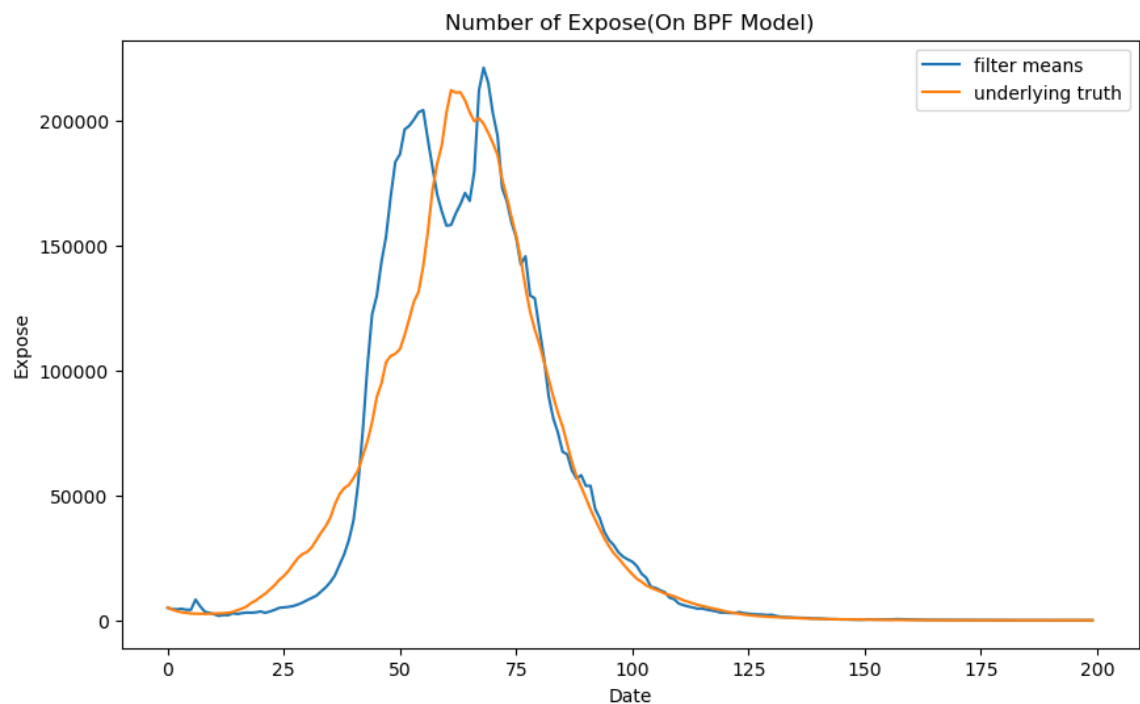
```

In [12]: chosen_trajectory=8

```

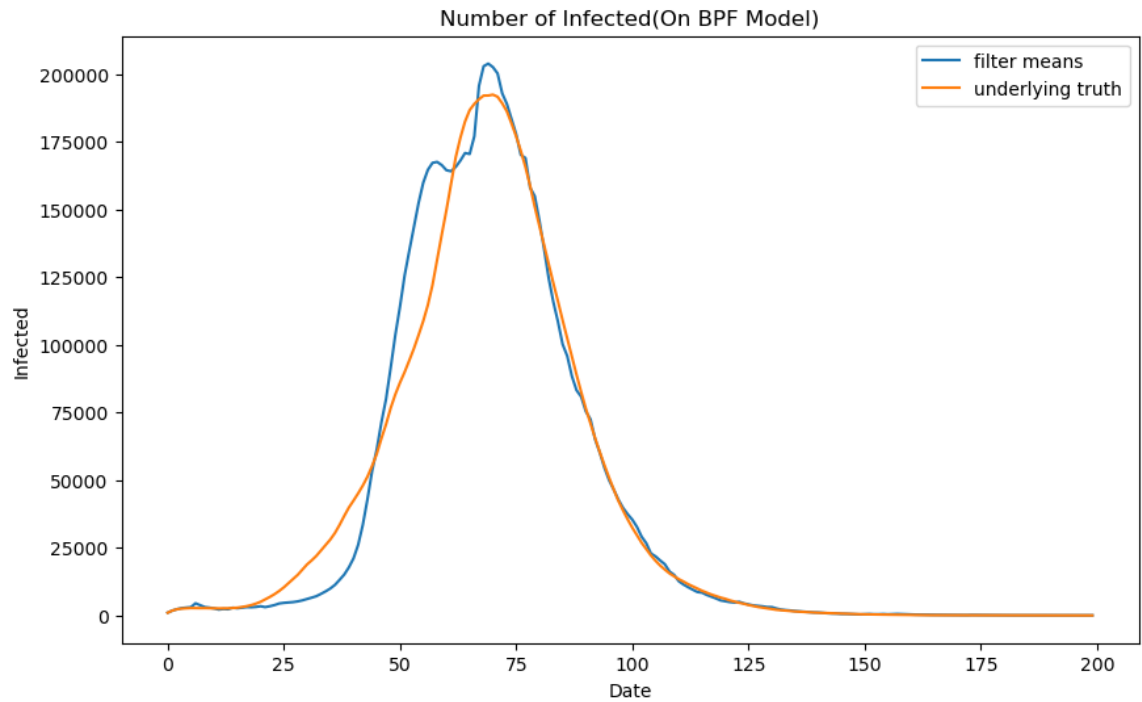
```
In [12]: chosen_trajectory=8
y2=y[0,chosen_trajectory,:]
alpha2=alpha[:,chosen_trajectory,:]
smc_res2=bpf(model, y=y2, numParticles=100)

# Plot Expose
index=1
plt.plot(smc_res2.alpha_filt[index,0,:],label='filter means')
plt.plot(alpha2[index,:],label='underlying truth')
plt.title("Number of Expose(On BPF Model)")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Expose')
