

## CSE 515T (Spring 2015) Midterm solutions

1. Consider two coins with unknown bias  $\theta_1$  and  $\theta_2$ , respectively. We place independent, identical beta priors on these quantities:

$$p(\theta_1) = \mathcal{B}(\theta_1; 2, 2); \quad p(\theta_2) = \mathcal{B}(\theta_2; 2, 2).$$

Imagine someone flips both coins and tells you that *exactly one* of the outcomes (but not which) was a “head.” Thus the observation was either HT or TH, but you are not told which. The below expressions are conditioned on “H” to indicate this observation.

- Give an expression for the posterior of the first coin’s bias given this observation,  $p(\theta_1 \mid \text{H})$ . Simplify the result as much as you can. Plot the prior and the posterior for  $\theta_1$  over the interval  $\theta_1 \in (0, 1)$ .
- Give an expression for the joint posterior  $p(\theta_1, \theta_2 \mid \text{H})$ . Plot the joint prior, the likelihood, and the joint posterior as three separate heat maps over the unit square  $(\theta_1, \theta_2) \in (0, 1)^2$ . Use a grid with at least 100 values along each of the two  $\theta$  axes.
- Summarize what the observation taught us about the bias of the coins.

### Solution

The outcome of the unseen experiment was either HT or TH. These are mutually exhaustive and independent events, so we may use the sum rule to derive the desired posterior:

$$p(\theta_1 \mid \text{H}) = \Pr(\text{HT})p(\theta_1 \mid \text{HT}) + \Pr(\text{TH})p(\theta_1 \mid \text{TH}).$$

Notice that both  $p(\theta_1 \mid \text{HT})$  and  $p(\theta_1 \mid \text{TH})$  can be computed explicitly as updated beta distributions, now that the coins in the outcomes have been identified:

$$\begin{aligned} p(\theta_1 \mid \text{HT}) &= \mathcal{B}(\theta_1; 3, 2) \\ p(\theta_1 \mid \text{TH}) &= \mathcal{B}(\theta_1; 2, 3). \end{aligned}$$

What is  $\Pr(\text{HT})$ ? It can be calculated explicitly, but a simpler approach is to appeal to symmetry between the coins to conclude  $\Pr(\text{HT}) = \Pr(\text{TH}) = 1/2$ . Thus

$$p(\theta_1 \mid \text{H}) = \frac{1}{2}(p(\theta_1 \mid \text{HT}) + p(\theta_1 \mid \text{TH})) = \frac{1}{2}(\mathcal{B}(\theta_1; 3, 2) + \mathcal{B}(\theta_1; 2, 3)).$$

We can simplify this expression further. The posterior is proportional to

$$p(\theta_1 \mid \text{H}) \propto \theta_1^2(1 - \theta_1) + \theta_1(1 - \theta_1)^2 = \theta_1(1 - \theta_1) \propto \mathcal{B}(\theta_1; 2, 2).$$

Therefore the distribution of  $\theta_1$  has not changed given our observation! The prior (and posterior!) for  $\theta_1$  is plotted in Figure 1.

There are two ways to compute the joint posterior over  $(\theta_1, \theta_2)$ : the easy way and the hard way. The easy way is to use the sum rule again to write

$$\begin{aligned} p(\theta_1, \theta_2 \mid \text{H}) &= \Pr(\text{HT})p(\theta_1, \theta_2 \mid \text{HT}) + \Pr(\text{TH})p(\theta_1, \theta_2 \mid \text{TH}) \\ &= \frac{1}{2}\mathcal{B}(\theta_1; 3, 2)\mathcal{B}(\theta_2; 2, 3) + \frac{1}{2}\mathcal{B}(\theta_1; 2, 3)\mathcal{B}(\theta_2; 3, 2). \end{aligned}$$

