

# The Hedging Hyperprior

Stephen D. Simon 1\*

Department of Biomedical and Health Informatics,  
University of Missouri-Kansas City

Yu Jiang

Division of Epidemiology, Biostatistics, and Environmental Health,  
University of Memphis

and

Byron J. Gajewski

Department of Biostatistics,  
University of Kansas Medical Center

June 27, 2017

## Abstract

In many settings, Bayesian models rely on flat or non-informative prior distributions, but sometimes you might prefer to use an informative prior distribution. An informative prior will provide you with added precision, but if your data is inconsistent with **your** prior distribution, you may wish that you had used a flat prior instead. In this paper, we present the hedging hyperprior, a simple modification of the informative prior through the use of a hyperparameter. The hedging hyperprior gives you the best of both worlds. It maintains the added precision that the informative prior gives you when the data is consistent with your prior distribution, but sharply downweights the informative prior distribution when the data is inconsistent. The hedging hyperprior is easy to use for any prior distribution that can be reformulated to include a parameter representing the prior sample size. We illustrate how the hedging hyperprior works for the beta-binomial model and show its equivalence to the modified power prior.

*Keywords:* Beta-binomial model, Informative prior, Modified power prior

---


\*The authors gratefully acknowledge *please remember to list all relevant funding sources in the unblinded version*



# 1 Introduction

Stephen Senn is a delightful source for quotes that are humorous but which slyly illustrate important statistical principles. One of these quotes, found in (Statistical Issues in Drug Development, 2nd edition, Senn SJ (2007) page 46.) is especially relevant to Bayesian data analysts. "A Bayesian is one who, vaguely expecting a horse and catching a glimpse of a donkey, strongly concludes he has seen a mule."

The horse, of course, is the prior distribution, and donkey is the data, and a Bayesian will combine these, even when they are so different to produce something that is neither horse nor donkey. You avoid this problem, of course, if you use a flat or non-informative prior as it puts all (or almost all) of the weight on the data. There are situations, however, where you can and should use an informative prior.

In previous research involving patient accrual models (Gajewski et al 2008  we developed a Bayesian predictive model that asked the researchers for their opinion about how long a proposed clinical trial would take and then asked how confident they were about this estimate on a scale of 1 to 10. Their level of confidence would provide a preliminary estimate of the precision of the prior distribution. An answer of 5, for example, would produce a suggested prior distribution with a precision comparable to 50% of the planned sample size. An answer of 2 would provide a much more variable prior distribution with a precision comparable to 20% of the proposed sample size.

We did not, however, allow the researcher to abdicate total responsibility and suggest a value of 0, corresponding to a flat or non-informative prior. In the setting of patient accrual, a flat prior would be nonsensical. It would be like saying, "I'm not sure, maybe the trial will take 10 days and maybe it will take 10 years. Either value would be quite reasonable." A researcher who could not narrow down the time frame any more than this would be unqualified to conduct the research.

In the context of patient accrual, an informative prior has the additional advantage in that it allows you to predict the duration of the trial before any data is collected and stabilizes the estimated predicted duration of the trial when only a small amount of accrual data is available.

A second  where informative priors make sense is when you have information about

the placebo response rate that you'd expect to see in your current trial based on the placebo response rate that you observed in previous trials. Every drug trial is different because every new drug is different, but often the placebo arm is run using the exact same methods the exact same types of patients every time. You can incorporate the information of these placebo patients using a Bayesian meta-analytic approach. You could use this approach even for drugs that are substantially different because you allow the treatment response rate to have a vague or uninformative prior. A useful example of this approach appears in Walley 2015.

You could argue in this setting that use of a flat prior is unethical. You are failing to utilize valuable information from the previous placebo groups, information that could minimize patient risk by reducing the sample size of the current trial and reducing the number of patients who receive the placebo.

While informative prior distributions can be very useful, they can lead to serious problems if the results of the clinical trial are substantially inconsistent with the prior distribution. We propose a simple modification of the Bayesian model, the hedging hyperprior, that adds of a single hyperparameter to allow downweighting of the informative prior when there is a large discrepancy between the prior distribution and the data.

In section 2 of this article, we show a simple case of the beta-binomial problem with an informative prior, both when the data is reasonably consistent with the prior distribution and when it is markedly different. In section 3, we show how adding a uniform hyperprior can provide the extra precision that you want when the data is consistent with the prior distribution but which simply downweights the strength of the prior distribution when the data is inconsistent with the prior distribution. In section 4, we summarize the modified power prior and show that the hedging hyperprior is equivalent to the modified power prior.

## 2 A beta-binomial model with an informative prior distribution

Consider a Bayesian model for a binomial distribution with a beta prior.

