
1RT705: Project Assignment

Advanced Probabilistic Machine Learning 2020

Paul Häusner
Nathaniel Ahy
Patrick Petersson

Abstract

1 The goal of this project is to implement the by Microsoft developed Trueskill model
2 in order to rank teams or players based on results of the games they played against
3 each other. In the implementation two different approaches are tested out. The first
4 one is based on Bayesian networks and uses Gibbs sampling in order to estimate
5 the posterior distribution. The second approach uses factor graphs and is based on
6 their belief propagation with moment-matching in order to estimate non-Gaussian
7 distributions.

1 Q.1 Modeling

9 The four random variables in our model are shown with their distribution in equation (1) to (3). The
10 Bayesian model itself is given in equation 4.

$$s_1 \sim \mathcal{N}(s_1; \mu_1, \sigma_1^2), s_2 \sim \mathcal{N}(s_2; \mu_2, \sigma_2^2) \quad (1)$$

$$t|s_1, s_2 \sim \mathcal{N}(t; s_1 - s_2, \sigma_t^2), \quad (2)$$

$$y|t = \begin{cases} 1 & \text{if } t \geq 0, \\ -1 & \text{else} \end{cases} \quad (3)$$

$$p(s_1, s_2, t, y) = p(s_1)p(s_2)p(t|s_1, s_2)p(y|t) \quad (4)$$

11 Thus, the model contains 5 hyper-parameters: The mean and standard deviation for each of the skills
12 $\mu_1, \mu_2, \sigma_1^2, \sigma_2^2$ respectively, as well as the standard deviation for the result of the match σ_t^2 .

13 Q.2 Bayesian Network

14 The Bayesian model that was introduced in the last section is represented graphically as a Bayesian
15 network as shown in Figure 1.

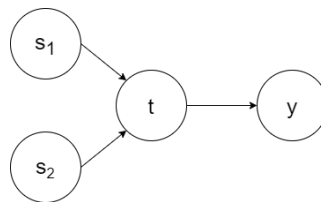


Figure 1: Bayesian Network of the model

16 From this Bayesian network we can identify the following independent sets of random variables
17 using the rules of D-separation:

- 18 • $s_1 \perp\!\!\!\perp s_2 \mid \emptyset$: The variables that show the skills of a player are independent when the
19 remaining variables in the network are non-observed. This follows from the head-to-head
20 rule.

- $(s_1, s_2) \perp\!\!\!\perp y \mid t$: The skills of the player depend solely on the result of the game if it is observed and not the games outcome. This follows from the head-to-tail node with intermediate observed variables.

24 Q.3 Computing with the Model

25 For this section, all relevant matrices are found in the appendix. The first distribution to compute is
 26 $p(s_1, s_2 | t, y)$. This is equivalent to computing $p(s_1, s_2 | t)$ as can be seen in the Bayesian network in
 27 the previous section with the conditional independence of y and s_1, s_2 derived by the head-to-tail
 28 node with an observed variable t in between them. This enables us to apply Corollary 1 from Lecture
 29 2. By this corollary we get that

$$p(s_1, s_2 | t) \stackrel{def}{=} p(\mathbf{s} | t) = \mathcal{N}(\mathbf{s}; \boldsymbol{\mu}_{\mathbf{s}|t}, \boldsymbol{\Sigma}_{\mathbf{s}|t}) \quad (5)$$

30 where we compute $\boldsymbol{\mu}_{\mathbf{s}|t}$ and $\boldsymbol{\Sigma}_{\mathbf{s}|t}$ using Corollary 1 along with the matrices in Appendix B leading
 31 us to the closed form solution of the mean and the variance.

$$\boldsymbol{\mu}_{\mathbf{s}|t} = \boldsymbol{\Sigma}_{\mathbf{s}|t} \begin{pmatrix} \sigma_1^{-2} \mu_1 + \sigma_t^{-2} t \\ \sigma_2^{-2} \mu_2 - \sigma_t^{-2} t \end{pmatrix} \quad (6)$$

$$\boldsymbol{\Sigma}_{\mathbf{s}|t} = \frac{1}{(\sigma_2^{-2} + \sigma_t^{-2})(\sigma_1^{-2} + \sigma_t^{-2}) - \sigma_t^{-4}} \begin{pmatrix} \sigma_2^{-2} + \sigma_t^{-2} & \sigma_t^{-2} \\ \sigma_t^{-2} & \sigma_1^{-2} + \sigma_t^{-2} \end{pmatrix}. \quad (7)$$

