

Linear Modeling - Bayesian Approach

TF4063

Fadjar Fathurrahman

The material in this note is based on Rogers2017.

1 Coin game: introduction

Coin game: The stall owner tosses a coin ten times for each customer. If the coin lands heads on six or fewer occasions, the customer wins back their \$ 1 stake plus an additional \$ 1. Seven or more and the stall owner keeps their money. The binomial distribution describes the probability of a certain number of successes (heads) in N binary events.

The probability of y heads from N tosses where each toss lands heads with probability r is given by:

$$P(Y = y) = \binom{N}{y} r^y (1 - r)^{N-y} \quad (1)$$

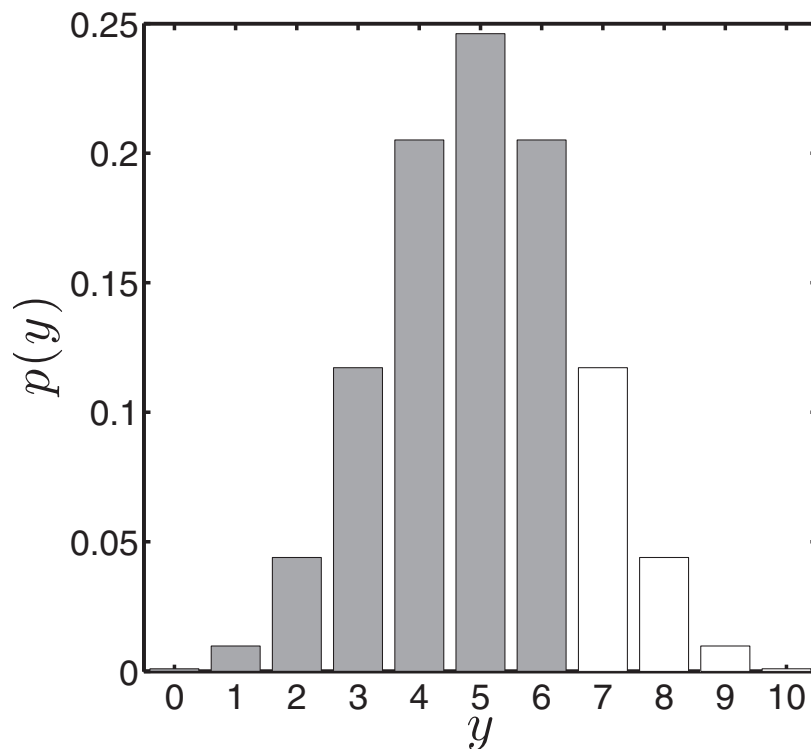


Figure 1: The binomial density function when $N = 10$ and $r = 0.5$. Regions for $y_N \leq 6$ are shaded.

Assume the coin is fair: $r = 0.5$. Probability of winning the game:

$$\begin{aligned} P(Y \leq 6) &= 1 - P(Y > 6) \\ &= 1 - (P(Y = 7) + P(Y = 8) + P(Y = 9) + P(Y = 10)) \\ &= 0.8281 \end{aligned}$$

This seems like a pretty good game - you'll double your money with probability 0.8281.

It is also possible to compute the expected return from playing the game. Recall that the expectation value of $f(X)$ computed for probability distribution $P(x)$ is:

$$\mathbb{E}_{P(x)} \{f(X)\} = \sum_x f(x)P(x)$$

where the summation is over all possible values that the random variable can take. Let X be the random variable that takes value of 1 if we win and 0 if we lose, so:

$$P(X = 1) = P(Y \leq 6)$$

If we win, $X = 1$, we get a return of \$ 2, so $f(1) = 2$. If we lose, we get a return of nothing: $f(0) = 0$. Hence our expected return is:

$$f(1)P(X = 1) + f(0)P(X = 0) = 2 \times P(Y \leq 6) + 0 \times P(Y > 6) = 1.6562$$

Given that it cost \$ 1 to play, you win, on average, \$ 1.6562 - \$ 1 per game. If you play 100 times, we expect to walk away with a profit of \$ 65.62.

Given these odds of success, it seems sensible to play. However, whilst waiting you notice that the stall owner looks reasonably wealthy and very few customers seem to be winning. Perhaps the assumptions underlying the calculations are wrong. These assumptions are:

1. The number of heads can be modelled as a random variable with a binomial distribution, and the probability of a head on any particular toss is r .
2. The coin is fair - the probability of heads is the same as the probability of tails, $r = 0.5$.

It seems hard to reject the binomial distribution - events are taking place with only two possible outcomes and the tosses do seem to be

independent. This leaves r , the probability that the coin lands heads. Our assumption was that the coin was fair the probability of heads was equal to the probability of tails. Maybe this is not the case? To investigate this, we can treat r as a parameter (like \mathbf{w} and σ^2 last week) and fit it to some data.