Start with a prior distribution which is  $\text{beta}(4,16)$ . This informative prior, shown in

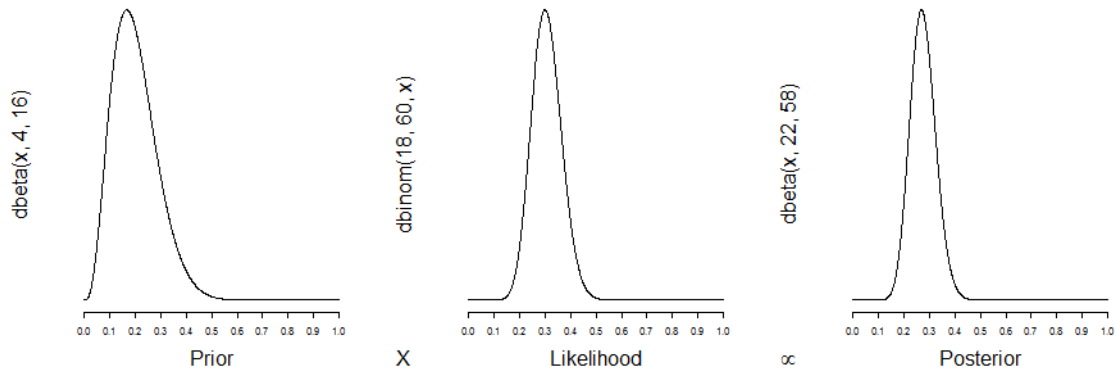


Figure 1: Prior, likelihood, and posterior for a nice beta-binomial model

Figure 1, has a mean of 0.2 and places most of the weight on values for the binomial proportion between 0.1 and 0.4. The precision of the beta distribution is equal to the sum of the two parameters, alpha and beta, which is 20 in this case. This is commonly called the prior sample size. A prior sample size of 20 implies that when you collect 20 observations, the Bayesian model will give equal weight to those 20 observations and the prior distribution.

Suppose that you collect data from a clinical trial and find that among 60 subjects, 18 were classified as "successes" and 42 as "failures" or a proportion of successes equal to 0.3. That's just a bit higher than your prior mean, but still a fairly reassuring result.

The likelihood for this data, also shown in Figure 1, reaches a maximum at 0.3 but the likelihood for values between 0.2 and 0.4 are also reasonable. You combine the prior and the likelihood to get a posterior distribution that is  $\text{beta}(4+18=22, 42+18=58)$ . One of the key benefits of the Bayesian model is the extra precision that the prior distribution provides. The posterior sample size in this example is 80, which is the sum of the prior sample size (20) and the size of the data itself (60).

The posterior distribution,  $\text{beta}(22, 58)$ , has mean of 0.275 and a 95% credible interval of 0.18 to 0.38. The posterior mean is between the mean of the prior distribution and the mean of the data and closer to the mean of the data because the sample size of your data is larger than the sample size of your prior distribution.

This posterior distribution has a sample size of 80 which is equal to the sum of the

Table 1: The classic beta-binomial model.

Beta-binomial formulation

$$\pi \sim \text{beta}(\alpha, \beta)$$
$$x \sim \text{Binomial}(n, \pi)$$

JAGS code for the beta-binomial model (store in jbb.txt).

```
model {  
  pi ~ dbeta(a,b)  
  x ~ dbin(pi,n)  
}
```

R code

```
library("rjags")  
dat.bb <- list(a=4,b=16,x=18,n=60)  
mod.bb <- jags.model("jbb.txt",dat.bb)  
update(mod.bb,1000)  
out.bb <- coda.samples(mod.bb, "pi", 1000)
```

sample size of the prior distribution and the sample size of the data.

You don't really need statistical software to solve this problem, but to set the stage for later Bayesian models, let's see how you would use R and JAGS to solve this simple beta-binomial model. The code shown here is easily adapted to BUGS or STAN. Table 1 shows four lines of JAGS code that you store in a simple text file. Name it jbb.txt and store it in the current working directory of your R session. Then run the R code also shown in Table 1.

You are probably quite happy with yourself at this point. You have a posterior mean which is close to what the data tells you. You don't mind that the mean is a bit smaller because you had some information, possibly from historical data, possibly from interviews with experts in the area, and possibly straight from your gut, that told you that the proportion might be a bit smaller. What really makes you happy, though, is that you have some extra precision because you were smart enough to use an informative prior distribution. Thank you, Reverend Bayes.

Now imagine yourself in a setting that is not quite as nice. A setting that is quite

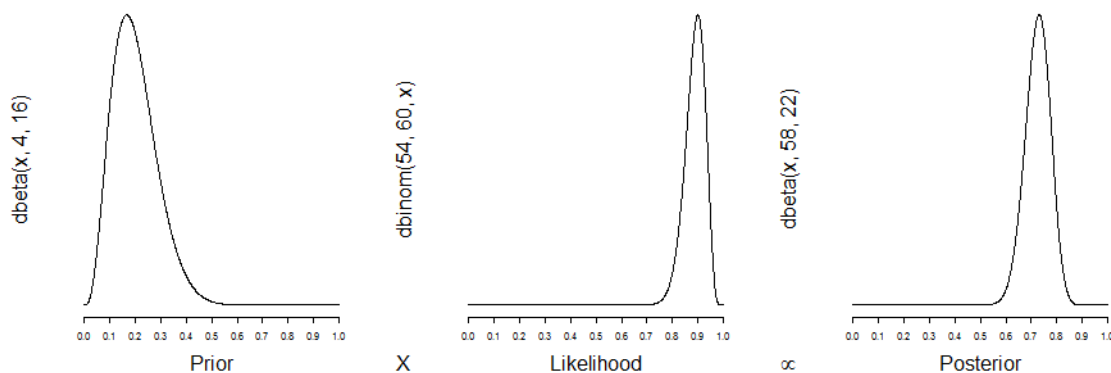


Figure 2: Prior, likelihood, and posterior for a nasty beta-binomial model

bad. More than bad, actually, disastrously bad. Prepare yourself because the next few paragraphs are going to be rather upsetting.