plt.show()
```



```
In [13]: # Plot Infected
```

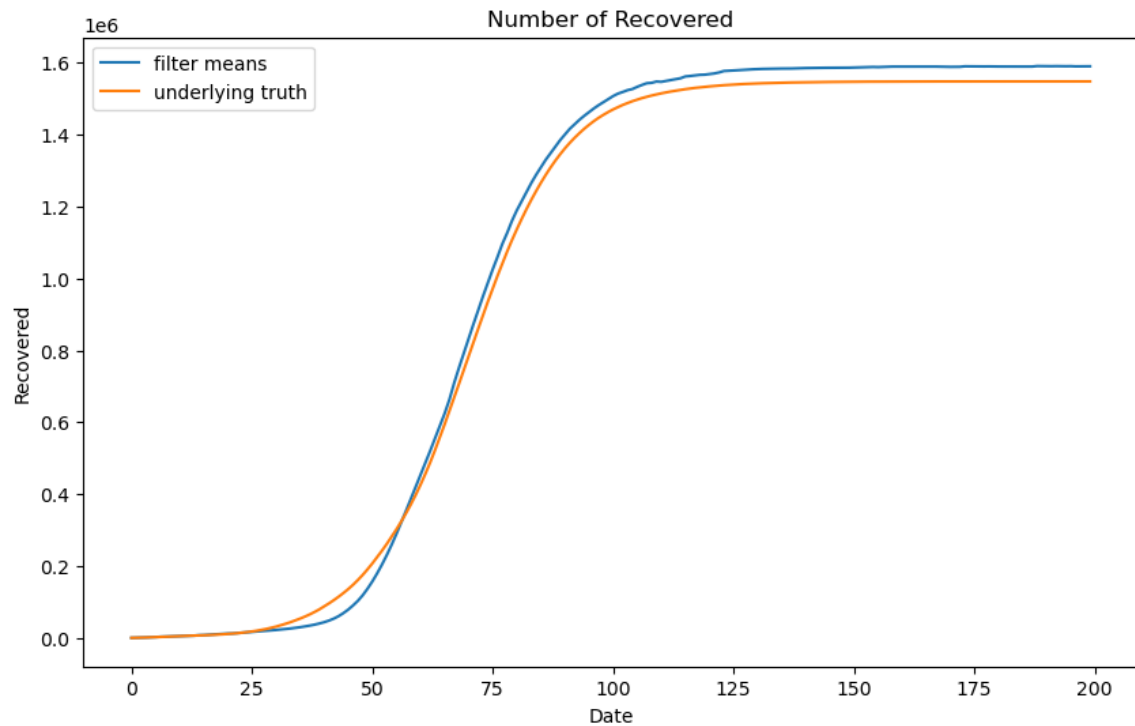
```
In [13]: # Plot Infected
index=2
plt.plot(smc_res2.alpha_filt[index,0,:],label='filter means')
plt.plot(alpha2[index,:],label='underlying truth')
plt.title("Number of Infected(On BPF Model)")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Infected')
plt.show()
```



```
In [14]: # Plot Recovered
```