32 For finding the full conditional distribution of the outcome, we utilize Bayes' Theorem which states
 33 the proportionality relation shown in Equation 8:

$$p(t | s_1, s_2, y) \propto p(y | t) p(t | s_1, s_2). \quad (8)$$

34 Here we know from the equations introduced in section Q.1 that the factor $p(t | s_1, s_2)$ is simply given
 35 from $\mathcal{N}(t; s_1 - s_2, \sigma_t^2)$, and as such we need to evaluate the factor $p(y | t)$. Since the random variable
 36 y is defined as $\text{sgn}(t)$, we know that $y = 1 \Leftrightarrow t \geq 0$. Thus, $p(y = 1 | t) = \mathbb{I}_{\{y=\text{sign}(t)\}}$, where \mathbb{I} is an
 37 indicator function. Put together, we get:

$$p(t | s_1, s_2, y) = \begin{cases} \mathbb{I}_{\{y=\text{sign}(t)\}} \mathcal{N}(t; s_1 - s_2, \sigma_t) \cdot \frac{1}{\int_0^\infty \mathcal{N}(t; s_1 - s_2, \sigma_t) dt}, & y = 1 \\ \mathbb{I}_{\{y=\text{sign}(t)\}} \mathcal{N}(t; s_1 - s_2, \sigma_t) \cdot \frac{1}{\int_{-\infty}^0 \mathcal{N}(t; s_1 - s_2, \sigma_t) dt}, & y = -1, \end{cases} \quad (9)$$

38 i.e. this probability will take on a truncated Gaussian which we will denote as $\mathcal{TG}(t; a, b, \mu, \sigma)$ where
 39 a , and b , are the lower and upper bounds of truncation respectively, μ , and σ are the mean and
 40 standard deviation parameters of the underlying *normal* distribution.

41 The marginal probability that player 1 wins the game, $p(y = 1)$, is given by $p(t > 0)$. This probability
 42 is obtained by marginalizing out s_1 and s_2 from the joint distribution $p(s_1, s_2, t)$ using Corollary
 43 2 from the lectures. For this corollary we use the known distributions $p(\mathbf{s}) = \mathcal{N}(\mathbf{s}; \boldsymbol{\mu}_{\mathbf{s}}, \boldsymbol{\Sigma}_{\mathbf{s}})$ and
 44 $p(t | \mathbf{s}) = \mathcal{N}(t; \mathbf{A}\mathbf{s} + \mathbf{b}, \boldsymbol{\Sigma}_{t|\mathbf{s}})$. To find $p(t)$ from this corollary we compute $\boldsymbol{\mu}_{\mathbf{s}|t}$ and $\boldsymbol{\Sigma}_{\mathbf{s}}|t$ with the
 45 equations given in Corollary 2. Therefore, we can as such obtain the following marginal distribution:

$$p(t) = \mathcal{N}(t; \boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) \stackrel{def}{=} \mathcal{N}(t; \mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2 + \sigma_t^2) \quad (10)$$

46 Finally, the probability that player 1 wins is calculated as

$$p(y = 1) = p(t > 0) = \int_0^\infty \mathcal{N}(t; \mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2 + \sigma_t^2) dt \quad (11)$$

47 Q.4 A first Gibbs sampler

48 In order to implement a Gibbs sampler with an observed variable y it is required that we can sample
 49 from both of the conditional probabilities:

$$p(\mathbf{s} | t, y = 1) = \mathcal{N}(\mathbf{s}; \boldsymbol{\mu}_{\mathbf{s}|t}, \boldsymbol{\Sigma}_{\mathbf{s}|t}) \quad (12)$$

$$p(t | s_1, s_2, y = 1) = \mathcal{TG}(t; 0, \infty, s_1 - s_2; \sigma_t^2) \quad (13)$$

These equations and parameters have been evaluated in the previous section: Q.3. For the prior distribution $s_1 \sim \mathcal{N}(s_1; 0, 1)$ is used. The second skill variable, s_2 , uses the same prior distribution and $\sigma_t^2 = 2$ is chosen as a hyper-parameter.

In Figure 2a the samples for s_1 and s_2 are shown after observing the first game where $y = 1$. As an initial value for the sampling $t_* = 20$ has been chosen (this has intentionally been chosen to be slightly more off than we expect it to be in order to visualize the burn in). In Figure 2b sampling is shown again where the new prior that was obtained from the first sampling is shown. It seems take longer in the second run to reach the stationary distribution compared to the initial run we executed. We select the burning period therefore to have a size of $b = 20$ since then we have reached the stationary part in both cases and can be sure to only sample from the stationary distribution.

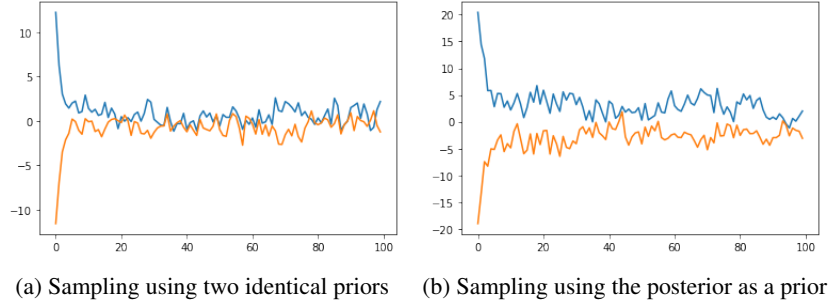


Figure 2: Gibbs sampling

In Figure 3 different choices of samples drawn with their respective required times are shown. We can see that when choosing a very small number of samples like 50 or 100, the estimated Gaussian curve doesn't fit the sample data too well. For higher numbers like 1,000 to 10,000 the Gaussian approximation looks a lot better. In order to obtain an overall algorithm with a justifiable runtime, we choose the sample size $K = 1,000$. In Figure 4 the prior and posterior distribution of the skill