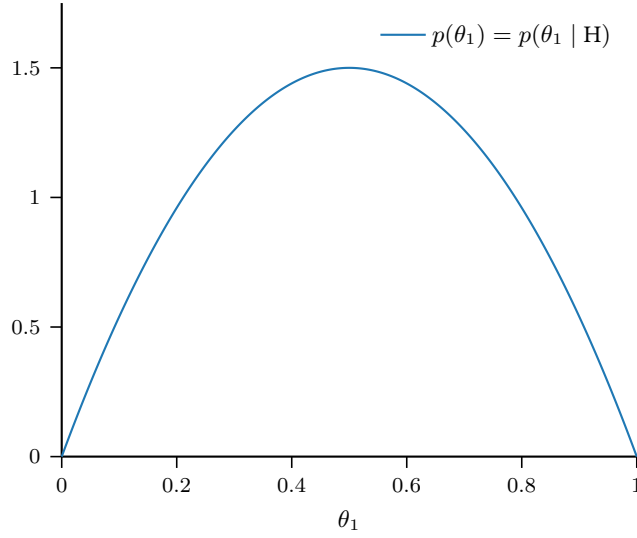


Figure 1: The prior and posterior (given the observation H) of  $\theta_1$ .

If we go with the hard way, we begin by computing the joint prior. By independence, we have:

$$p(\theta_1, \theta_2) = \mathcal{B}(\theta_1; 2, 2)\mathcal{B}(\theta_2; 2, 2).$$

To derive the likelihood, we again note that our observation could have been generated by two mutually independent events: HT or TH. Given  $\theta_1$  and  $\theta_2$ , the total probability of these events is:

$$\Pr(H | \theta_1, \theta_2) = \theta_1(1 - \theta_2) + (1 - \theta_1)\theta_2;$$

the first term accounts for HT and the second term for TH. The posterior is now:

$$\begin{aligned} p(\theta_1, \theta_2 | H) &= \frac{1}{Z} \Pr(H | \theta_1, \theta_2) p(\theta_1, \theta_2) \\ &= \frac{1}{Z} (\theta_1(1 - \theta_2) + (1 - \theta_1)\theta_2) \mathcal{B}(\theta_1; 2, 2) \mathcal{B}(\theta_2; 2, 2). \end{aligned}$$

The normalization constant is

$$Z = \Pr(H) = \int_0^1 \int_0^1 (\theta_1(1 - \theta_2) + (1 - \theta_1)\theta_2) \mathcal{B}(\theta_1; 2, 2) \mathcal{B}(\theta_2; 2, 2) d\theta_1 d\theta_2.$$

This integral is tractable and equals  $1/2$ .<sup>1</sup> In fact, this is always true for any arbitrary mean- $1/2$  beta priors on  $\theta_1, \theta_2$ : if our best guess is that each coin is unbiased, then the outcomes HT/TH always have equal combined probability as the outcomes HH/TT.

The posterior is now

$$p(\theta_1, \theta_2 | H) = 2(\theta_1(1 - \theta_2) + (1 - \theta_1)\theta_2) \mathcal{B}(\theta_1; 2, 2) \mathcal{B}(\theta_2; 2, 2).$$

This is equivalent to the expression we derived with “the easy way.”

---

<sup>1</sup><http://goo.gl/wRofuX>

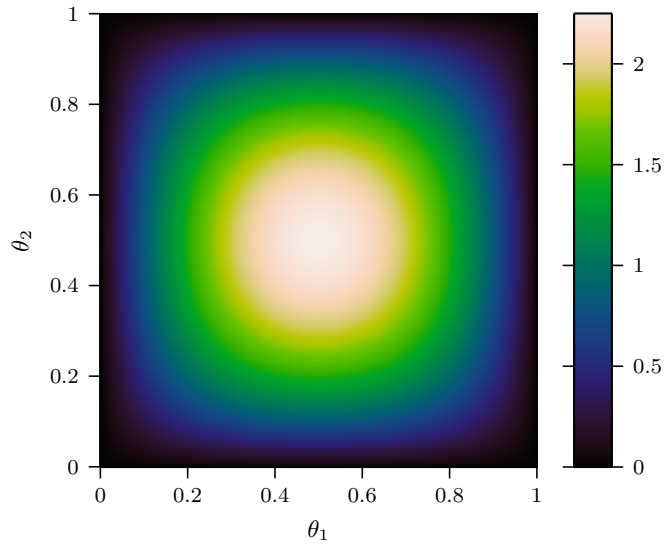


Figure 2: The joint prior  $p(\theta_1, \theta_2)$ .

