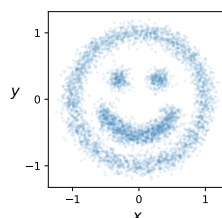# Example sheet 2
Bayesian inference
Data Science—DJW—2022/2023

All Bayesian questions are the same: first write out the likelihood of the observed data $\Pr(x|\Theta = \theta)$, then the prior likelihood $\Pr_\Theta(\theta)$, then apply Bayes's rule to get the posterior likelihood $\Pr_\Theta(\theta \,|\, X = x)$. Questions 3–8 are repetitions of this idea, in progressively more complex settings. You needn't attempt them all: instead, work through enough of the earlier questions for you to be confident in answering question 8.

What's more important than the algebra is getting used to reformulating high-level questions into mathematical questions about distributions of random variables. Even if you don't answer a question, read it carefully, and add it to your repertoire of "how data scientists ask questions".

**Question 1.** Define a function `rxy()` that produces a random pair of values $(X, Y)$ which, when shown in a scatterplot, produces a smiley face like this. Also plot the marginal distributions of $X$ and $Y$.



**Question 2.** Consider this code for generating random variables $X$ and $Y$:

```
x = np.random.uniform()
y = np.random.geometric(p=x)
```

Derive the marginal likelihood $\Pr_Y(y)$, and the conditional likelihood $\Pr_X(x \,|\, Y = y)$.

**Question 3.** I sample $x_1, \ldots, x_n$ from $\mathrm{Uniform}[0, \theta]$. The parameter $\theta$ is unknown, and I shall use $\Theta \sim \mathrm{Pareto}(b_0, \alpha_0)$ as my prior, where $b_0 > 0$ and $\alpha_0 > 1$ are known. This has the cumulative distribution function

$$\mathbb{P}(\Theta \leq \theta) = \begin{cases} 1 - \left(b_0/\theta\right)^{\alpha_0} & \text{if } \theta \geq b_0, \\ 0 & \text{if } \theta < b_0. \end{cases}$$

(a) Calculate the prior likelihood for $\Theta$.
(b) Show that the posterior distribution of $(\Theta \,|\, x_1, \ldots, x_n)$ is Pareto, and find its parameters.
(c) Find a 95% posterior confidence interval for $\Theta$.
(d) Find a different 95% posterior confidence interval. Which is better? Why?

**Question 4.** I have a collection of numbers $x_1, \ldots, x_n$ which I take to be independent samples from the $\mathrm{Normal}(\mu, \sigma_0^2)$ distribution. Here $\sigma_0$ is known, and $\mu$ is unknown. Using the prior distribution $M \sim \mathrm{Normal}(\mu_0, \rho_0^2)$ for $\mu$, show that the posterior density is

$$\Pr_M(\mu \,|\, x_1, \ldots, x_n) = \kappa e^{-(\mu - c)^2 / 2\tau^2}$$

where $\kappa$ is a normalizing constant, and where you should find formulae for $c$ and $\tau$ in terms of $\sigma_0$, $\mu_0$, and $\rho_0$, and the $x_i$. Hence deduce that the posterior distribution is $\mathrm{Normal}(c, \tau^2)$. *[Note: 'M' is the upper-case form of the Greek letter 'μ'.]*

**Question 5 (Leaky priors).** I repeatedly attempt a task, and each time I attempt it I succeed with probability $\theta$ and fail with probability $1 - \theta$. The parameter $\theta$ is unknown, so I model it as

a random variable $\Theta$. Ever the optimist, my prior for $\Theta$ is heavily biased in favour of large values for $\theta$:

$$\Pr_\Theta(\theta) = \varepsilon 1_{\theta \leq 1/2} + (2 - \varepsilon)1_{\theta > 1/2}$$

for some known small value $\varepsilon > 0$; this implies $\mathbb{P}(\Theta \leq {}^1\!/_2) = \varepsilon/2$.

But I experience an unbroken run of $n$ failures. How big does $n$ need to be, for me to concede there's a 50% posterior probability that $\Theta \leq {}^1\!/_2$? How big would it need to be, if $\varepsilon = 0$?