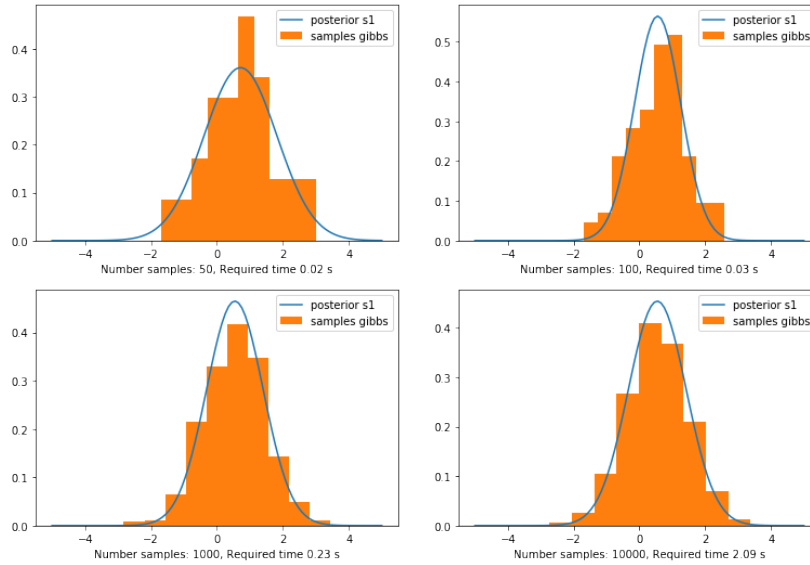


Figure 3: Different number of gibbs samples used to approximate the posterior distribution

variables s_1 and s_2 are shown. The prior distribution is just normally distributed with mean 0 and a standard deviation of 1 for both random variables. The posterior distribution is obtained after observing one data point where $y = 1$. This means that $t > 0$ and thus, player 1 won the actual game. To reflect this the mean of the skill for player 1 is increased and the mean for player 2 is decreased. This can be seen in the graph as the posterior distribution for s_1 is shifted to the right, while the one

for s_2 is shifted to the left. Since now we have more certainty about the actual skill of each player the standard deviation is decreased for each variable.

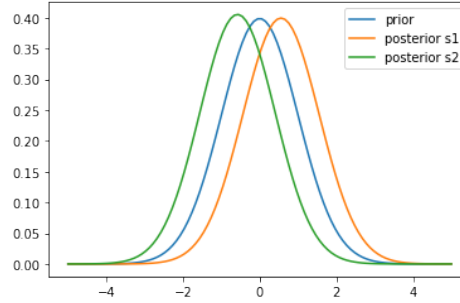


Figure 4: The prior distribution of s_1 and s_2 as well as their respective posterior distribution after observing a game with the outcome $y = 1$

Q.5 Assumed Density filtering

For each team we set the prior distribution of the skill to for each player: $p(s_x) \sim N(s_x; 0, 1)$, where $x \in \{1, 2\}$. The hyper-parameter σ_t is arbitrarily set to 2.

The final ranking of the teams is given in Figure 5 where the different distributions for the teams are shown after the whole dataset has been analyzed. We can see that the first and the second ranked team have a distinctively larger mean than the other teams. Those two teams are Juventus and Napoli. Those correspond to the champion and the vice-champion of the league. In the figure the ranking of the teams based on their mean is also given as a list. We can see that all the teams have a fairly similar standard deviation. This indicates how certain we are how that team performs and how likely that is that they will win or loose a game against a team that is either a lot better or worse respectively.

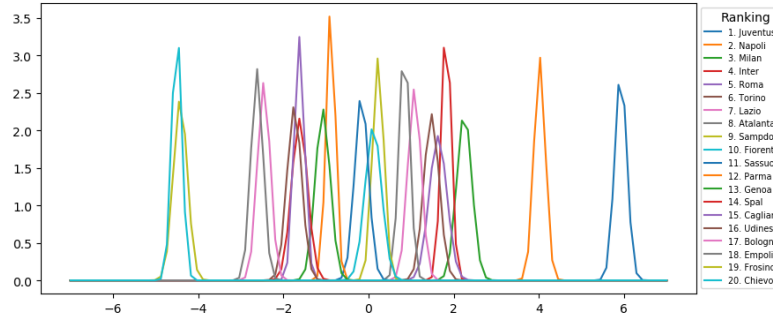


Figure 5: The distribution of the posterior skill for all teams from the dataset

When the ordering of the matches in the dataset is changed the result is most likely different. The reason for this is that with the ordered matches every team has seen an equal amount of matches after each game day. Thus, all the posterior distributions that represent the skill of the player have a similar standard deviation. When randomly shuffling the matches, this is not the case any more and we might have to estimate more often how teams perform about which we don't know the performance yet.

Q.6 Using the model for predictions

When using the predictor that is based on our TrueSkill model the obtained prediction rate is $r = 0.65$. When using a completely random guesser we would expect to achieve $r = 0.5$. However, the provided dataset is slightly biased. It is more likely for a team to win when it plays at home (it is listed as team1 in the dataset). Thus, when we predict for every match that the home team wins we obtain a

prediction rate of $r = 0.61$. Our model doesn't currently take into account this data bias. Furthermore, we can see that the performance of our predictor still outperforms the random predictors slightly.