The prior, likelihood, and posterior are plotted below. From the posterior, we can see that joint probabilities corresponding to jointly low or high values: these combinations would correspond to a higher probability of seeing either HH or TT observations. Despite the marginals for  $\theta_1$  and  $\theta_2$  remaining unchanged, the H observation has entangled the previously independent beliefs in the anticorrelated posterior.

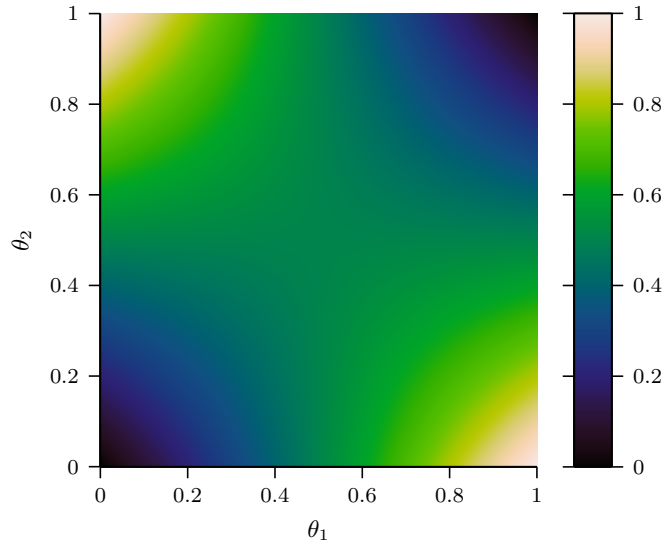


Figure 3: The joint likelihood  $p(H \mid \theta_1, \theta_2)$ .

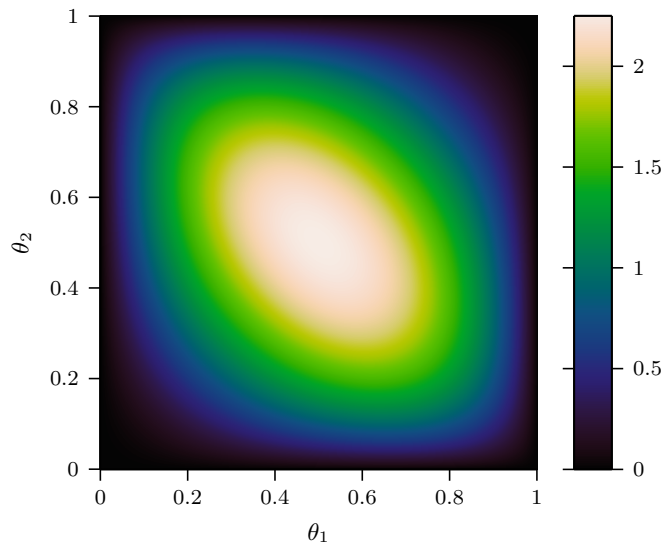


Figure 4: The joint posterior  $p(\theta_1, \theta_2 \mid H)$ .

2. Consider the three-dimensional parameter vector  $\boldsymbol{\theta} = [\theta_1, \theta_2, \theta_3]^\top$ , with the following joint multivariate Gaussian prior:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}\left(\begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}; \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 2 & 0 \\ 2 & 9 & 0 \\ 0 & 0 & 16 \end{bmatrix}\right).$$

We are going to consider a decision problem with action space  $\mathcal{A} = \{1, 2, 3\}$ . The result of choosing an action  $a \in \mathcal{A}$  will be to observe the exact value of  $\theta_a$ , the  $a$ th element of  $\boldsymbol{\theta}$ .