In [14]: # Plot Recovered

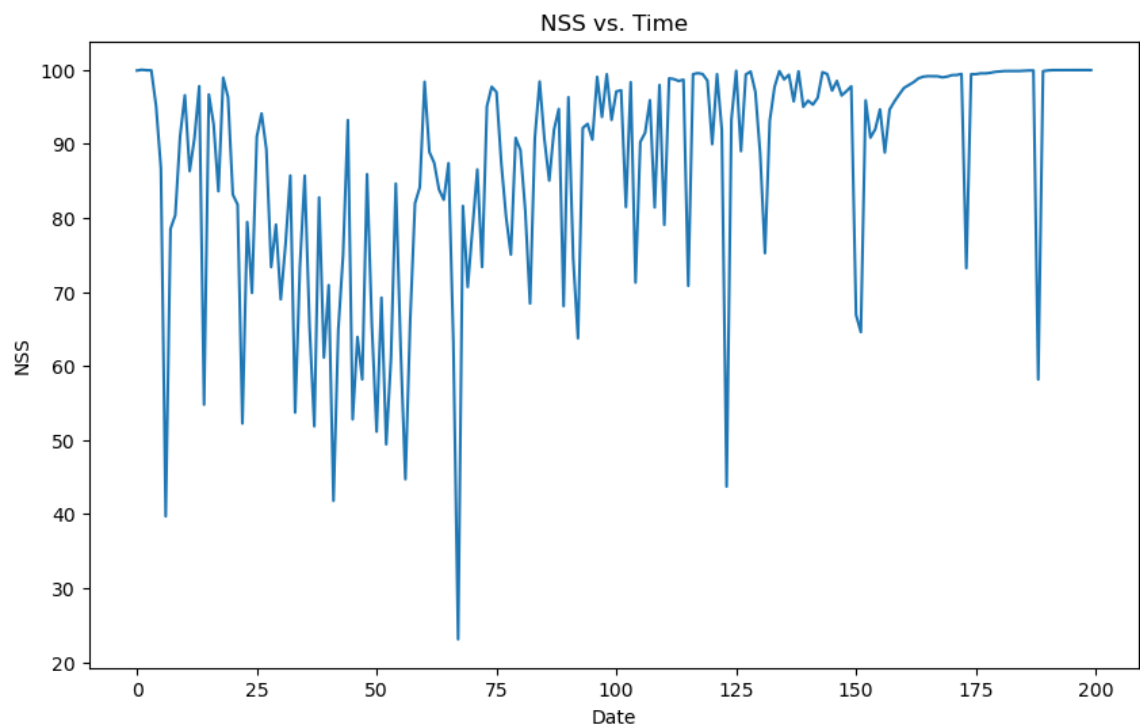
```
plt.plot(2500000-np.sum(smc_res2.alpha_filt[0:3,0,:],axis=0),label='filter means')
plt.plot(2500000-np.sum(alpha2[0:3,:],axis=0),label='underlying truth')
plt.title("Number of Recovered")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Recovered')
plt.show()
```



In [15]: plt.plot(smc\_res2.N\_eff)



```
In [15]: plt.plot(smc_res2.N_eff)
plt.title("NSS vs. Time")
#plt.legend()
plt.xlabel('Date')
plt.ylabel('NSS')
plt.show()
```



### 3.5 Estimating the data likelihood and learning a model parameter

In this section we consider the real data and learning the model using this data. For simplicity we will only look at the problem of estimating the  $\rho$  parameter and assume that others are fixed.

You are more than welcome to also study the other parameters.

Before we begin to tweak the parameters we run the particle filter using the current parameter values to get a benchmark on the log-likelihood.