Q.7 Factor graph

We begin by simply drawing the factor graph corresponding to the random variables and their connections in the project, shown in Figure 6.

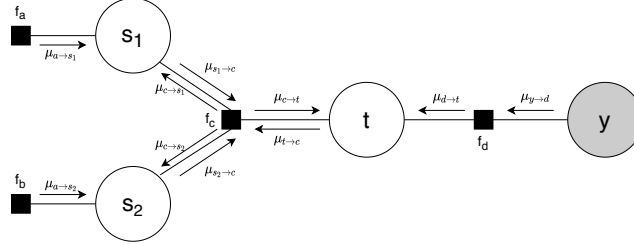


Figure 6: The random variables and the messages presented in a factor graph.

To compute the distributions $p(s_1|y)$ and $p(s_2|y)$ we use message passing, through observing the factor graph of this section. We let

$$f_a(s_1) = \mathcal{N}(s_1; \mu_1, \sigma_1^2), \quad f_b(s_2) = \mathcal{N}(s_2; \mu_2, \sigma_2^2)$$

$$f_c(t) = \mathcal{N}(t; s_1 - s_2, \sigma_t^2), \quad f_d(y) = \delta(y),$$

where δ is an indicator function. This, along with the formulas provided in Lecture 5, allows us to derive the following messages in a trivial manner:

$$\mu_{a \rightarrow s_1}(s_1) = f_a(s_1), \quad \mu_{b \rightarrow s_2}(s_2) = f_b(s_2)$$

$$\mu_{s_1 \rightarrow c}(s_1) = f_a(s_1), \quad \mu_{s_2 \rightarrow c}(s_2) = f_b(s_2)$$

$$\mu_{c \rightarrow t}(t) = \iint \mathcal{N}(t; s_1 - s_2, \sigma_t^2) \mathcal{N}(s_1; \mu_1, \sigma_1^2) \mathcal{N}(s_2; \mu_2, \sigma_2^2) ds_2 ds_1 \quad (14)$$

$$\stackrel{\text{Corollary 2}}{=} \mathcal{N}(t; \mu_1 - \mu_2, \sigma_t^2 + \sigma_1^2 + \sigma_2^2), \quad (15)$$

where the identities we use for the above corollary are in Appendix C.

Q.8 A Message passing algorithm

We have that

$$\mu_{c \rightarrow t}(t) \mu_{y \rightarrow t}(t) = \delta(y = \text{sign}(t)) \mathcal{N}(t; \mu_1 - \mu_2, \sigma_t^2 + \sigma_1^2 + \sigma_2^2),$$

and moment match $\hat{q}(t; \hat{\mu}_q, \hat{\sigma}_q^2)$ to this distribution. Moreover, we make the definition:

$$\hat{p}(t; \hat{\mu}, \hat{\sigma}^2) \stackrel{\text{def}}{=} \frac{\hat{q}(t; \hat{\mu}_q, \hat{\sigma}_q^2)}{\mu_{c \rightarrow t}(t)}$$

$$= \mathcal{N}(t; \hat{\mu}, \hat{\sigma}^2) \approx \mu_{t \rightarrow c}(t).$$

Finally

$$\mu_{c \rightarrow s_1}(s_1) \approx \int \mathcal{N}(t; \hat{\mu}, \hat{\sigma}^2) \int [\mathcal{N}(t; s_1 - s_2, \sigma_t^2) \mathcal{N}(s_2; \mu_2, \sigma_2^2) ds_2] dt \quad (16)$$

$$\stackrel{\text{Corollary 2}}{=} \int \mathcal{N}(t; \hat{\mu}, \hat{\sigma}^2) \mathcal{N}(t; s_1 - \mu_2, \sigma_t^2 + \sigma_2^2) dt, \quad (17)$$

$$= \int \mathcal{N}(s_1; t + \mu_2, \sigma_t^2 + \sigma_2^2) \mathcal{N}(t; \hat{\mu}, \hat{\sigma}^2) \quad (18)$$

where the matrices in the Corollary 2 are found in the Appendix C.2. Applying Corollary 2 with matrices found in Appendix C.3 with the following: gives us that

$$\mu_{c \rightarrow s_1}(s_1) \approx \mathcal{N}(s_1; \hat{\mu} + \mu_2, \sigma_t^2 + \sigma_2^2 + \hat{\sigma}^2). \quad (19)$$

Finally, the desired posterior probability is computed by the equality below:

$$p(s_1|y) \propto \mu_{c \rightarrow s_1}(s_1) \mathcal{N}(s_1; \mu_1, \sigma_1^2).$$

Similarly we get that the posterior distribution for s_2 is:

$$p(s_2|y) \propto \mathcal{N}(s_2; \hat{\mu} - \mu_1, \sigma_t^2 + \sigma_1^2 + \hat{\sigma}^2) \mathcal{N}(s_2; \mu_2, \sigma_2^2).$$

With this we can now estimate the posterior distribution using moment-matching and compare this distribution with the one obtained from the Gibbs sampling as shown in Figure 7.