Consider the following loss functions,  $\ell_1$  and  $\ell_2$ :

$$\ell_1(\boldsymbol{\theta}, a) = \begin{cases} 1 & \theta_a > 0 \\ 0 & \theta_a \leq 0 \end{cases} \quad \ell_2(\boldsymbol{\theta}, a) = \min(0, \theta_a).$$

For each:

- Write a generic expression for the expected loss of action  $a$  in terms of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . Evaluate any integrals you encounter.
- Give a numerical value for the expected loss of each action, using the values of  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  provided above.
- State the Bayes action.

### Solution

First, we note that the loss functions only depend on  $\boldsymbol{\theta}$  through  $\theta_a$ , so we must only consider the marginal belief about  $\theta_a$  when contemplating action  $a$ . By applying the marginalization formula for multivariate Gaussians, this belief is:

$$p(\theta_a) = \mathcal{N}(\theta_a; \mu_a, \Sigma_{aa}).$$

For loss  $\ell_1$ , we may calculate the expected loss of each action:

$$\mathbb{E}[\ell_1(\boldsymbol{\theta}, a)] = \int_{-\infty}^{\infty} \ell_1(\boldsymbol{\theta}, a) p(\theta_a) d\theta_a = \int_0^{\infty} \mathcal{N}(\theta_a; \mu_a, \Sigma_{aa}) d\theta_a = 1 - \Phi(0; \mu_a, \Sigma_{aa}).$$

Using this result, we may numerically calculate the expected loss for each action:

$$\mathbb{E}[\ell_1(\boldsymbol{\theta}, 1)] = 0.5 \quad \mathbb{E}[\ell_1(\boldsymbol{\theta}, 2)] = 0.631 \quad \mathbb{E}[\ell_1(\boldsymbol{\theta}, 3)] = 0.691.$$

The Bayes action is  $a = 1$ , with the lowest expected loss.

For loss  $\ell_2$ , we proceed in the same way:

$$\mathbb{E}[\ell_2(\boldsymbol{\theta}, a)] = \int_{-\infty}^{\infty} \ell_2(\boldsymbol{\theta}, a) p(\theta_a) d\theta_a = \int_{-\infty}^0 \theta_a \mathcal{N}(\theta_a; \mu_a, \Sigma_{aa}) d\theta_a.$$

We may compute this definite integral; I used a table of Gaussian integrals<sup>2</sup> and the identity

$$\phi\left(\frac{a - \mu}{\sigma}\right) = \sigma \mathcal{N}(a; \mu, \sigma^2)$$

<sup>2</sup>[http://en.wikipedia.org/wiki/List\\_of\\_integrals\\_of\\_Gaussian\\_functions](http://en.wikipedia.org/wiki/List_of_integrals_of_Gaussian_functions)

to derive

$$\mathbb{E}[\ell_2(\boldsymbol{\theta}, a)] = \mu_a \Phi(0; \mu_a, \Sigma_{aa}) - \Sigma_{aa} \mathcal{N}(0; \mu_a, \Sigma_{aa}).$$

Using this result, we may numerically calculate the expected loss for each action:

$$\mathbb{E}[\ell_2(\boldsymbol{\theta}, 1)] = -0.399 \quad \mathbb{E}[\ell_2(\boldsymbol{\theta}, 2)] = -0.763 \quad \mathbb{E}[\ell_2(\boldsymbol{\theta}, 3)] = -0.791.$$

The Bayes action is now  $a = 3$ , with the lowest expected loss.

3. Consider a  $d$ -dimensional vector  $\boldsymbol{\theta}$  with an arbitrary multivariate Gaussian distribution:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

- Give a general expression for the distribution of the following (scalar) value  $\tau$ .

$$\tau = \theta_1 + 2\theta_2 + \cdots d\theta_d$$

- Consider again the specific distribution of the three-dimensional vector  $\boldsymbol{\theta}$  from the last problem, as well as the action space  $\mathcal{A}$  with the same observation mechanism: after choosing  $a \in \mathcal{A}$ , we will observe the corresponding value  $\theta_a$ . Suppose we may select one action and then must predict  $\tau$  under a squared loss function:

$$\ell(\tau, \hat{\tau}) = (\tau - \hat{\tau})^2.$$

Using the distribution from the last problem. what is the expected loss of each of the three available actions? Which is the Bayes action?