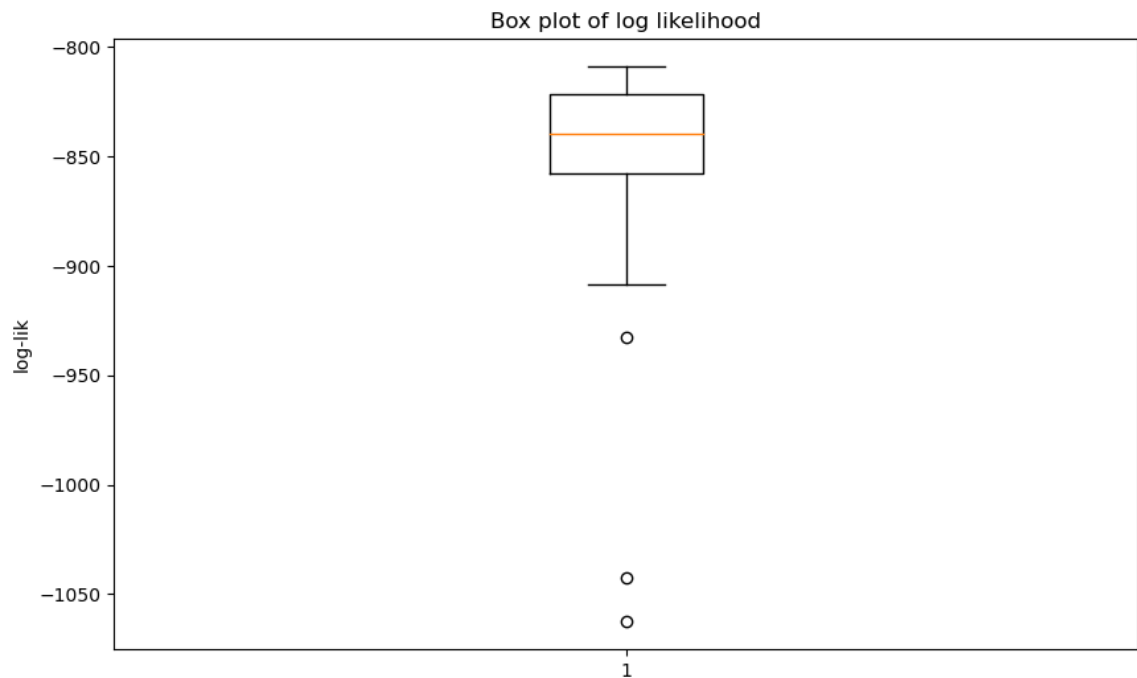
**Q6:** Run the bootstrap particle filter using  $N = 200$  particles on the real dataset and calculate the log-likelihood. Rerun the algorithm 20 times and show a box-plot of the log-likelihood.

```
In [16]: loglik=np.zeros(20)
```

```
In [16]: loglik=np.zeros(20)

#chosen_trajectorie=8
#y3=y[0,chosen_trajectorie,:]
y3=y_sthlm
for i in range(20):
    smc_res3=bpf(model, y=y3, numParticles=200)
    loglik[i]=smc_res3.logZ

plt.boxplot(loglik)
plt.title("Box plot of log likelihood")
plt.ylabel("log-lik")
plt.show()
```



**Q7:** Make a grid of the  $\rho$  parameter in the interval  $[0.1, 0.9]$ . Use the bootstrap particle filter to calculate the log-likelihood for each value. Run the bootstrap particle filter using  $N = 200$  multiple times (at least 20) per value and use the average as your estimate of the log-likelihood. Plot the log-likelihood function and mark the maximal value.

