

MCMC Example

Davide G.
in collaboration with UCL HR



AUCIL

Program

- o Data vector
- o Bayesian Statistic
- Synthetic Data + REAL world
 measurement
- o MCMC implementation (emcee)



Dala Vector

âUCI

How much stressed is a PhD. student as a function of time and weather?

```
Model
Parameters:
```

- 1. Thesis anxiety
- 2. What am I doing?
- 3. Pub attendance
- 4. Finish the paper!
- 5. Sport activities

```
12
    variables of the data vector function:
    time from 0 to 36 months,
13
    weather conditions
14
    from - 0.3 (rain/cold/pollution --> Mordor ) to + 0.3 (sun --> Shire)
15
    .....
16
17
    18
19
    define data vector function,
20
    theoretical model
    21
    def phd_stress_function(time, weather, theta_ar, doing_sport):
22
23
        """ different factors
24
25
        thesis_anxiety = theta_ar[0] * np.exp(time / 36.0)
        what am i doing = theta ar[1] * np.exp(- time / 36.0)
26
        pub_attendance = theta_ar[2] * ((time - 18.0) / 36.0) ** 2.0 + weather
27
        finish_the_paper = theta_ar[3] * (np.sin((5.0 * 2.0) * np.pi * time / 36.0) + 1.0)
28
29
        if doing sport == True:
30
           sport_activities = theta_ar[4] * (- time / 36.0) + 1.0 + weather
31
           return thesis_anxiety + what_am_i_doing - pub_attendance \
32
                  + finish_the_paper - sport_activities
33
34
        return thesis_anxiety + what_am_i_doing - pub_attendance + finish_the_paper
35
```

â II CL

Unavoidable Bayesian Statistic Stide

- o Prior: previous knowledge / degree of belief
- Likelihood: is the probability of the evidence given the parameters (wikipedia)
- Posterior: is the probability of the parameters given the evidence (again wikipedia)

Likelihood
$$P(c \mid x) = \frac{P(x \mid c)P(c)}{P(x)}$$
Posterior Probability

Predictor Prior Probability

Predictor Prior Probability

Pythonically speaking:

```
66
                                                                   define prior, likelihood, posterior functions
                                                              67
   41
                                                                  def lnprob_f( x, meas, params, icov):
                                                              68
42
   def lnprior(x,ranges):
                                                              69
43
                                                              70
                                                                      time ar, weather, theta ranges, doing sport = params
      par list = x
44
                                                              71
45
46
      for i in range(len(par_list)):
                                                                      lp = lnprior(x,theta_ranges)
                                                              72
         if (par list[i] < ranges[i,0]) or (par list[i] > ranges[i,1]):
47
                                                              73
             return -np.inf
48
                                                              74
                                                                      if not np.isfinite(lp):
49
                                                                          return - np.inf
                                                              75
50
      return 0.0
                                                              76
51
   lnprob = lp + lnlike(x, meas, params, icov)
52
                                                              77
   53
                                                              78
54
   def lnlike(x, meas, params, inv_cov):
                                                                      # Check for Inprob returning NaN.
                                                              79
55
                                                                      if np.any(np.isnan(lnprob)):
                                                              80
56
      par_list = x
                                                                          # Find indexes of Inprob array with NaN values.
                                                              81
57
                                                                          indxs_of_bad_p = np.where(np.isnan(lnprob) == True)[0]
58
      time_ar, weather, theta_ranges, doing_sport = params
                                                              82
59
                                                                          # Print some debugging stuff.
                                                              83
      theta ar = np.array(par_list)
60
                                                              84
                                                                          print("NaN value of lnprob for parameters: ")
             = phd stress function(time ar, weather, theta ar, doing sport)
                                                                          print(x[indxs_of_bad_p])
                                                              85
      diff
             = model - meas
62
                                                                          # Finally raise exception.
                                                              86
63
                                                                          raise ValueError("Inprob returned NaN.")
                                                              87
      return - np.dot(diff, np.dot(inv_cov, diff)) / 2.0
64
                                                              88
```