**Question 6.** I have a collection of numbers

$$[4.3, 2.8, 3.9, 4.1, 9, 4.5, 3.3]$$

which look like they mostly come from a Gaussian distribution, but with the occasional outlier. Model the data as

$$X \text{ is } \begin{cases} \text{Normal}(\mu, 0.5^2) & \text{with probability } 99\% \\ \text{Cauchy} & \text{with probability } 1\%. \end{cases}$$

Use a Normal$(0, 5^2)$ prior distribution for $\mu$. Give pseudocode to plot the posterior distribution. *[Note. The Cauchy random variable occasionally generates wildly huge values. The library function* `scipy.stats.cauchy.pdf(x)` *computes its pdf.]*

**Question 7.** In lecture notes section 2.6 we investigated a dataset of police stop-and-search actions. Let the outcome for record $i$ be $y_i \in \{0, 1\}$, where 1 denotes that the police found something and 0 denotes that they found nothing. Consider the probability model $Y_i \sim \text{Binom}(1, \beta_{\text{eth}_i})$ where $\text{eth}_i$ is the recorded ethnicity for the individual involved in record $i$, and where the parameters $\beta_{\text{As}}$, $\beta_{\text{Blk}}$, $\beta_{\text{Mix}}$, $\beta_{\text{Oth}}$, $\beta_{\text{Wh}}$ are unknown. As a prior distribution, suppose that the five $\beta$ parameters are all independent Beta$({}^1\!/_2, {}^1\!/_2)$ random variables.

(a)   Write down the joint prior density for $(\beta_{\text{As}}, \beta_{\text{Blk}}, \beta_{\text{Mix}}, \beta_{\text{Oth}}, \beta_{\text{Wh}})$.

(b)   Find the joint posterior distribution of $(\beta_{\text{As}}, \beta_{\text{Blk}}, \beta_{\text{Mix}}, \beta_{\text{Oth}}, \beta_{\text{Wh}})$ given the $y$ data.

**Question 8.** I am prototyping a diagnostic test for a disease. In healthy patients, the test result is Normal$(0, 2.1^2)$. In sick patients it is Normal$(\mu, 3.2^2)$, but I have not yet established a firm value for $\mu$. In order to estimate $\mu$, I trialled the test on 30 patients whom I know to be sick, and the mean test result was 10.3. I subsequently apply the test to a new patient, and get the answer 8.8. I wish to know whether this new patient is healthy or sick.

(a)   In this question there are two unknown quantities: $\mu$, and $h \in \{\text{healthy, sick}\}$ the status of the new patient. Model the former as a random variable $M$ with prior distribution Normal$(5, 3^2)$ and the latter as a random variable $H$ with prior distribution

$$\Pr_H(h) = 0.99 \times 1_{h = \text{healthy}} + 0.01 \times 1_{h = \text{sick}}.$$

Write down the joint prior likelihood for $(M, H)$.

(b)   In this question the data consists of 31 values, test results $x_1, \ldots, x_{30}$ from the known sick patients and test result $y$ from the new patient. Write down the data likelihood $\Pr(x_1, \ldots, x_{30}, y \mid \mu, h)$.

(c)   Find the posterior density of $(M, H)$. Leave your answer as an unnormalized density function. It should simplify to be a function of $\bar{x}$ and $y$, where $\bar{x}$ is the mean test result for the known sick patients.

(d)   Give pseudocode to compute the posterior distribution of $H$, i.e. compute $\mathbb{P}(H = h \mid \text{data})$ for both $h = \text{healthy}$ and $h = \text{sick}$.

**Question 9.** In the lecture notes on linear modelling, we proposed a linear model for temperature increase:

$$\texttt{temp} \approx \alpha + \beta_1 \sin(2\pi \texttt{t}) + \beta_2 \cos(2\pi \texttt{t}) + \gamma(\texttt{t} - 2000).$$

Suggest a probability model for `temp`. Suggest Bayesian prior distributions for the unknown parameters $\alpha$, $\beta_1$, $\beta_2$, and $\gamma$. Give pseudocode to find a 95% confidence interval for $\gamma$.

# Hints and comments

**Question 1.** Try extending the Gaussian mixture model from section 1. For plotting, here's some code. It assumes that you have stored your samples in a numpy array of shape $n \times 2$, one row per sample point, columns for $x$ and $y$.

```
fig,((ax_x,dummy),(ax_xy,ax_y)) = plt.subplots(2,2, figsize=(4,4),
    sharex='col', sharey='row', gridspec_kw='height_ratios':[1,2], 'width_ratios':[2,1])
dummy.remove()
ax_xy.scatter(xy[:,0], xy[:,1], s=3, alpha=.1)
ax_x.hist(???, density=True, bins=60) # fill in the ???
ax_y.hist(???, density=True, bins=60, orientation='horizontal') # fill in the ???
plt.show()
```

**Question 2.** There are two versions of the Geometric distribution; look up the numpy help page to see which one is being used here. For the marginal likelihood, write out the joint likelihood and integrate. For the conditional likelihood, the calculation is similar to exercise 4.2.2 from lecture notes.