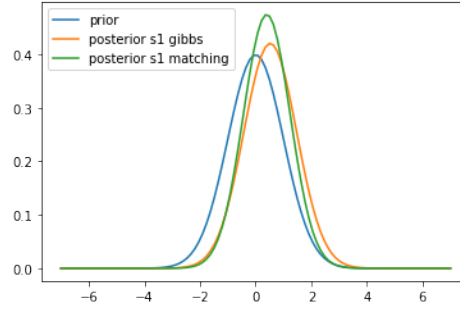


Figure 7: Prior distribution and posterior distribution estimated using a gibbs sampler and moment-matching

112

113 Q.9 Our own data

The developed Trueskill method is tested on a dataset containing tennis matches from the ATP tournament.¹ We only take a look at the matches that have taken place in the last 5 years in order to have a manageable amount of different players and matches between them.

Our model achieves a prediction accuracy of 65% on this dataset. We disregarded Gibbs sampling as the run time of the algorithm increases significantly due to the relatively high number of needed samples and a high number of total matches.

120 Q.10 Project extension

As already described in section Q.6 where we try to predict the outcome of the upcoming matches the data includes a small bias towards the chance of the home team winning. In reality we can see that the probability that the home team wins is about 60%. We slightly adjust our model to take that factor into account. Therefore, when we predict the outcome of the next game we shift the mean of distribution to the right by a factor $\Delta\mu$ such that $p(y = 1) = 0.6$ given that both teams have the same mean. Thus we chose an approximate $\Delta\mu$ such that the following holds for teams with an equal mean

$$p(y = 1) = p(t > 0) = \int_0^\infty N(t; \Delta\mu; \sigma_1^2 + \sigma_2^2 + \sigma_t^2) = 0.6 \quad (20)$$

The value $\Delta\hat{\mu}$ can be looked up in a table for the Gaussian normal distribution. We choose $\Delta\hat{\mu} = 0.25\sqrt{(\sigma_1^2 + \sigma_2^2 + \sigma_t^2)}$. The marginal distribution of y is now given by

$$p(y = 1) = p(t > 0) = \int_0^\infty N(t; \mu_1 - \mu_2 + \Delta\hat{\mu}; \sigma_1^2 + \sigma_2^2 + \sigma_t^2). \quad (21)$$

Adding this a-priori advantage for the home team to our model increases the performance of the predictions from 65% to approximately 70%.

¹<https://www.kaggle.com/sijovm/atpdata>

132 A Code

133 A.1 Code for Quiz

```
134 #project test
135 m1=1
136 sigma1=1#sigma is the variance not the volatility here!
137 m2=-1
138 sigma2=4
139 sigma3=5
140 #task 1
141 div=(sigma2**(-1)+sigma3**(-1))*(sigma1**(-1)+sigma3**(-1))-sigma3**(-2)
142 sigma_st=np.array([[sigma2**(-1)+sigma3**(-1), sigma3**(-1)],
143                   [sigma3**(-1), sigma1**(-1)+sigma3**(-1)]])/div
144 t=3
145 print(sigma_st)
146 multiplier=np.array([[sigma1**(-1)*m1+sigma3**(-1)*t],[sigma2**(-1)*m2-sigma3**(-1)*t]])
147 mu_st=np.matmul(sigma_st,multiplier)
148 print(f"mu_st is:{mu_st} and sigma_st is:{sigma_st}")
149
150 #task 2
151 py_1=1-norm.cdf(0,m1-m2,np.sqrt(sigma1+sigma2+sigma3))#takes in square root
152 print(f"p(y=1) is:{py_1}")
```