### Solution

Define the (row) vector  $\mathbf{d}^\top = [1, 2, \dots, d]$ . We first notice that  $\tau$  is simply a linear transformation of  $\boldsymbol{\theta}$ :

$$\tau = \mathbf{d}^\top \boldsymbol{\theta};$$

therefore  $\tau$  has a multivariate Gaussian distribution:

$$p(\tau) = p(\mathbf{d}^\top \boldsymbol{\theta}) = \mathcal{N}(\tau; \mathbf{d}^\top \boldsymbol{\mu}, \mathbf{d}^\top \boldsymbol{\Sigma} \mathbf{d}).$$

In the second part of the question, we must consider estimating  $\tau$  under a squared loss function  $\ell(\tau, \hat{\tau}) = (\tau - \hat{\tau})^2$ . A general result from Bayesian decision theory is that the Bayes action is to estimate  $\hat{\tau}$  as the (posterior) mean of  $\tau$ . For example, given the initial belief from the last problem, we would predict  $\hat{\tau} = \mathbf{d}^\top \boldsymbol{\mu} = 8$ . What is the *expected* loss when predicting the mean  $\hat{\tau} = \mathbb{E}[\tau]$ ?

$$\mathbb{E}[\ell(\tau, \mathbb{E}[\tau])] = \mathbb{E}[(\tau - \mathbb{E}[\tau])^2] = \text{var}[\tau] = \mathbf{d}^\top \boldsymbol{\Sigma} \mathbf{d}.$$

The expected squared loss is simply the variance of  $\tau$ ! Conveniently, we have a closed-form expression for this variance.

The problem asks us to consider how we would proceed with estimating  $\tau$  if we could observe one of the entries of the vector  $\boldsymbol{\theta}$  before making our prediction  $\hat{\tau}$ . If we wish to minimize our expected loss, we should minimize the variance of  $\tau$  with our observation. Observing an entry of  $\boldsymbol{\theta}$  is a conditioning observation of a multivariate Gaussian. We have a closed-form expression for the posterior covariance of  $\boldsymbol{\theta}$  after observing any entry  $\theta_a$ . Remarkably, the posterior covariance of  $\boldsymbol{\theta}$  does not depend on the actual value we observe, only the index of the entry we choose,  $a$ . The

posterior covariance matrices  $\Sigma_{\theta|\theta_a}$  for each available action  $a \in \mathcal{A}$  are:

$$\begin{aligned}\Sigma_{\theta|\theta_1} &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 16 \end{bmatrix} \\ \Sigma_{\theta|\theta_2} &= \begin{bmatrix} \frac{5}{9} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 16 \end{bmatrix} \\ \Sigma_{\theta|\theta_3} &= \begin{bmatrix} 1 & 2 & 0 \\ 2 & 9 & 0 \\ 0 & 0 & 0 \end{bmatrix}.\end{aligned}$$

The expected losses of our final prediction of  $\hat{\tau}$  given  $\theta_a$  is now given by

$$\mathbb{E}[\ell(\tau, \hat{\tau}) \mid \theta_a] = \mathbf{d}^\top \Sigma_{\theta|\theta_a} \mathbf{d}.$$

For our particular problem, the expected final loss after each potential action is:

$$\mathbb{E}[\ell(\tau, \hat{\tau}) \mid \theta_1] = 164 \quad \mathbb{E}[\ell(\tau, \hat{\tau}) \mid \theta_2] = 144 \frac{5}{9} \quad \mathbb{E}[\ell(\tau, \hat{\tau}) \mid \theta_3] = 45.$$

The Bayes action is  $a = 3$ . Despite the fact that  $\theta_3$  is uncorrelated with the other two entries, collapsing its large variance from 16 to zero has the effect of reducing the variance of  $\tau$  (and therefore our expected loss) by  $3^2 \cdot 16 = 144$ .