89

90

return Inprob

Main Code: synthetic data + REAL Life PhD. student opinion

101

```
spoiler alert: true values parameters """
102
     true_theta_ar = np.array([1.0, 1.8, 2.7, 0.2, 0.6])
103
104
     doing_sport = False
105
     """ when we asked the PhD. students their opinions and what was the weather like (coordinates)
106
     time_ar = np.linspace(1.0,36.0,72)
107
     weather = np.cos(6.0 * np.pi * time_ar / 36.0 + np.pi/3.0) * 0.3
108
109
     """ generate N independent / uncorrelated PhD. students opinions (SYNTHETIC DATA !!)"""
110
     N_phd_students
111
                           = 200
     true stress
                     = phd_stress_function(time_ar, weather, true_theta_ar, doing_sport)
112
                      = np.ones(true_stress.size) * 0.1
     std normal dist
113
     N_survey_measurements = np.random.normal(true_stress, std_normal_dist,(N_phd_students,true_stress.size ))
114
115
     """ covariance matrix derived from synthetic data """
116
                           = np.cov(N_survey_measurements.T)
117
     cov matrix
118
         special PhD. student (MEASUREMENT ON REAL DATA)"""
119
      selected_phd_opinion = np.random.normal(true_stress, std_normal_dist * 0.1, (true_stress.size ))
120
```

MCMC set up:

ÎUCL

- o how many walkers (aka hobbits)
- o burn in steps (training for the long hike to Mordor)
- o Steps (to Mordor)
- Initial position of the walkers (Hobbiton, Shire, NW13 OLA)
- Sampler: magically tells the walkers where to go in the parameter space (aka Gandalf)

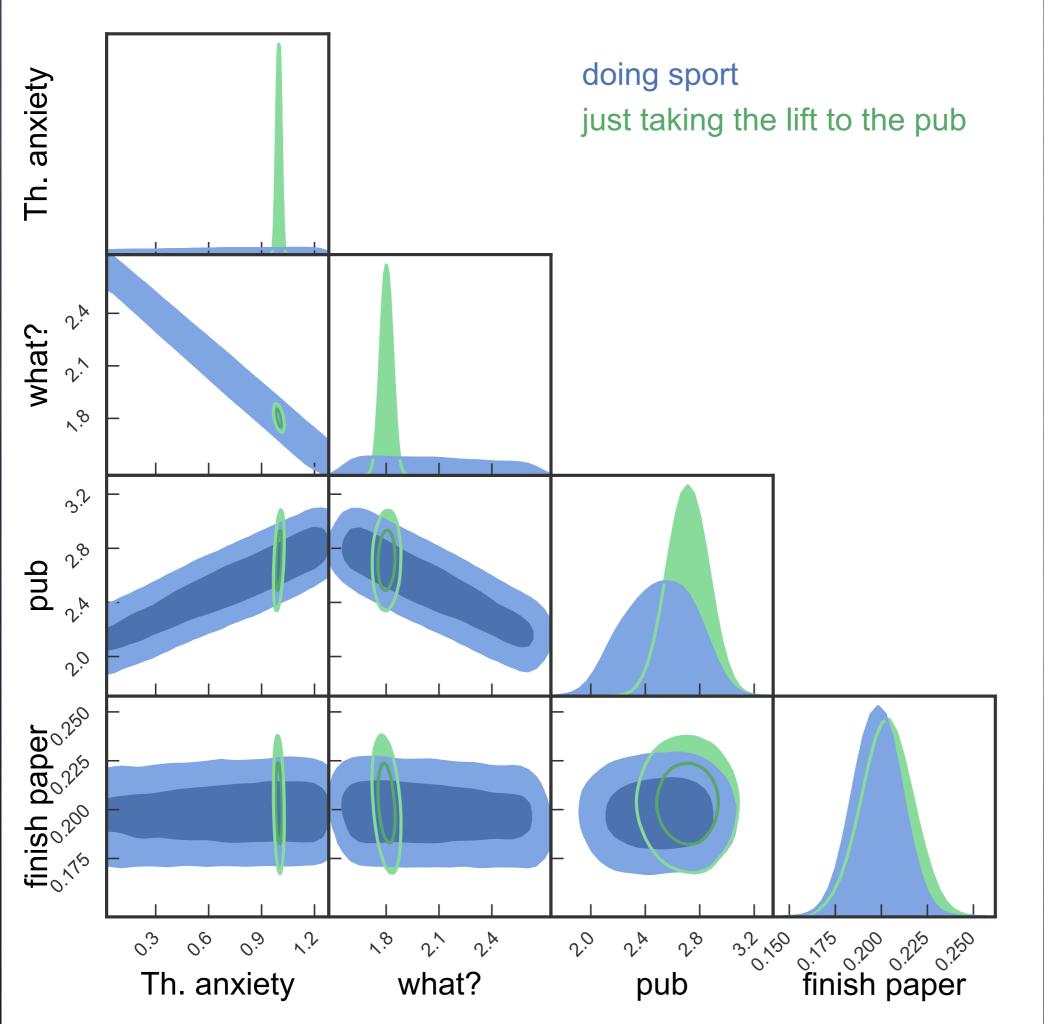
Result: CHAIN (it has kept track of the hobbits wandering around)