156 A.2 Script for report results

```
157 import pandas as pd
158 import numpy as np
159 from scipy.stats import norm, truncnorm, multivariate_normal
160 from matplotlib import pyplot as plt
161 import time
162
163
164 def sample_s(m1, m2, sigma1, sigma2, sigma3, t):
165     div=(sigma2**(-1)+sigma3**(-1))*(sigma1**(-1)+sigma3**(-1))-sigma3**(-2)
166     sigma_st=np.array([[sigma2**(-1)+sigma3**(-1), sigma3**(-1)],
167                       [sigma3**(-1), sigma1**(-1)+sigma3**(-1)]])/div
168     multiplier=np.array([[sigma1**(-1)*m1+sigma3**(-1)*t],[sigma2**(-1)*m2-sigma3**(-1)*t]])
169     mu_st=np.matmul(sigma_st,multiplier)
170     # sample s1, s2
171     return multivariate_normal.rvs(mean=mu_st.flatten(), cov=sigma_st)
172
173
174 # compare question 2 in the quiz
175 def sample_t(s1, s2, y, sigma_t):
176     if y == 1:
177         return truncnorm.rvs(0, np.inf, loc=(s1-s2), scale=np.sqrt(sigma_t))
178     else:
179         return truncnorm.rvs(np.NINF, 0, loc=(s1-s2), scale=np.sqrt(sigma_t))
180
181
182 # used to predict new games
183 def marginal_y(mu1, mu2, sigma1, sigma2, sigmat):
184     return 1-norm.cdf(0,mu1-mu2, np.sqrt(sigma1+sigma2+sigmat))
185
186
187 # calculate the marginal of y includeing the data bias
188 def marginal_y_biased(mu1, mu2, sigma1, sigma2, sigmat):
189     # 0.3 standard deviations equals to 11.79%
190     offset_mu = 0.25 * np.sqrt(sigma1+sigma2+sigmat)
191     return 1-norm.cdf(0,mu1-mu2+offset_mu, np.sqrt(sigma1+sigma2+sigmat))
192
```

```

193
194
195 def Gauss_mult(m1, m2, s1, s2):
196     s = 1 / (1 / s1 + 1 / s2)
197     m = (m1 / s1 + m2 / s2) * s
198     return m, s
199
200
201 def Gauss_div(m1, m2, s1, s2):
202     m, s = Gauss_mult(m1, m2, s1, -s2)
203     return m, s
204
205
206 def Gauss_trunc(a, b, m1, s1):
207     b_norm = (b - m1) / np.sqrt(s1)
208     a_norm = (a - m1) / np.sqrt(s1)
209     m = truncnorm.mean(a_norm, b_norm, loc=m1, scale=np.sqrt(s1))
210     s = truncnorm.var(a_norm, b_norm, loc=m1, scale=np.sqrt(s1))
211     return m, s
212
213
214 # the gibbs sampler should return mean1, std1, mean2 and std2 for the two players
215 # burn in and number of samples empirically chosen
216 def gibbs_sampler(prior_1, prior_2, y, sigma_t=2, t_=50, b=20, K=1000):
217     # lists to store the samples
218     s1 = []
219     s2 = []
220
221     # here is where the magic happens
222     for i in range(K):
223         s = sample_s(prior_1[0], prior_2[0], prior_1[1], prior_2[1], sigma_t, t_)
224         s1_ = s[0]
225         s1.append(s1_)
226         s2_ = s[1]
227         s2.append(s2_)
228         t_ = sample_t(s1_, s2_, y, sigma_t)[0]
229
230     # return the sample values without the burn period
231     return s1[b:K], s2[b:K]
232
233
234 # Q4.2
235 def estimate_posterior(prior_1, prior_2, y, sigma_t=2):
236     s1, s2 = gibbs_sampler(prior_1, prior_2, y, sigma_t=sigma_t)
237     # estimate mean and std for the two players
238     s1_m = estimate_mean(s1)
239     s1_s = estimate_std(s1, s1_m)
240     s2_m = estimate_mean(s2)
241     s2_s = estimate_std(s2, s2_m)
242     # return results posterior_1 and posterior_2
243     return (s1_m, s1_s), (s2_m, s2_s)
244
245
246 # Q4.2
247 def estimate_mean(X):
248     return 1/len(X) * sum(X)
249
250
251 def estimate_std(X, mean):
252     return 1/len(X) * sum([(x - mean) ** 2 for x in X])
253
254
255 def q4():
256     prior = (0.0, 1.0)
257     k = 100

```



```

258     s1, s2 = gibbs_sampler(prior, prior, y=1, b=0, K=k)
259     s1_m = estimate_mean(s1)
260     s1_s = estimate_std(s1, s1_m)
261     s2_m = estimate_mean(s2)
262     s2_s = estimate_std(s2, s2_m)
263     plt.plot(range(0, k), s1,
264              range(0, k), s2)
265     plt.show()
266     sigma_t = 2
267     k = 100
268     s1, s2 = gibbs_sampler((s1_m, s1_s), (s2_m, s2_s), 1, b=0, K=k)
269     plt.plot(range(0, k), s1,
270              range(0, k), s2)
271     plt.show()
272     prior = (0.0, 1.0)
273     # test out different k values and for each plot the distribution and the
274     # estimated gaussian
275     for k_ in [50, 100, 150, 500, 1000, 5000, 10000, 25000]:
276         start = time.time()
277         s1, s2 = gibbs_sampler(prior, prior, y=1, K=(k_+20))
278         s1_m = estimate_mean(s1)
279         s1_s = estimate_std(s1, s1_m)
280         s2_m = estimate_mean(s2)
281         s2_s = estimate_std(s2, s2_m)
282         end = time.time()
283         elapsed = end - start
284         x = np.linspace(-5, 5, 100)
285         plt.plot(x, norm.pdf(x, s1_m, s1_s), label="posterior s1")
286         plt.hist(s1, density=1, label="samples gibbs")
287         plt.xlabel("Number of gibbs samples: {}, Required time {:.2f} s".format(k_,
288                                     elapsed))
289         plt.legend()
290         plt.show()
291     prior = (0.0, 1.0)
292     # run gibbs sampler and estimate posterior gaussians
293     s1, s2 = gibbs_sampler(prior, prior, y=1, K=1000)
294     s1_m = estimate_mean(s1)
295     s1_s = estimate_std(s1, s1_m)
296     s2_m = estimate_mean(s2)
297     s2_s = estimate_std(s2, s2_m)
298     # plot 4 gaussian distributions
299     x = np.linspace(-5, 5, 100)
300     plt.plot(x, norm.pdf(x, 0, 1), label="prior")
301     plt.plot(x, norm.pdf(x, s1_m, s1_s), label="posterior s1")
302     plt.plot(x, norm.pdf(x, s2_m, s2_s), label="posterior s2")
303     plt.legend()
304     plt.show()
305
306
307 def output_result(x):
308     # output the binary result based on a winning probability
309     if x >= 0.5:
310         y_ = 1
311     else:
312         y_ = -1
313     return y_
314
315
316 def get_y(row):
317     if row["score1"] == row["score2"]:
318         return 0
319     elif row["score1"] > row["score2"]:
320         return 1
321     else:
322         return -1