(hint: use `np.linspace` to create a grid of parameter values)

```
In [32]: # Choose sample size
```

```

In [32]: # Choose sample size
rho_candidate_num=100

rho_candidate=np.linspace(0.1, 0.9, num=rho_candidate_num, endpoint=True, rets
average_loglik=np.zeros(rho_candidate_num)

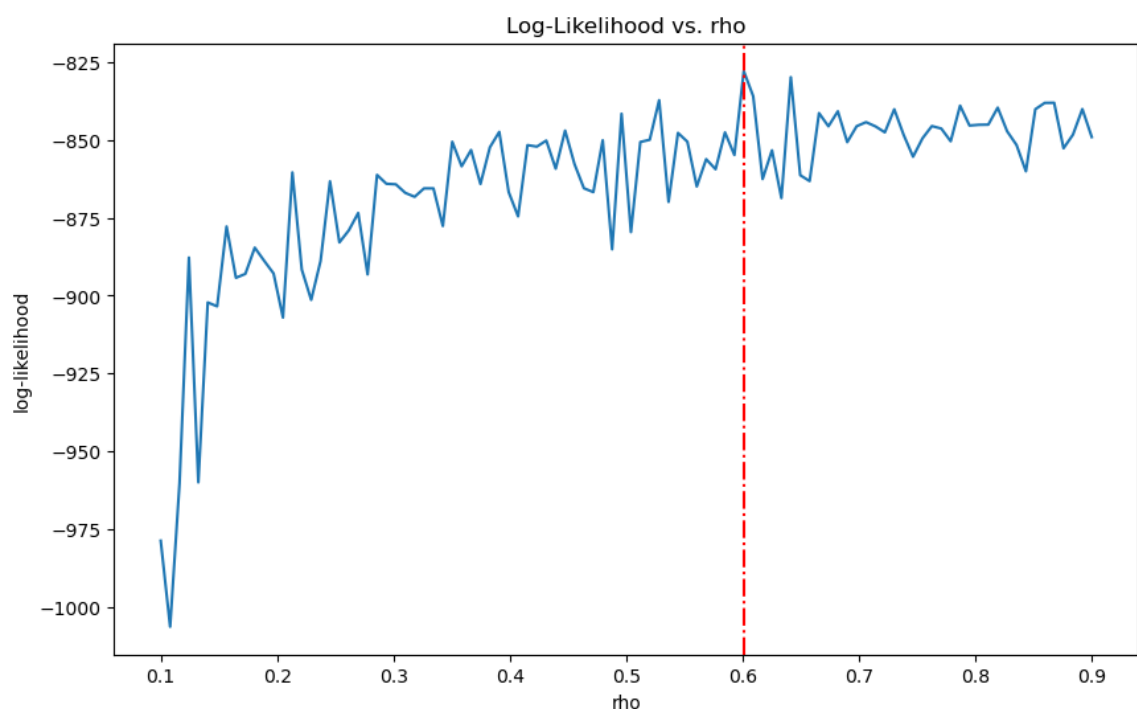
for i in range(rho_candidate_num):
    params4 = Param(pse, pei, pir, pic, rho_candidate[i], sigma_epsilon, init_
    model4 = SEIR(params4)
    loglik=np.zeros(20)
    for j in range(20):
        smc_res4=bpf(model4, y=y_sthlm, numParticles=200)
        loglik[j]=smc_res4.logZ
    average_loglik[i]=np.mean(loglik)

print("The best rho in the candidates is:",rho_candidate[np.argmax(average_log

plt.plot(rho_candidate,average_loglik)
plt.axvline(rho_candidate[np.argmax(average_loglik)],linestyle='-.',color='red
plt.title("Log-Likelihood vs. rho")
#plt.legend()
plt.xlabel('rho')
plt.ylabel('log-likelihood')
plt.show()