**Question 3.** For part (a), just differentiate the cdf to get the pdf, i.e. the likelihood. Write it out using indicator function notation, $1_{\theta \geq b_0}$. This is often a good idea, when we're working with parameters that affect boundaries.

For the rest: **all Bayesian calculations start in exactly the same way.** First write out the likelihood of the observed data $\Pr(x_1, \ldots, x_n \mid \Theta = \theta)$, then (1) write down the prior likelihood $\Pr_\Theta(\theta)$, (2) apply Bayes's rule which says that the posterior likelihood is

$$\Pr_\Theta(\theta \mid x_1, \ldots, x_n) = \kappa \Pr_\Theta(\theta) \Pr(x_1, \ldots, x_n \mid \Theta = \theta).$$

In this question, write out the likelihood of the data using indicator notation, as in example sheet 1 question 4. Once you have the posterior density, gather together the $\theta$ terms, and you should end up with the density of another Pareto.

For the posterior confidence interval: the definition of a posterior confidence interval is in lecture notes section 7.4. You just have to solve the equations for `lo` and `hi`, using the cumulative distribution function for the Pareto.

**Question 4. All Bayesian calculations start in exactly the same way.** First write out the likelihood of the observed data $\Pr(x_1, \ldots, x_n \mid M = \mu)$, then (1) write down the prior likelihood $\Pr_M(\mu)$, (2) apply Bayes's rule which says that the posterior likelihood is

$$\Pr_M(\mu \mid x_1, \ldots, x_n) = \kappa \Pr_M(\mu) \Pr(x_1, \ldots, x_n \mid M = \mu).$$

Remember, this is a density function for a random variable $M$, and the argument is $\mu$. Write your answer to gather together all the $\mu$ terms as much as you can. This involves expanding quadratic terms and completing the square. Any terms that don't involve $\mu$ can be amalgamated with the constant factor $\kappa$. What you end up with should look like a Normal density function, as a function of $\mu$, and this lets you conclude that the posterior distribution is Normal.

When a question asks "find the posterior distribution", you should start by calculating the posterior density, leaving it unnormalized i.e. including a constant factor, call it $\kappa$. Then (a) if you recognize this as a standard density function, as in this case, just give its name; (b) if it's easy to find $\kappa$ using "densities sum to one" then do so; (c) otherwise leave your answer as an unnormalized density function.

**Question 5.** Let $x$ be the number of successes, modelled as $X \sim \text{Bin}(n, \theta)$. We observe $x = 0$. Use the usual Bayesian method to find the posterior likelihood $\Pr_\Theta(\theta \mid X = 0)$; you can find the normalizing constant $\kappa$ with simple integration, splitting the integral over $0 \leq \theta \leq 1$ into two parts, $0 \leq \theta \leq 1/2$ and $1/2 < \theta \leq 1$. Once you've calculated the posterior likelihood, the posterior

probability that the question is referring to is $\mathbb{P}(\Theta \leq {}^1\!/2 \,|\, X = x)$, which you can find by integrating the posterior likelihood.

**Question 6. All Bayesian computations start in exactly the same way.** First write out the likelihood of the data, $\Pr(x_1, \ldots, x_n \,|\, M = \mu)$. The probability model here is very similar to a Gaussian mixture model, which we analysed in mock exam question 1. You'll need the cdf for the Cauchy, but you don't actually need to know a formula for it: just write $\mathrm{cdf}_{\mathrm{Cauchy}}(x)$ and $\mathrm{pdf}_{\mathrm{Cauchy}}(x)$. Then, (1) take a sample $\mu_1, \ldots, \mu_n$ from the prior distribution, (2) compute weights by evaluating the likelihood of the data at each one of these sampled $\mu$-values, and rescaling so they sum to one.

For plotting the posterior distribution, see the examples in section 7.3.

**Question 7.** This is a Bayesian question with multiple unknown parameters. You need to start with a joint prior density for all of them,

$$\Pr(\beta_{\mathrm{As}}, \beta_{\mathrm{Blk}}, \beta_{\mathrm{Mix}}, \beta_{\mathrm{Oth}}, \beta_{\mathrm{Wh}}).$$

See the mathematical solution to exercise 7.3.2 in lecture notes.