```

```

323
324
325 def q56(shuffle=False, gibbs=True, advanced_predictor=False):
326     print("Using Serie A data")
327     data = pd.read_csv('SerieA.csv')
328     # create new row with y values
329     data['y'] = data.apply(lambda x: get_y(x), axis=1)
330     # filter out the games where the teams draw
331     data = data.loc[data['y'] != 0]
332     # create a list of all the teams and store their skill representation
333     d_teams = data['team1'].unique()
334     teams = pd.DataFrame(index=d_teams)
335     # set starting values for the mean and the std of the skill representation
336     teams['mean'] = 0.0
337     teams['std'] = 1.0
338     if shuffle:
339         print("Shuffle data in advance for Q.5")
340         data = data.sample(frac=1).reset_index(drop=True)
341     correct_predictions = 0
342     random_predictions = 0
343     # run through all the matches
344     for index, row in data.iterrows():
345         # fetch current values for team1 and team2
346         s1 = teams.loc[row['team1']]
347         s2 = teams.loc[row['team2']]
348         # get priors
349         prior_1 = (s1['mean'], s1['std'])
350         prior_2 = (s2['mean'], s2['std'])
351         sigma_t = 2
352         # Q6: predict who should win based on the model
353         if advanced_predictor:
354             y_ = output_result(marginal_y_biased(prior_1[0], prior_2[0], prior_1[1],
355             prior_2[1], sigma_t))
356         else:
357             y_ = output_result(marginal_y(prior_1[0], prior_2[0], prior_1[1],
358             prior_2[1], sigma_t))
359         # fetch the actual result
360         y = row['y']
361         if y == 1:
362             random_predictions += 1
363         if y == y_:
364             correct_predictions += 1
365         # run sampling
366         if gibbs:
367             posterior_1, posterior_2 = estimate_posterior(prior_1, prior_2, y,
368             sigma_t=sigma_t)
369         else:
370             posterior_1, posterior_2 = moment_matching(prior_1, prior_2, y,
371             st=sigma_t)
372         # rewrite values mean and std
373         s1['mean'] = posterior_1[0]
374         s1['std'] = posterior_1[1]
375         s2['mean'] = posterior_2[0]
376         s2['std'] = posterior_2[1]
377
378     if advanced_predictor:
379         print("Using advanced predictor")
380     if gibbs:
381         print("Using gibbs sampling")
382     else:
383         print("Using moment matching")
384     print(f"The prediction rate is {correct_predictions / data.shape[0]}\nUsing
385     random guessing it is {random_predictions / data.shape[0]}")
386     teams = teams.sort_values(["mean", "std"], ascending=False)
387     x = np.linspace(-10, 10, 100)

```

```

388     for num, (index, row) in enumerate(teams.iterrows()):
389         l = str(num + 1) + ". " + index
390         plt.plot(x, norm.pdf(x, row['mean'], row['std']), label=l)
391     plt.legend(title="Ranking", loc="upper left", bbox_to_anchor=(1,1),
392               fontsize='x-small')
393     plt.show()
394
395
396 def moment_matching(prior1, prior2, y, st=2):
397     m1 = prior1[0]
398     s1 = prior1[1]
399     m2 = prior2[0]
400     s2 = prior2[1]
401     #we call this c
402     mus2_ft_m=m2
403     mus2_ft_s=s2
404     mus1_ft_m=m1
405     mus1_ft_s=s1
406     muft_t_m=mus1_ft_m-mus2_ft_m
407     muft_t_s=mus1_ft_s+mus2_ft_s+st
408     # truncated gaussian approximation via moment-matching
409     if y == -1:
410         a = np.NINF
411         b = 0
412     else:
413         a = 0
414         b = 1000
415     # q/mufxt_t approximates mut_fxt
416     q_m, q_s = Gauss_trunc(a, b, muft_t_m, muft_t_s)
417     p_m, p_s = Gauss_div(q_m, muft_t_m, q_s, muft_t_s)
418     #see notes
419     muft_s1_m=m2+p_m
420     muft_s1_s=p_s+st+s2
421     s1_m, s1_s = Gauss_mult(muft_s1_m, m1, muft_s1_s, s1)
422     muft_s2_m=m1-p_m
423     muft_s2_s=p_s+st+s1
424     s2_m, s2_s = Gauss_mult(muft_s2_m, m2, muft_s2_s, s2)
425     return (s1_m, s1_s), (s2_m, s2_s)
426
427
428 def q78():
429     prior = (0.0, 1.0)
430     # run gibbs sampler and estimate posterior gaussians
431     s1, s2 = gibbs_sampler(prior, prior, y=1, K=1000)
432     s1_m = estimate_mean(s1)
433     s1_s = estimate_std(s1, s1_m)
434     s2_m = estimate_mean(s2)
435     s2_s = estimate_std(s2, s2_m)
436     # run matching
437     s1_, s2_ = moment_matching(prior, prior, y=1)
438     # plot 4 gaussian distributions
439     x = np.linspace(-4, 4, 100)
440     plt.plot(x, norm.pdf(x, 0, 1), label="prior")
441     plt.plot(x, norm.pdf(x, s1_m, s1_s), label="posterior s1 gibbs")
442     plt.plot(x, norm.pdf(x, s1_[0], s1_[1]), label="posterior s1 moment-matching")
443     # plt.plot(x, norm.pdf(x, s2_m, s2_s), label="posterior s2")
444     plt.legend()
445     plt.show()
446     plt.plot(x, norm.pdf(x, 0, 1), label="prior")
447     plt.plot(x, norm.pdf(x, s1_[0], s1_[1]), label="posterior s1 moment-matching")
448     plt.plot(x, norm.pdf(x, s2_[0], s2_[1]), label="posterior s2 moment-matching")
449     plt.legend()
450     plt.show()
451
452