```

The best rho in the candidates is: 0.601010101010101



**Q8:** Run the bootstrap particle filter on the full dataset with the optimal  $\rho$  value. Present a plot of the estimated Infected, Exposed and Recovered states.

```

In [33]: params5 = Param(pse, pei, pir, pic, rho_candidate[np.argmax(average_loglik)],
model5 = SEIR(params5)

smc_res5=bpf(model5, y=y_sthlm, numParticles=200)

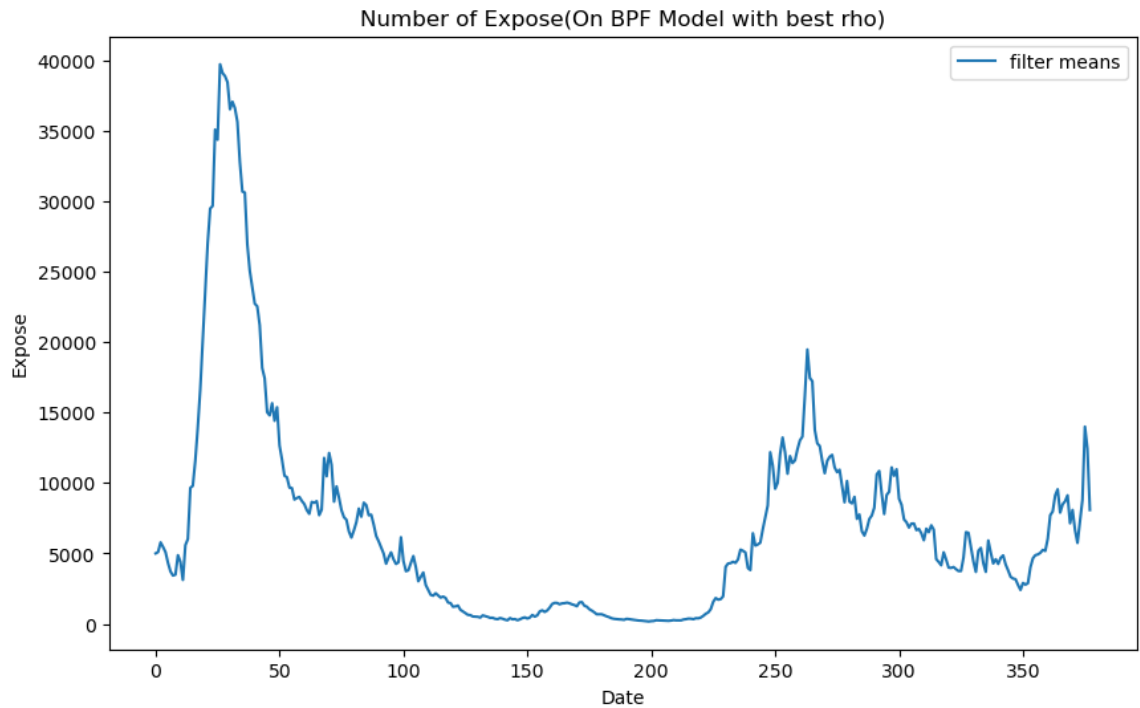
```

```

In [34]: # Plot Expose

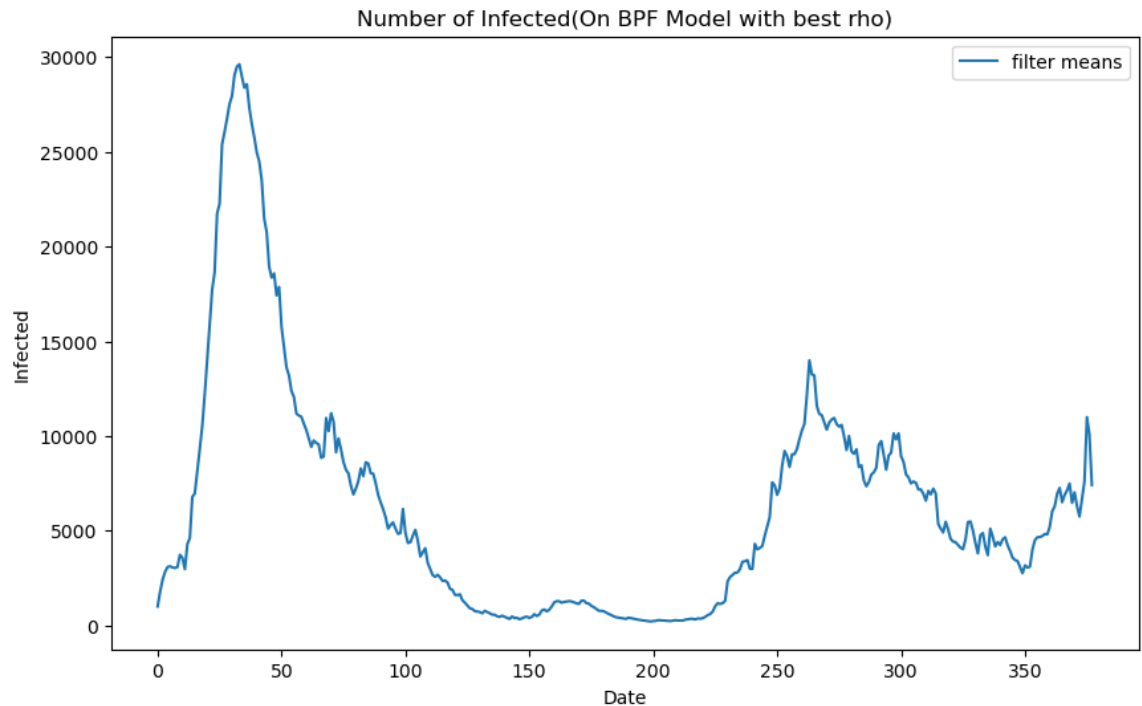
```

```
In [34]: # Plot Expose
index=1
plt.plot(smc_res5.alpha_filt[index,0,:],label='filter means')
plt.title("Number of Expose(On BPF Model with best rho)")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Expose')
plt.show()
```