Bayes's rule, in its general form, says that

$$\Pr_{\Theta}(\theta \,|\, x) = \kappa \Pr_{\Theta}(\theta) \Pr_X(x | \Theta = \theta)$$

where $\theta$ denotes *all* the unknown parameters and $x$ denotes *all* the dataset. Again, see exercise 7.3.2 in lecture notes. Leave your answer with $\kappa$.

After you've found the joint posterior density function, see if you can recognize it from the list of standard random variables.

**Question 8. This is a question about multiple unknowns,** using both the mathematical and the computational solutions. See exercise 7.3.2 from lecture notes.

For part (b), for the likelihood $\Pr_Y(y \,|\, \mu, h)$, see the Gaussian mixture model in exercise 4.3.5 in lecture notes.

For part (c), your formula for the posterior distribution will involve equations very similar to question 4.

Part (d) is a question about using marginalization to ignore nuisance parameters. See exercise 7.3.2 from lecture notes for an example of marginalization, and exercise 7.3.3 for a similar calculation about posteriors over binary outcomes.

**Question 9.** You should implement your proposed Bayesian model, and find a numerical value for the confidence interval. You can find a code skeleton at `https://github.com/damonjw/datasci/blob/master/ex2.ipynb`.

It's up to you to invent whatever probability distribution you like for `temp`; the simplest choice is to assume Gaussian errors as in section 2.4, and to pluck the noise parameter out of thin air. If you truly are uncertain about the noise parameter, then treat it as a random variable and invent a prior distribution for it.

It's up to you to invent whatever priors you like for the unknown parameters. It may seem totally arbitrary, but that's Bayesianism for you.

*These questions are not intended for supervision (unless your supervisor directs you otherwise). Some of require careful maths, some are best answered with coding, some are philosophical.*

**Question 10.** Consider this code for generating random variables $X \to Y \to Z$:

```
x = np.random.uniform()
y = np.random.binomial(n=1, p=x)
z = np.random.normal(loc=y, scale=ε)
```

Show that

$$\Pr_Y(1 \mid X = x, Z = z) = \frac{x}{x + (1-x)e^{(1-2z)/2\varepsilon^2}}.$$

How does $\Pr_Y(1 \mid X = x, Z = z)$ depend on $x$ and $z$ when $\varepsilon \approx 0$? What if $\varepsilon$ is very large?

*[If we want to predict $Y$, and we have $x$ and $z$ available, should we use $\Pr_Y(y \mid X = x, Z = z)$, or $\Pr_Y(y \mid X = x)$, or $\Pr_Y(y \mid Z = z)$? The obvious answer is that we should use the first, since it uses all available data.*
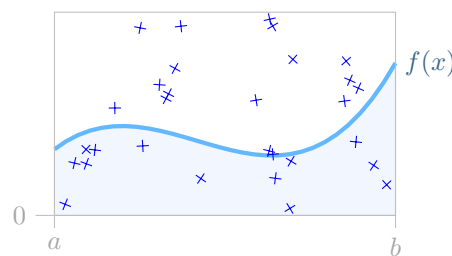
*[But suppose we're interested in predicting $Y$, and we've trained a predictor on $(x, y, z)$ data generated according to the code above, but in deployment the data comes from a slightly different model – which of the three predictors is robust to this change in environment? If the first line of code is different for the new data environment, then the first and second predictors still work correctly. If the second line of code is different, then all bets are off. If the third line of code is different, only the second predictor still works. So, for robust prediction, we might prefer the second predictor. It's called the 'causal predictor' since it only uses the input variable that directly causes the response we're interested in.*

*[The challenge is that, in typical machine learning tasks, we don't know which of our predictor variables are causal and which aren't.]*

**Question 11.** Suppose we're given a function $f(x) \geq 0$ and we want to evaluate

$$\int_{x=a}^{b} f(x) \, dx.$$

Here's an approximation method: (i) draw a box that contains $f(x)$ over the range $x \in [a, b]$, (ii) scatter points uniformly at random in this box, (iii) return $A \times p$ where $A$ is the area of the box and $p$ is the fraction of points that are under the curve. Explain why this is a special case of Monte Carlo integration.