```

```

453 def q9():
454     print("Using the ATP dataset")
455     data = pd.read_csv('ATP.csv')
456     # only consider events from the last 5 years
457     data = data.loc[data["tourney_date"] > 20150000, ["winner_id", "winner_name",
458         "loser_id", "loser_name", "tourney_date"]]
459     d_teams = np.unique(np.append(data['winner_id'].unique(),
460         data['loser_id'].unique()))
461     teams = pd.DataFrame(index=d_teams)
462     # set starting values for the mean and the std of the skill representation
463     teams['mean'] = 0.0
464     teams['std'] = 1.0
465     correct_predictions = 0
466     # run through all the matches
467     for index, row in data.iterrows():
468         # Switch between Gibbs and moment-matching
469         gibbs = False
470         # fetch current values for team1 and team2
471         s1 = teams.loc[row['winner_id']]
472         s2 = teams.loc[row['loser_id']]
473         # get priors
474         prior_1 = (s1['mean'], s1['std'])
475         prior_2 = (s2['mean'], s2['std'])
476         sigma_t = 2
477         # predict next game
478         y_ = output_result(marginal_y(prior_1[0], prior_2[0], prior_1[1],
479             prior_2[1], sigma_t))
480         # in this dataset there is no home or away s1 is per definition the winner
481         # of the game
482         y = 1
483         if y == y_:
484             correct_predictions += 1
485         # run sampling
486         if gibbs:
487             posterior_1, posterior_2 = estimate_posterior(prior_1, prior_2, y,
488                 sigma_t=sigma_t)
489         else:
490             posterior_1, posterior_2 = moment_matching(prior_1, prior_2, y,
491                 st=sigma_t)
492         # rewrite values mean and std
493         s1['mean'] = posterior_1[0]
494         s1['std'] = posterior_1[1]
495         s2['mean'] = posterior_2[0]
496         s2['std'] = posterior_2[1]
497     print(f"The prediction rate is {correct_predictions / data.shape[0]}")
498
499
500 def main():
501     q4()
502     q56()
503     print("\n")
504     q56(shuffle=True)
505     print("\n")
506     q56(gibbs=False)
507     print("\n")
508     q56(gibbs=False, advanced_predictor=True)
509     print("\n")
510     q78()
511     q9()
512
513
514 if __name__ == '__main__':
515     main()

```

517 **B Corollary 1 Matrices**

518 The matrices used for the application of Corollary 1 in Section Q.3 are given in standard notation as:

$$\mathbf{A} = \begin{pmatrix} 1 & -1 \end{pmatrix} \quad (22)$$

$$\mathbf{\Sigma}_s = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix} \quad (23)$$

$$\mathbf{\Sigma}_{t|s} = \sigma_t^2 \quad (24)$$

$$\boldsymbol{\mu}_s = (\mu_1 \quad \mu_2) . \quad (25)$$

519 **C Corollary 2 Matrices**

520 **C.1 First usage of Corollary 2**

521 The first usage of Corollary 2 was in (14) and we used the following matrices:

$$\boldsymbol{\mu}_s = (\mu_1, \mu_2)^T, \quad \mathbf{\Sigma}_s \stackrel{s_1 \perp s_2}{=} \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix} \quad (26)$$

$$\mathbf{A} = (1, -1), \quad \mathbf{b} = 0, \quad \mathbf{\Sigma}_{t|s} = \sigma_t^2 \quad (27)$$

522 **C.2 Second usage of Corollary 2**

523 The matrices created for the second application of Corollary 2 in (16) is given in standard notation as:

$$\mathbf{A} = -1, \quad \mathbf{b} = s_1 \quad (28)$$

$$\boldsymbol{\mu}_a = \mu_2, \quad \mathbf{\Sigma}_{b|a} = \sigma_t^2. \quad (29)$$

524 **C.3 Third usage of Corollary 2**

525 The third usage of Corollary 2 was in (19) and utilized the following vectors.

$$A = 1, \quad b = \mu_2 \quad (30)$$

$$\mu_a = \hat{\mu}, \quad \mathbf{\Sigma}_{b|a} = \sigma_t^2 + \sigma_2^2. \quad (31)$$