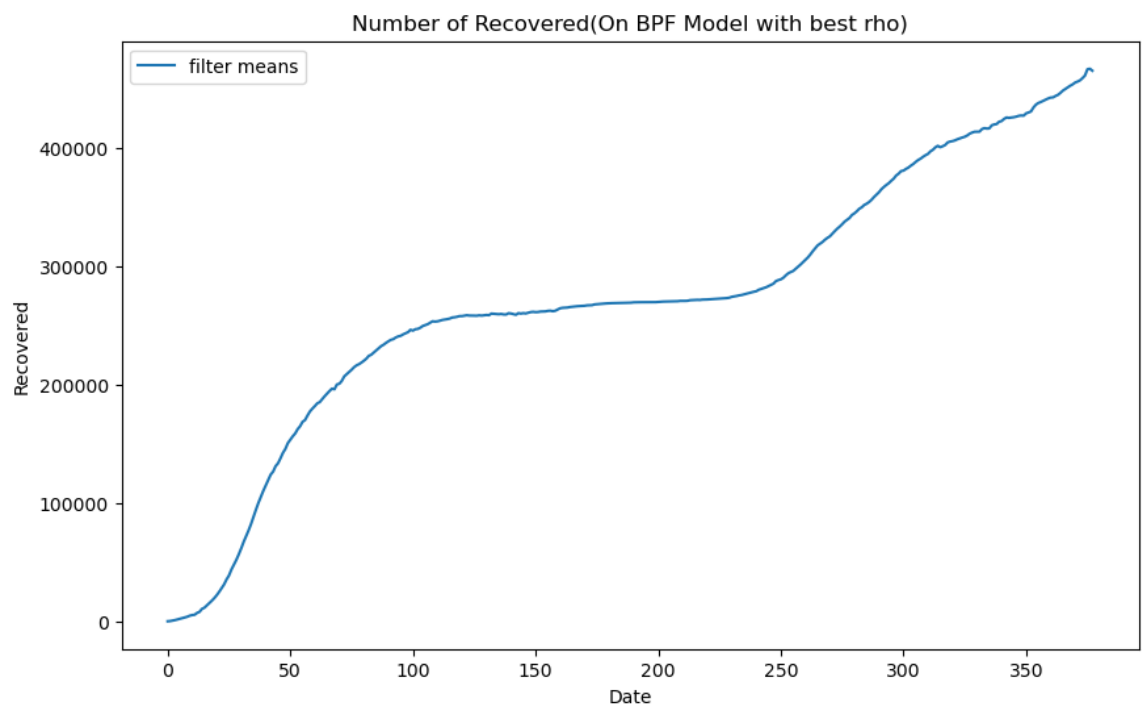


```
In [35]: # Plot Infected
```

```
In [35]: # Plot Infected
index=2
plt.plot(smc_res5.alpha_filt[index,0,:],label='filter means')
plt.title("Number of Infected(On BPF Model with best rho)")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Infected')
plt.show()
```



```
In [36]: # Plot Recovered
plt.plot(2500000-np.sum(smc_res5.alpha_filt[0:3,0,:],axis=0),label='filter means')
plt.title("Number of Recovered(On BPF Model with best rho)")
plt.legend()
plt.xlabel('Date')
plt.ylabel('Recovered')
plt.show()
```



```
In [ ]:
```

In [ ]:

Date