Counting heads: maximum likelihood

Suppose that from previous play we obtain the following results:

H, T, H, H, H, H, H, H, H, H

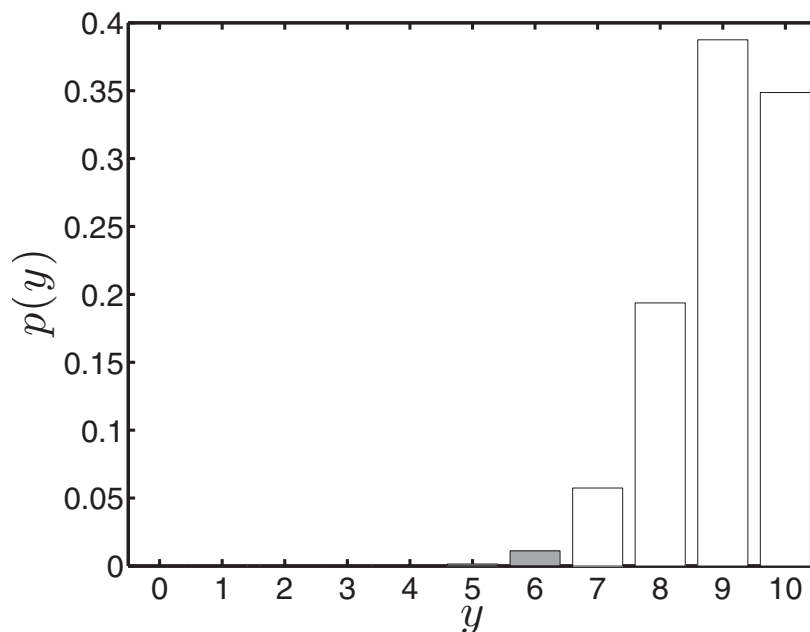


Figure 2: The binomial density function when $N = 10$ and $r = 0.9$.

It is possible to compute the maximum likelihood value of r as follows. The likelihood is given by the binomial distribution:

$$P(Y = y|r, N) = \binom{N}{y} r^y (1-r)^{N-y}$$

Take natural logarithm:

$$L = \log P(Y = y|r, N) = \log \binom{N}{y} + y \log r + (N-y) \log(1-r)$$

Differentiate this expression, equate to zero, and solve for the maxi-

maximum likelihood estimate of the parameter:

$$\begin{aligned}\frac{\partial L}{\partial r} &= \frac{y}{r} - \frac{N-y}{1-r} = 0 \\ y(1-r) &= r(N-y) \\ y &= rN \\ r &= \frac{y}{N}\end{aligned}$$

We obtain, for our data, $r = 0.9$. The expected return is now:

$$2 \times P(Y \leq 6) + 0 \times P(Y > 6) = 0.0256$$

Given that it costs \$ 1 to play, we expect to make $0.0256 - 1 = -0.9744$ per game - a loss of approximately \$ 0.97. $P(Y \leq 6) = 0.0128$ suggests that only about 1 person in every 100 should win, but this does not seem to be reflected in the number of people who are winning. Although the evidence from this run of coin tosses suggests $r = 0.9$, it seems too biased given that several people have won.

The Bayesian way

The value of r computed in the previous section was based on just ten tosses. Given the random nature of the coin toss, if we observed several sequences of tosses it is likely that we would get a different r each time. Thought about this way, r feels a bit like a random variable, R . Maybe we can learn something about the *distribution* of R rather than try and find a particular value. We saw in the previous section that obtaining an exact value by counting is heavily influenced by the particular tosses in the short sequence. No matter how many such sequences we observe there will always be some uncertainty in r - considering it as a random variable with an associated distribution will help us measure and understand this uncertainty. In particular, defining the random variable Y_N to be the number of heads obtained in N tosses, we would like the distribution of r conditioned on the value of Y_N :

$$p(r|y_N)$$

Given this distribution, it would be possible to compute the expected probability of winning by taking the expectation of $P(Y_{\text{new}} \leq 6|r)$

with respect to $p(r|y_N)$:

$$P(Y_{\text{new}} \leq 6|y_N) = \int P(Y_{\text{new}} \leq 6|r)p(r|y_N) dr$$

where Y_{new} is a random variable describing the number of heads in a future set of ten tosses.

Bayes' rule allows us to reverse the conditioning of two (or more) random variables, e.g. compute $p(a|b)$ from $p(b|a)$. Here we're interested in $p(r|y_N)$, which, if we reverse the conditioning, is $p(y_N|r)$ - the probability distribution function over the number of heads in N independent tosses where the probability of a head in a single toss is r . This is the binomial distribution function that we can easily compute for any y_N and r . In our context, Bayes' rule is:

$$p(r|y_N) = \frac{P(y_N|r)p(r)}{P(y_N)} \quad (2)$$

The likelihood, $P(y_N|r)$: how likely is it that we would observe our data (in this case, the data is y_N) for a particular value of r (our model)? For our example, this is the binomial distribution. This value will be high if r could have feasibly produced the result y_N and low if the result is very unlikely. For example, Figure 3 shows the likelihood $P(y_N|r)$ as a function of r for two different scenarios. In the first, the data consists of ten tosses ($N = 10$) of which six were heads. In the second, there were $N = 100$ tosses, of which 70 were heads.