*Do NOT give a wishy-washy qualitative argument along the lines of "there are random points, and we're evaluating an integral, so it's a type of Monte Carlo". Monte Carlo has a precise meaning: $\mathbb{E}\, h(X) \approx n^{-1} \sum_i h(x_i)$. In your answer you should (a) explain the random variable in question, (b) specify the h function, (c) give an explanation along the lines of section 5.1 of lecture notes.*

**Question 12 (Sequential Bayes).** I have a biased coin, with unknown probability of heads $\theta$. I toss it $n$ times, with outcomes $x_1, x_2, \ldots, x_n$ where $x_n = 1$ indicates heads and $x_n = 0$ indicates tails. My prior belief is $\Theta \sim \text{Uniform}[0, 1]$. Here are two approaches to applying Bayes's rule:

- *One-shot Bayes.* Use Bayes's rule to compute the posterior of $\Theta$, given data $(x_1, \ldots, x_n)$, using prior $\Theta \sim \text{Uniform}[0, 1]$, and assumimg that coin tosses are independent.

- *Sequential Bayes.* Use Bayes's rule to compute the posterior of $\Theta$ given data $x_1$, using the uniform prior; let the posterior density be $p_1(\theta)$. Apply Bayes's rule again to compute the posterior of $\Theta$ given data $x_2$, but this time using $p_1(\theta)$ as the prior; let the posterior density be $p_2(\theta)$. Continue applying Bayes's rule in this way, until we have found $p_n(\theta)$.

State the posterior distribution found by one-shot Bayes. Prove by induction on $n$ that sequential Bayes gives the same answer.

*Sequential Bayes and one-shot Bayes give the same answer for any inference problem, not just this coin-tossing example. Can you prove the general case?*
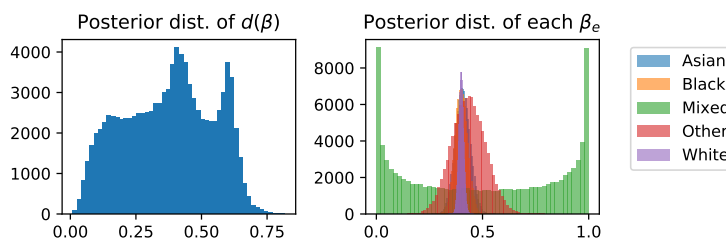
**Question 13.** In the setting of question 7, I wish to measure the amount of police bias. Given a 5-tuple of parameters $\beta = (\beta_{\mathrm{As}}, \beta_{\mathrm{Blk}}, \beta_{\mathrm{Mix}}, \beta_{\mathrm{Oth}}, \beta_{\mathrm{Wh}})$, I define the overall bias score to be

$$d(\beta) = \max_{e,e'}\left|\beta_e - \beta_{e'}\right|.$$

If $d(\beta)$ is large, then there is *some* pair of ethnicities with very unequal treatment.

As a Bayesian I view $\beta$ as a random variable taking values in $[0,1]^5$, therefore $d(\beta)$ is a random variable also. To investigate its distribution, I sample $\beta$ from the posterior distribution that I found in question 7, I compute $d(\beta)$, and I plot a histogram. The output, shown on the left, is bizarre. To help me understand what's going on, I plot histograms of each of the individual $\beta_e$ coefficients, shown on the right.

Explain the results. *[Hint. Explore the Beta distribution numerically. For what parameters does it have a bimodal distribution? What are the posterior distributions in this question?]*



Posterior dist. of $d(\beta)$ — Posterior dist. of each $\beta_e$

**Question 14.** Consider the outlier model from question 6. How likely is it that the datapoint with value 9 is an outlier? *[Hint. Treat this as a two-parameter problem, like question 8.]*

**Question 15.** I have a coin, which might be biased. I toss it $n$ times and get $x$ heads.

I am uncertain whether or not the coin is biased. Let $m \in \{\mathrm{fair}, \mathrm{biased}\}$ indicate which of the two cases is correct; and if it is biased let $\theta$ be the probability of heads. The probabilty of observing $x$ heads is thus

$$\Pr(x \mid m, \theta) = \begin{cases} \binom{n}{x}\theta^x(1-\theta)^{n-x} & \text{if } m = \text{biased} \\ \binom{n}{x}(1/2)^x(1-1/2)^{n-x} & \text{if } m = \text{unbiased} \end{cases}$$

As a Bayesian I shall represent my uncertainty about $m$ with a prior distribution, $\Pr_M(\mathrm{fair}) = p$, $\Pr_M(\mathrm{biased}) = 1 - p$. If it is biased, my prior belief is that the probability of heads is $\Theta \sim \mathrm{Uniform}[0,1]$.