Again in python (1):

```
""" walkers """
136
     n hobbits
137
                      = 200
138
     """ burn in """
139
     warm_up_steps = 200
140
141
     """ steps tracked in the chain """
142
     steps to mordor = 1000
143
144
     """ run in parallel using MPI """
145
     mpi_or_not
                     = False
146
147
     """ initial guess for the maximum likelihood parameters """
148
      shire_perimeter = np.array([[0.8,1.2], [1.5,2.5], [2.0,3.5], [0.0,0.6], [0.3,1.0]])
149
     if doing sport == False:
150
          shire perimeter = np.array([[0.8,1.2], [1.5,2.5], [2.0,3.5], [0.0,0.6]])
151
152
     """ starting point of walkers based on the initial guess """
153
154
      hobbiton
                      = np.random.rand(par_num * n_hobbits).reshape((par_num,n_hobbits))
155
156
     for i in range(par_num):
          hobbiton[i] = shire_perimeter[i,0] + (shire_perimeter[i,1] - shire_perimeter[i,0]) * hobbiton[i]
157
158
                      = hobbiton.T
159
      hobbiton
160
     """ additional parameters to pass to the model function, including coordinates """
161
                      = time_ar, weather, middle_earth, doing_sport
162
      params
```

Again in python (2): ±UCL

```
if mpi_or_not == False:
167
168
          """ initialisate the MAGIC sampler """
169
          gandalf = emcee.EnsembleSampler(n_hobbits, par_num, lnprob_f,
170
                                              args=[selected_phd_opinion, params, np.linalg.inv(cov_matrix)], threads = 1)
171
172
          "burn in"
173
          pos, prob, state = gandalf.run_mcmc(hobbiton, warm_up_steps)
174
          gandalf.reset()
175
176
          "run mcmc"
177
          gandalf.run_mcmc(pos, steps_to_mordor)
178
179
          "reshape"
180
181
          samples = gandalf.chain[:, :, :].reshape((-1, par_num))
182
          "intervals"
183
          mcmc_list = []
184
          mcmc_list = map(lambda v: (v[1], v[2] - v[1], v[1] - v[0]),
185
                          zip(*np.percentile(samples, [16, 50, 84],
186
                                              axis=0)))
187
```



1-2D
marginalised
posterior
distribution plot

References:

- https://github.com/ davidegua/ mcmc_example
- http://dfm.io/ emcee/current/

Gualdi et al. in prep. (2018)