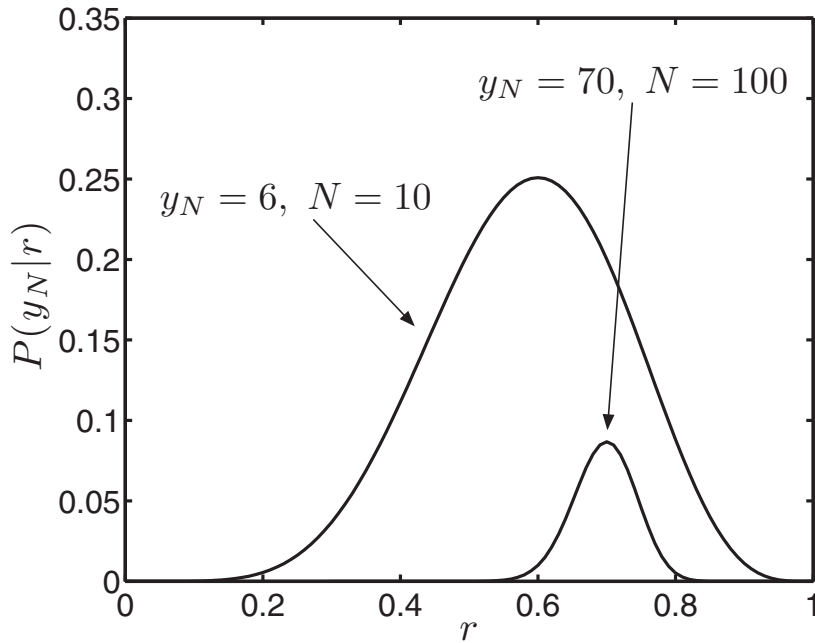


Figure 3: Examples of the likelihood $p(y_N|r)$ as a function of r for two scenarios.

This plot reveals two important properties of the likelihood. Firstly, it is not a probability density. If it were, the area under both curves would have to equal 1. We can see that this is not the case without working out the area because the two areas are completely different. Secondly, the two examples differ in how much they appear to tell us about r . In the first example, the likelihood has a nonzero value for a large range of possible r values (approximately $0.2 \leq r \leq 0.9$). In the second, this range is greatly reduced (approximately $0.6 \leq r \leq 0.8$). This is very intuitive: in the second example, we have much more data (the results of 100 tosses rather than 10) and so we should know more about r .

The prior distribution, $p(r)$: allows us to express any belief we have in the value of r before we see any data. To illustrate this, we shall consider the following three examples:

1. We do not know anything about tossing coins or the stall owner.
2. We think the coin (and hence the stall owner) is fair.
3. We think the coin (and hence the stall owner) is biased to give more heads than tails.

We can encode each of these beliefs as different prior distributions. r can take any value between 0 and 1 and therefore it must be modelled as a *continuous random variable*. Figure 4 shows three density functions that might be used to encode our three different prior beliefs.

Belief number 1 is represented as a uniform density between 0 and 1 and as such shows no preference for any particular r value. Number 2 is given a density function that is concentrated around $r = 0.5$, the value we would expect for a fair coin. The density suggests that we do not expect much variance in r : it's almost certainly going to lie between 0.4 and 0.6. Most coins that any of us have tossed agree with this. Finally, number 3 encapsulates our belief that the coin (and therefore the stall owner) is biased. This density suggests that $r > 0.5$ and that there is a high level of variance. This is fine because our belief is just that the coin is biased: we don't really have any idea how biased at this stage.

The three functions shown in Figure 4 have not been plucked from thin air. They are all examples of beta probability density functions. The beta density function is used for continuous random variables

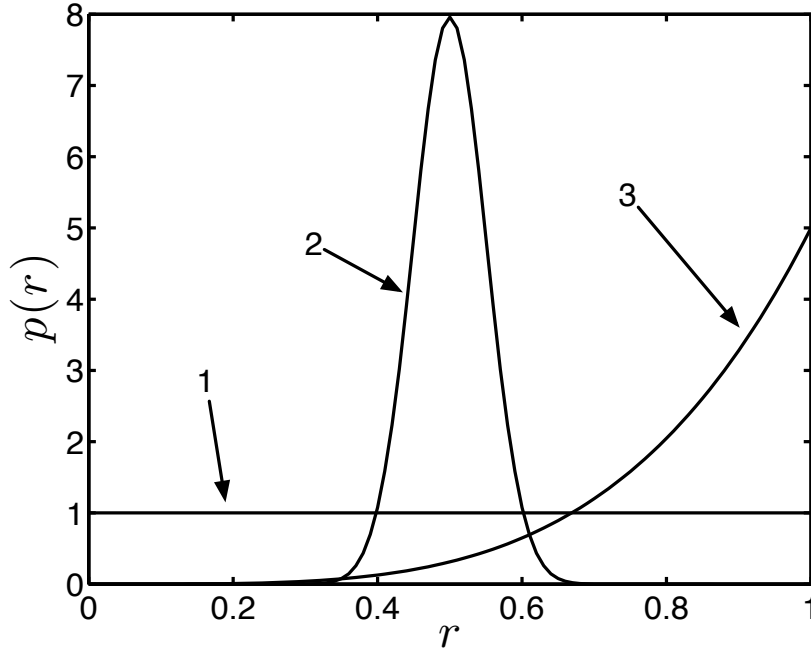


Figure 4: Examples of prior densities for r , $p(r)$, for three different scenarios.

constrained to lie between 0 and 1 - perfect for our example. For a random variable R with parameters α and β , it is defined as:

$$p(r) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} r^{\alpha-1} (1-r)^{\beta-1} \quad (3)$$

$\Gamma(a)$ is known as the gamma function. In Equation (3) the gamma functions ensure that the density is normalized (that is, it integrates to 1 and is therefore a probability density function). In particular

$$\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} = \int_{r=0}^{r=1} r^{\alpha-1} (1-r)^{\beta-1} dr$$

ensuring that:

$$\int_{r=0}^{r=1} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} r^{\alpha-1} (1-r)^{\beta-1} dr = 1$$

The two parameters α and β control the shape of the resulting density function and must both be positive. Our three beliefs as plotted in Figure 4 correspond to the following pairs of parameter values:

1. Know nothing: $\alpha = 1, \beta = 1$.
2. Fair coin: $\alpha = 50, \beta = 50$.
3. Biased : $\alpha = 5, \beta = 1$.

The problem of choosing these values is a big one. For example, why should we choose $\alpha = 5, \beta = 1$ for a biased coin? There is no easy answer to this. We shall see later that, for the beta distribution, they

can be interpreted as a number of previous, hypothetical coin tosses. For other distributions no such analogy is possible and we will also introduce the idea that maybe these too should be treated as random variables. In the mean time, we will assume that these values are sensible and move on.