(a)   Write down the prior distribution for the pair $(M, \Theta)$, assuming independence as usual.

(b)   Find the posterior distribution of $(M, \Theta)$ given $x$.

(c)   Find $\mathbb{P}(M = \text{unbiased} \mid x)$, i.e. the posterior probability that the coin is unbiased.

*This is a Bayesian question, and it's answered in the same way as any other Bayesian question: write down the prior density $\Pr_{M,\Theta}(m,\theta)$, write down the data density $\Pr(x \mid m,\theta)$, and multiply them together (times a constant factor) to get the posterior $\Pr_{M,\Theta}(m,\theta \mid x)$. To keep track of all the cases, it may be helpful to use indicator functions, both for $\Pr_M$ and for $\Pr(x \mid m,\theta)$.*

*Part (c) is about nuisance parameters, as in exercise 7.4 in lecture notes (look at the mathematical solution of that exercise). Once we've found the posterior density, say $\Pr_{M,\Theta}(m,\theta) = \kappa f(m,\theta)$ where $\kappa$ is the normalizing constant, we have to integrate out $\theta$ to find the marginal distribution, as in exercise 7.4:*

$$\mathbb{P}(M = \text{fair} \mid x) = \int_\theta \kappa f(\text{fair}, \theta)\, d\theta \qquad \mathbb{P}(M = \text{biased} \mid x) = \int_\theta \kappa f(\text{biased}, \theta)\, d\theta.$$

*Then solve for $\kappa$, using the "densities sum to one" rule, as in exercise 7.5 from lecture notes.*

*This question is an illustration of Bayesian model selection, which you can read about in section 7.4 of lecture notes.*

**Question 16.** (a)   Suppose we have a single observation $x$, drawn from $\text{Normal}(\mu + \nu, \sigma^2)$, where $\mu$ and $\nu$ are unknown parameters, and $\sigma^2$ is known. Explain why the maximum likelihood estimates for $\mu$ and $\nu$ are non-identifiable.

(b)   For $\mu$ use $\text{Normal}(\mu_0, \rho_0^2)$ as prior, and for $\nu$ use $\text{Normal}(\nu_0, \rho_0^2)$, where $\mu_0$, $\nu_0$, and $\rho_0$ are known. Find the posterior density of $(\mu, \nu)$. Calculate the parameter values $(\hat{\mu}, \hat{\nu})$ where the posterior density is maximum. (These are called *maximum a posteriori estimates* or *MAP estimates*.)

(c)   An engineer friend tells you "Bayesianism is the Apple of inference. You just work out the posterior, and everything Just Works™, and you don't need to worry about irritating things like non-identifiability." What do you think?

**Question 17.**   Here's my answer to question 1:

```
1   k = np.random.choice(4, p=[.6,.3,.05,.05], size=n)
2   t = np.random.uniform(size=n)
3   x = np.column_stack([np.sinπ(2**t), 0.55*np.sinπ(2**(0.4*t+0.3)), −0.3*np.ones(n), 0.3*np.ones(n)])
4   y = np.column_stack([np.cosπ(2**t), 0.55*np.cosπ(2**(0.4*t+0.3)), 0.3*np.ones(n), 0.3*np.ones(n)])
5   xy = np.column_stack([x[np.arange(n), k], y[np.arange(n), k]])
6   xy = np.random.normal(loc=xy, scale=.08)
```

Compute the distribution of $(X \mid Y = 0.3)$. Give your answer as a histogram.

*You will need to derive your own method for sampling, along the lines of the derivation of computational Bayes in section 5.2. The difference here is that instead of using Bayes's rule*

$$\text{Pr}_X(x \mid Y = y) = \kappa \, \text{Pr}_{X,Y}(x, y) = \kappa \, \text{Pr}_X(x) \, \text{Pr}_Y(y \mid X = x)$$

*you will need to use a version more suited to the generation method used here,*

$$Pr_{X,Y}(x, y) = \sum_k \int_t \text{Pr}(x, y, k, t)\, dt = \sum_k \int_t \text{Pr}_K(k)\, \text{Pr}_T(t)\, \text{Pr}_X(x \mid k, t)\, \text{Pr}_Y(y \mid k, t)\, dt \,.$$

*You should end up with a Monte Carlo integration that uses $(K, T, X)$ samples.*