4. Consider the following data:

$$\mathbf{x} = [0.54, 1.84, -2.26, 0.86, 0.32]^\top;$$

$$\mathbf{y} = [-1.31, -0.43, 0.34, 3.58, 2.77]^\top.$$

Consider the Bayesian linear regression model with  $\phi(x) = [1, x]^\top$ . Use the prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \mathbf{I})$ .

Plot the posterior probability that the slope of the regression line is positive as a function of the standard deviation of the observation noise  $\sigma$  (the noise variance is then  $\sigma^2$ ). Use a grid of at least 100 points in the range  $\sigma \in (0.01, 10)$ .

### Solution

Using the given linear regression model, we assume

$$y = \phi(x)^\top \mathbf{w} + \varepsilon = w_1 + w_2 x + \varepsilon.$$

The second entry of the weight vector  $\mathbf{w}$ ,  $w_2$ , therefore serves as the slope of the regression line.

The Bayesian linear regression model gives the following posterior for  $\mathbf{w}$  given observations  $\mathcal{D}$  and a specified noise variance  $\sigma^2$ :

$$p(\mathbf{w} \mid \mathcal{D}, \sigma^2) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_{\mathbf{w}|\mathcal{D}}, \boldsymbol{\Sigma}_{\mathbf{w}|\mathcal{D}}),$$

where

$$\boldsymbol{\mu}_{\mathbf{w}|\mathcal{D}} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{y};$$

$$\boldsymbol{\Sigma}_{\mathbf{w}|\mathcal{D}} = \mathbf{I} - \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{X},$$

where we have plugged the given prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \mathbf{I})$  into the general result.

Given a value of  $\sigma$ , the formulas above give the posterior over  $\mathbf{w}$ . To determine the probability that  $w_2$  is positive, we simply take the marginal posterior distribution and evaluate the normal CDF:

$$\Pr(w_2 > 0 \mid \mathcal{D}, \sigma^2) = 1 - \Phi(0; (\boldsymbol{\mu}_{\mathbf{w}|\mathcal{D}})_2, (\boldsymbol{\Sigma}_{\mathbf{w}|\mathcal{D}})_{22}).$$

This quantity is plotted as a function of  $\sigma$  in Figure 5. The larger the noise, the less confident we become about the sign of the slope.

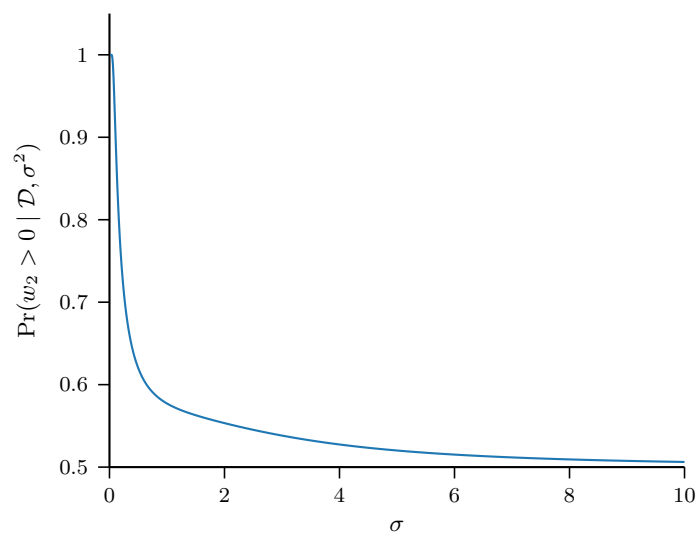


Figure 5: The posterior probability that the slope of the regression line is positive as a function of the noise standard deviation  $\sigma$ . Note the limits on the  $y$ -axis.