The marginal distribution of y_N - $P(y_N)$ The third quantity in our equation, $P(y_N)$, acts as a normalizing constant to ensure that $p(r|y_N)$ is a properly defined density. It is known as the *marginal distribution* of y_N because it is computed by integrating r out of the joint density $p(y_N, r)$:

$$P(y_N) = \int_{r=0}^{r=1} p(y_N|r) dr$$

This joint density can be factorized to give:

$$P(y_N) = \int_{r=0}^{r=1} P(y_N|r)p(r) dr$$

which is the product of the prior and likelihood integrated over the range of values that r may take.

$p(y_N)$ is also known as the **marginal likelihood**, as it is the likelihood of the data, y_N , averaged over all parameter values. We shall see in later that it can be a useful quantity in model selection, but, unfortunately, *in all but a small minority of cases, it is very difficult to calculate.*

The posterior distribution - $p(r|y_N)$ This posterior is the distribution in which we are interested. It is the result of updating our prior belief $p(r)$ in light of new evidence y_N . The shape of the density is interesting - it tells us something about how much information we have about r after combining what we knew beforehand (the prior) and what we've seen (the likelihood). Three hypothetical examples are provided in Figure 5 (these are purely illustrative and do not correspond to the particular likelihood and prior examples shown in Figures 3 and 4). (a) is uniform - combining the likelihood and the prior together has left all values of r equally likely. (b) suggests that r is most likely to be low but could be high. This might be the result of starting with a uniform prior and then observing more tails than heads. Finally, (c) suggests the coin is biased to land heads more often. As it is a density, the posterior tells us not just which values are likely but also provides an indication of the level of uncertainty we still have in r having observed some data.

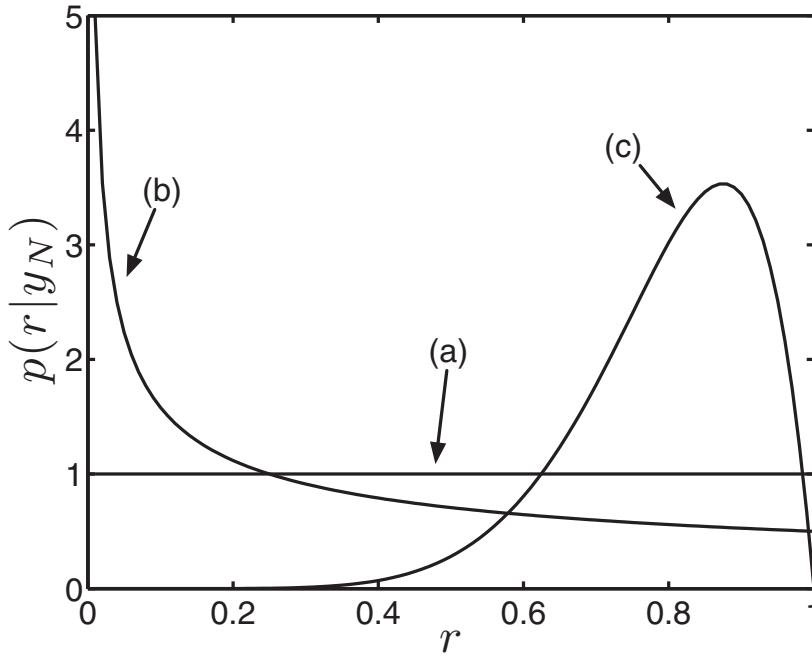


Figure 5: Examples of three possible posterior distributions $p(r|y_N)$.

As already mentioned, we can use the posterior density to compute expectations. For example, we could compute

$$\mathbb{E}_{p(r|y_N)} \{P(Y_{10} \leq 6)\} = \int_{r=0}^{r=1} P(Y_{10} \leq 6) p(r|y_N) dr$$

the expected value of the probability that we will win. This takes into account the data we have observed, our prior beliefs and the uncertainty that remains. It will be useful in helping to decide whether or not to play the game. We will return to this later, but first we will look at the kind of posterior densities we obtain in our coin example.

The exact posterior

The beta distribution is a common choice of prior when the likelihood is a binomial distribution. This is because we can use some algebra to compute the posterior density exactly. In fact, the beta distribution is known as the conjugate prior to the binomial likelihood. If the prior and likelihood are conjugate, the posterior will be of the same form as the prior. Specifically, $p(r|y_N)$ will give a beta distribution with parameters δ and γ , whose values will be computed from the prior and y_N . The beta and binomial are not the only conjugate pair of distributions and we will see an example of another conjugate prior and likelihood pair when we return to the Olympic data later.

Conjugate priors:

A likelihood-prior pair is said to be conjugate if they result in a posterior which is of the same form as the prior. This enables us to compute the posterior density analytically without having to worry about computing the denominator in Bayes' rule, the marginal likelihood.

Prior	Likelihood
Gaussian	Gaussian
beta	binomial
gamma	Gaussian
Dirichlet	multinomial

Using a conjugate prior makes things much easier from a mathematical point of view. However, as we mentioned in both our discussion on loss functions in Chapter 1 and noise distributions in Chapter 2, it is more important to base our choices on modeling assumptions than mathematical convenience. In the next chapter we will see some techniques we can use in the common scenario that the pair are non-conjugate.

Returning to our example, we can omit $p(y_N)$ from Equation 2, leaving:

$$p(r|y_N) \propto P(y_N|r)p(r)$$

Replacing the terms on the right hand side with a binomial and beta distribution gives:

$$p(r|y_N) \propto \left[\binom{N}{y_N} r^{y_N} (1-r)^{N-y_N} \right] \left[\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} r^{\alpha-1} (1-r)^{\beta-1} \right] \quad (4)$$

Because the prior and likelihood are conjugate, we know that $p(r|y_N)$ has to be a beta density. The beta density, with parameters δ and γ , has the following general form:

$$p(r) = Kr^{\delta-1}(1-r)^{\gamma-1}$$

where K is a constant. If we can arrange all of the terms, including r , on the right hand side of Equation (4) into something that looks like $r^{\delta-1}(1-r)^{\gamma-1}$, we can be sure that the constant must also be correct (it has to be $\Gamma(\delta + \gamma) / (\Gamma(\delta)\Gamma(\gamma))$ because we know that the posterior density is a beta density). In other words, we know what the normalizing constant for a beta density is so we do not need to compute $p(y_N)$.

Rearranging Equation (4) gives us:

$$\begin{aligned} p(r|y_N) &\propto \left[\binom{N}{y_N} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \right] \left[r^{y_N} (1-r)^{N-y_N} r^{\alpha-1} (1-r)^{\beta-1} \right] \\ &\propto r^{y_N + \alpha - 1} (1-r)^{N - y_N + \beta - 1} \\ &\propto r^{\delta - 1} (1-r)^{\gamma - 1} \end{aligned}$$

where $\delta = y_N + \alpha$ and $\gamma = N - y_N + \beta$. Therefor the posterior density is:

$$p(r|y_N) = \frac{\Gamma(\alpha + \beta + N)}{\Gamma(\alpha + y_N)\Gamma(\beta + N - y_N)} r^{\alpha + y_N - 1} (1-r)^{\beta + N - y_N - 1} \quad (5)$$

(note that when adding γ and δ , the y_N terms cancel). This is the posterior density of r based on the prior $p(r)$ and the data y_N . Notice how the posterior parameters are computed by adding the number of heads (y_N) to the first prior parameter (α) and the number of tails ($N - y_N$) to the second (β). This allows us to gain some intuition about the prior parameters α and β - they can be thought of as the number of heads and tails in $\alpha + \beta$ previous tosses. For example, consider the second two scenarios discussed in the previous section. For the fair coin scenario, $\alpha = \beta = 50$. This is equivalent to tossing a coin 100 times and obtaining 50 heads and 50 tails. For the biased scenario, $\alpha = 5, \beta = 1$, corresponding to six tosses and five heads. Looking at Figure 4, this helps us explain the differing levels of variability suggested by the two densities: the fair coin density has much lower variability than the biased one because it is the result of many more hypothetical tosses. The more tosses, the more we should know about r .

The analogy is not perfect. For example, α and β don't have to be integers and can be less than 1 (0.3 heads doesn't make much sense). The analogy also breaks down when $\beta = \beta = 1$. Observing one head and one tail means that values of $r = 0$ and $r = 1$ are impossible. However, density 1 in Figure 4, suggests that all values of r are equally likely. Despite these flaws, the analogy will be a useful one to bear in mind as we progress through our analysis.

2 Coin game: three scenarios

We will now investigate the posterior distribution $p(r|y_N)$ for the three different prior scenarios shown in Figure 4: no prior knowl-

edge, a fair coin and a biased coin.

In this scenario, we assume that we know nothing of coin tossing or the stall holder. Our prior parameters are $\alpha = 1, \beta = 1$, shown in Figure 6(a).

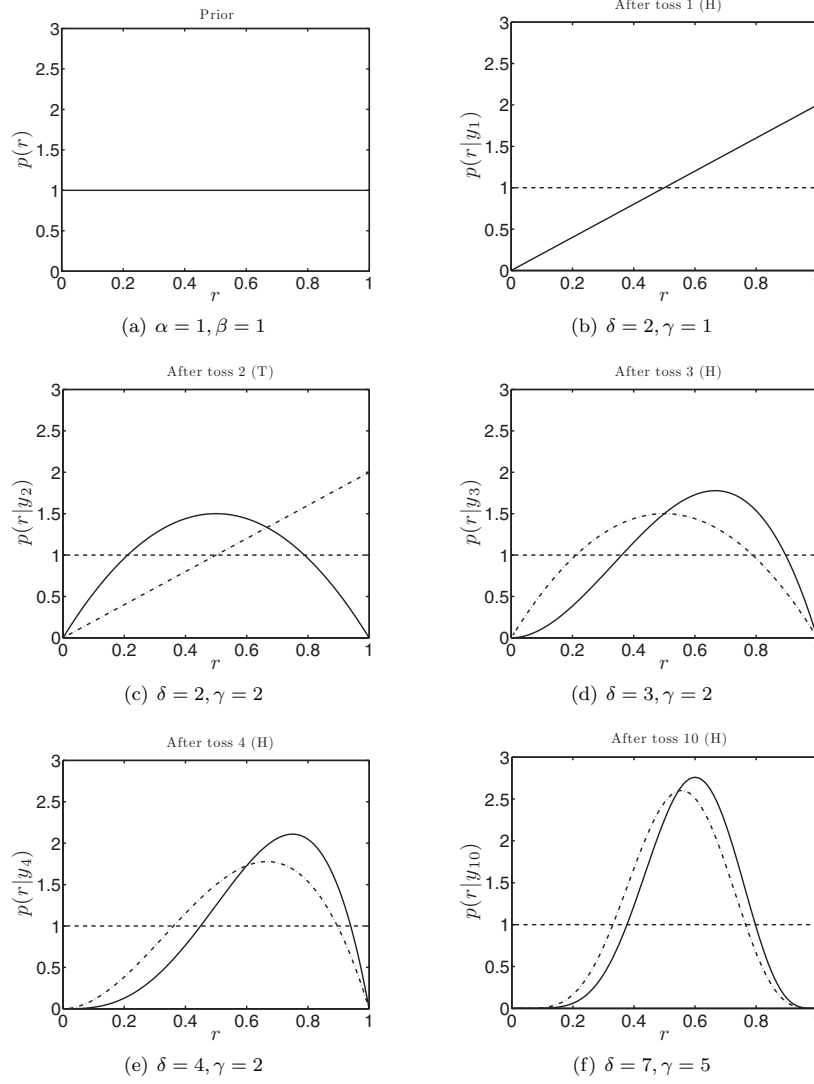


Figure 6: Evolution of $p(r|y_N)$ as the number of observed coin tosses increases.

To compare different scenarios we will use the expected value and variance of r under the prior. The expected value of a random variable from a beta distribution with parameters α and β (the density function of which we will henceforth denote $p(r) = \mathcal{B}(\alpha, \beta)$) is given as:

$$\mathbb{E}_{p(r)}\{R\} = \frac{\alpha}{\alpha + \beta} \quad (6)$$

The variance is given by:

$$\text{var}\{R\} = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)} \quad (7)$$

For scenario 1 ($\alpha = 1, \beta = 1$):

$$\begin{aligned}\mathbb{E}_{p(r)}\{R\} &= \frac{1}{2} \\ \text{var}\{R\} &= \frac{1}{12}\end{aligned}$$

Note that in our formulation of the posterior (Equation (5)) we are not restricted to updating our distribution in blocks of ten - we can incorporate the results of any number of coin tosses. To illustrate the evolution of the posterior, we will look at how it changes toss by toss.

A new customer hands over \$ 1 and the stall owner starts tossing the coin. The first toss results in a head. The posterior distribution after one toss is a beta distribution with parameters $\delta = \alpha + y_N$ and $\gamma = \beta + N - y_N$:

$$p(r|y_N) = \mathcal{B}(\delta, \gamma)$$

In this scenario, $\alpha = \beta = 1$, and as we had $N = 1$ toss and seen $y_N = 1$ head:

$$\begin{aligned}\delta &= \alpha + y_N = 1 + 1 = 2 \\ \gamma &= \beta + N - y_N = 1 + 1 - 1 = 1\end{aligned}$$

This posterior distribution is shown as the solid line in Figure 6(b) (the prior is also shown as a dashed line). This single observation has had quite a large effect - the posterior is very different from the prior. In the prior, all values of r were equally likely. This has now changed - higher values are more likely than lower values with zero density at $r = 0$. This is consistent with the evidence - observing one head makes high values of r slightly more likely and low values slightly less likely. The density is still very broad, as we have observed only one toss. The expected value of r under the posterior is

$$\mathbb{E}_{p(r|y_N)}\{R\} = \frac{\delta}{\delta + \gamma} = \frac{2}{3}$$

and we can see that observing a solitary head has increased the expected value of r from $1/2$ to $2/3$. The variance of the posterior is:

$$\text{var}\{R\} = \frac{\delta\gamma}{(\delta + \gamma)^2(\delta + \gamma + 1)} = \frac{1}{18}$$

which is lower than the prior variance $1/12$. So, the reduction in

variance tells us that we have less uncertainty about the value of r than we did (we have learnt something) and the increase in expected value tells us that what we've learnt is that heads are slightly more likely than tails.

The stall owner tosses the second coin and it lands tails. We have now seen one head and one tail and so $N = 2$, $y_N = 1$, resulting in

$$\begin{aligned}\delta &= \alpha + y_N = 1 + 1 = 2 \\ \gamma &= \beta + N - y_N = 1 + 2 - 1 = 2\end{aligned}$$

The posterior distribution is shown as the solid dark line in Figure 6(c). The lighter dash-dot line is the posterior we saw after one toss and the dashed line is the prior. The density has changed again to reflect the new evidence. As we have now observed a tail, the density at $r = 1$ should be zero and is ($r = 1$ would suggest that the coin always lands heads). The density is now curved rather than straight (as we have already mentioned, the beta density function is very flexible) and observing a tail has made lower values more likely. The expected value and variance are now

$$\begin{aligned}\mathbb{E}_{p(r|y_N)}\{R\} &= \frac{1}{2} \\ \text{var}\{R\} &= \frac{1}{20}\end{aligned}$$

The expected value has decreased back to $1/2$. Given that the expected value under the prior was also $1/2$, you might conclude that we haven't learnt anything. However, the variance has decreased again (from $1/18$ to $1/20$) so we have less uncertainty in r and have learnt something. In fact, we've learnt that r is closer to $1/2$ than we assumed under the prior.

The third toss results in another head. We now have $N = 3$ tosses, $y_N = 2$ heads and $N - y_N = 1$ tail. Our updated posterior parameters are

$$\begin{aligned}\delta &= \alpha + y_N = 1 + 2 = 3 \\ \gamma &= \beta + N - y_N = 1 + 3 - 2 = 2\end{aligned}$$

This posterior is plotted in Figure 6(d). Once again, the posterior is the solid dark line, the previous posterior is the solid light line and

the dashed line is the prior. We notice that the effect of observing this second head is to skew the density to the right, suggesting that heads are more likely than tails. Again, this is entirely consistent with the evidence - we have seen more heads than tails. We have only seen three coins though, so there is still a high level of uncertainty - the density suggests that r could potentially still be pretty much any value between 0 and 1. The new expected value and variance are

$$\begin{aligned}\mathbb{E}_{p(r|y_N)}\{R\} &= \frac{3}{5} \\ \text{var}\{R\} &= \frac{1}{25}\end{aligned}$$

Toss 4 also comes up heads ($y_N = 3, N = 4$), resulting in $\delta = 1 + 3 = 4$ and $\gamma = 1 + 4 - 3 = 2$. Figure 6(e) shows the current and previous posteriors and prior in the now familiar format. The density has once again been skewed to the right - we've now seen three heads and only one tail so it seems likely that r is greater than $1/2$. Also notice the difference between the $N = 3$ posterior and the $N = 4$ posterior for very low values of r - the extra head has left us pretty convinced that r is not 0.1 or lower. The expected value and variance are given by

$$\begin{aligned}\mathbb{E}_{p(r|y_N)}\{R\} &= \frac{2}{3} \\ \text{var}\{R\} &= \frac{2}{63} = 0.0317\end{aligned}$$

where the expected value has increased and the variance has once again decreased.

The remaining six tosses are made so that the complete sequence is

H, T, H, H, H, H, T, T, T, H

a total of six heads and four tails. The posterior distribution after $N = 10$ tosses ($y_N = 6$) has parameters $\delta = 1 + 6 = 7$ and $\gamma = 1 + 10 - 6 = 5$. This (along with the posterior for $N = 9$) is shown in Figure 6(f). The expected value and variance are

$$\begin{aligned}\mathbb{E}_{p(r|y_N)}\{R\} &= 0.5833 \\ \text{var}\{R\} &= 0.0187\end{aligned}$$

Our ten observations have increased the expected value from 0.5 to 0.5833 and decreased our variance from $1/12 = 0.0833$ to 0.0187.

However, this is not the full story. Examining Figure 6(f), we see that we can also be pretty sure that $r > 0.2$ and $r < 0.9$. The uncertainty in the value of r is still quite high because we have only observed ten tosses.

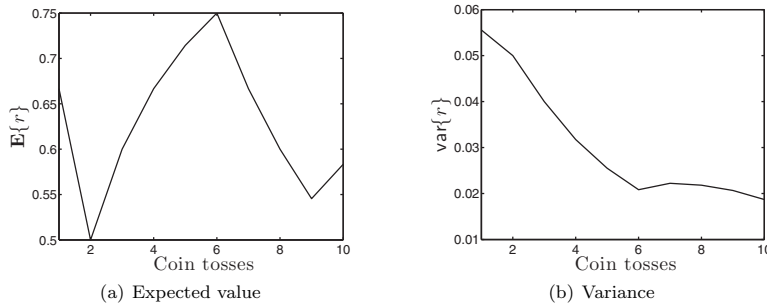


Figure 7: Evolution of expected value (a) and variance (b) of r as coin toss data is added to the posterior.

Figure 7 summarises how the expected value and variance change as the 10 observations are included. The expected value jumps around a bit, whereas the variance steadily decreases as more information becomes available. At the seventh toss, the variance increases. The first seven tosses are

H, T, H, H, H, H, T

The evidence up to and including toss 6 is that heads is much more likely than tails (5 out of 6). Tails on the seventh toss is therefore slightly unexpected. Figure 8 shows the posterior before and after the seventh toss. The arrival of the tail has forced the density to increase the likelihood of low values of r and, in doing so, increased the uncertainty.

Tasks: do the same for scenario 2 and 3.

3 Bayesian approach to Olympic 100m data

Let's revisit our linear model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n \quad (8)$$

where $\mathbf{w} = [w_0, \dots, w_K]^T$, $\mathbf{x}_n = [1, x_n, x_n^2, \dots, x_n^K]$. Combining all inputs into a single matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ we can write:

$$\mathbf{t} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon} \quad (9)$$

As opposed to the maximum likelihood approach that we have used before, we will now treat the model parameters, i.e. \mathbf{w} , as random

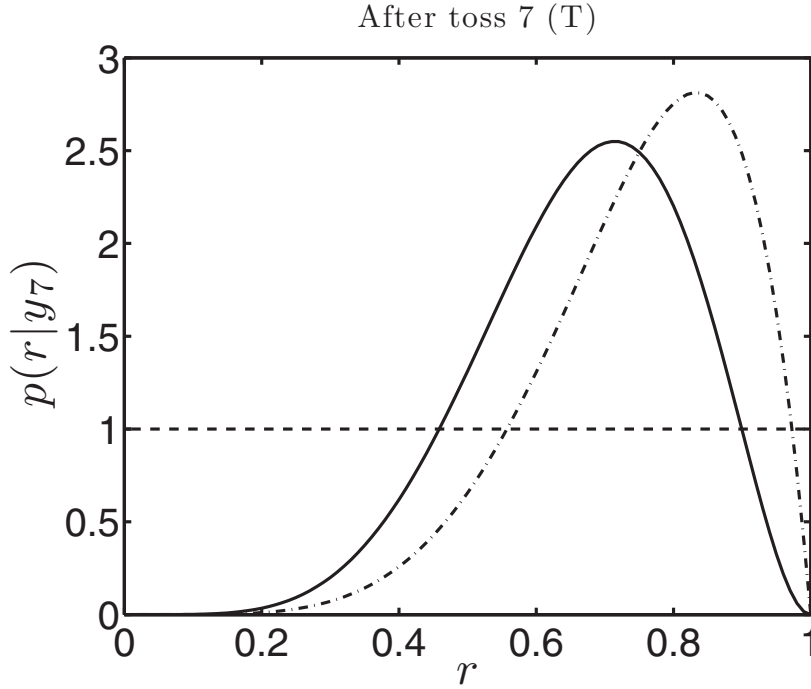


Figure 8: The posterior after six (light) and seven (dark) tosses.

variables. To simplify the analysis, however, we will assume that we know the true value of σ^2 .

The quantity that will be focusing on is

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2, \Delta) \quad (10)$$

which can be calculated using Bayes' rule

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2, \Delta) = \frac{p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2, \Delta)p(\mathbf{w}|\Delta)}{p(\mathbf{t}|\mathbf{X}, \sigma^2, \Delta)} \quad (11)$$

$$= \frac{p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2)p(\mathbf{w}|\Delta)}{p(\mathbf{t}|\mathbf{X}, \sigma^2, \Delta)} \quad (12)$$

where Δ corresponds to some set of parameters required to define the prior over \mathbf{w} that will be defined more precisely later.

Expanding the marginal likelihood:

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2, \Delta) = \frac{p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2)p(\mathbf{w}|\Delta)}{\int p(\mathbf{t}|\mathbf{X}, \sigma^2)p(\mathbf{w}|\Delta) d\mathbf{w}} \quad (13)$$

We are interested in making predictions which will involve taking an expectation w.r.t this posterior density. For a set of attributes \mathbf{x}_{new} corresponding to a new Olympic year, the density over the associated winning time t_{new} is given by

$$p(t_{\text{new}}|\mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t}, \sigma^2, \Delta) = \int p(t_{\text{new}}|\mathbf{x}_{\text{new}}, \mathbf{w}, \sigma^2)p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2, \Delta) d\mathbf{w} \quad (14)$$

Now, let's consider the terms involved in the Bayes' rule expression.

The likelihood

The likelihood $p(\mathbf{t}|\mathbf{X}, \sigma^2)$ is the quantity that we maximized previously, i.e.

$$p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N) \quad (15)$$

The likelihood is N -dimensional Gaussian density with mean $\mathbf{X}\mathbf{w}$ and variance $\sigma^2 \mathbf{I}_N$

The prior

Because we are interested in being able to produce an exact expression for the posterior, we will choose the prior as Gaussian:

$$p(\mathbf{w}|\mu_0, \Sigma_0) = \mathcal{N}(\mu_0, \Sigma_0) \quad (16)$$

where the parameters $\Delta = \{\mu_0, \Sigma_0\}$ will be chosen later.

The posterior

Because the likelihood and the posterior are Gaussians, the posterior will also be a (multivariate) Gaussian. It can be shown that the posterior is

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mu_{\mathbf{w}}, \Sigma_{\mathbf{w}}) \quad (17)$$

$$= \exp \left(-\frac{1}{2} (\mathbf{w} - \mu_{\mathbf{w}})^T \Sigma_{\mathbf{w}}^{-1} (\mathbf{w} - \mu_{\mathbf{w}}) \right) \quad (18)$$

with covariance matrix:

$$\Sigma_{\mathbf{w}} = \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right)^{-1} \quad (19)$$

and mean vector:

$$\mu_{\mathbf{w}} = \Sigma_{\mathbf{w}} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{t} + \Sigma_0^{-1} \mu_0 \right) \quad (20)$$

Given the new observation \mathbf{x}_{new} , we can make prediction as expectation

$$p(t_{\text{new}}|\mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t}, \sigma^2) = \int p(t_{\text{new}}|\mathbf{x}_{\text{new}}, \mathbf{w}, \sigma^2) p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) d\mathbf{w}$$

By our model this is defined as the product of \mathbf{x}_{new} and \mathbf{w} with some

additive Gaussian noise:

$$p(t_{\mathbf{new}} | \mathbf{x}_{\mathbf{new}}, \mathbf{X}, \mathbf{t}, \sigma^2) = \mathcal{N}(\mathbf{x}_{\mathbf{new}}^T \mathbf{w}, \sigma^2)$$

Because this expression and the posterior are both Gaussian, the result of the integral is another Gaussian. In general, if $p(\mathbf{w} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then the expectation is another Gaussian density

$$\mathcal{N}(\mathbf{x}_{\mathbf{new}} \boldsymbol{\mu}_{\mathbf{w}}, \mathbf{x}_{\mathbf{new}}^T \boldsymbol{\Sigma}_{\mathbf{w}} \mathbf{x}_{\mathbf{new}}) \quad (21)$$