Let's suppose the you have the same prior distribution,  $\text{beta}(4,16)$ , but now your data have 54 successes and only 6 failures. The data is telling you that the binomial proportion is probably around 0.9, but your prior distribution is concentrated around 0.2. When you combine the prior horse with the data donkey, you get a distinctly mulish  $\text{beta}(58, 22)$  with a posterior mean of 0.725 (Figure 2).

The 95% credible interval shows how mulish things have become. This interval (0.62 to 0.82) declares that the prior mean (0.2) is highly improbable. It also declares that the data mean (0.9) is highly improbable, so you're pretty much distrusting every aspect of this data analysis.

There's one more thing, though, that is very troubling about this analysis. Notice that the 95% credible interval in this case has the same width as the interval reported earlier. The equation that you thought you loved, prior sample size of 20 plus a data sample size of 60 gives a posterior sample size of 80, is now haunting you. Your credible interval is too precise because it benefits undeservedly from the precision of a prior distribution that is clearly flawed.

At this point, you want to crawl under a rock. Why oh why, you ask, did I use an informative prior. You want to take back your informative prior, but you can't. If you change your prior distribution after looking at the data, you're pretty sure that the

Reverend Bayes will strike you down with a lightning bolt.

That's how I think you'd react, but maybe you might be bothered at all by this. Perhaps a mule isn't all that bad. After all, you've probably told that bad joke about the statistician who noted the irrefutable demographic fact that the average person has one testicle and one ovary. Statisticians average across disparate groups all the time.

Or you might be thinking, now's not the time to get all wobbly. You didn't just pick that informative prior willy nilly. It was justified through a rigorous review of historical data or the careful elicitation of expert opinion. If you don't have the strength of will to stand by your informative prior, then maybe you aren't cut out to be a true Bayesian.

Whatever your perspective, wouldn't it be nice if you had something between the two extremes of a flat prior and an informative prior? Wouldn't it be nice to have a prior distribution that behaves like an informative prior when your data is consistent with the prior and behaves like a flat prior when your data is inconsistent with the prior?

### 3 The hedging hyperprior

The hedging hyperprior was originally proposed in Jiang et al. (2015) in the context of an informative gamma prior distribution on an exponential waiting time model. This paper will illustrate this same hyperprior in the setting of a beta-binomial model.

The hedging hyperprior adds a new hyperparameter,  $\tau$ , which (at least for now) is uniform on the interval 0 to 1. Use this hyperparameter to modify the two parameters of the beta distribution. When  $\tau$  is 1, the beta distribution remains unchanged from the classic beta-binomial formulation. When  $\tau$  is zero, the beta distribution changes to a flat prior. For values of  $\tau$  in between 0 and 1, the beta distribution becomes a weaker, but still informative prior distribution.

The JAGS code (also shown in 2) is only four lines longer than the code for a beta-binomial model. Note the inclusion of an extra line of code at the end to track the posterior sample size.

The output (Table 3) shows that the posterior mean (0.88) is close to the data mean. The posterior mean of  $\tau$ , 0.028, and the posterior mean sample size, 61.4, indicates that

Table 2: The beta-binomial model with the hdegging hyperprior.  
Hedging hyperprior formulation

$$\begin{aligned}\tau &\sim \text{Uniform}(0, 1) \\ \pi &\sim \text{beta}(1 + \tau(\alpha - 1), 1 + \tau(\beta - 1)) \\ x &\sim \text{Binomial}(n, \pi)\end{aligned}$$

JAGS code for the hedging hyperprior

```
model {
  tau ~ dunif(0,1)
  a0 <- 1+tau*(a-1)
  b0 <- 1+tau*(b-1)
  p ~ dbeta(a0,b0)
  x ~ dbin(p,n)
  post.n <- a0+b0+n
}
```

(R code remains the same)

the prior distribution was sharply downweighted. Note that a harmonic mean might be a better summary of the posterior sample size.

How does the hedging hyperprior perform with data inconsistent with the prior? Quite well actually. Consider the previous example with 54 successes out of 60 trials, but modify the beta(4, 16) prior by adding the hedging hyperparameter. The mean estimate of  $\pi$  is 0.89, which is very close to the proportion observed in the data. The mean estimate for  $\tau$  (0.068) demonstrates that the prior distribution has been sharply downweighted, and the average posterior sample size (61.4) is very close to the sample size of the data, indicating that your precision is not unfairly enhanced by a flawed prior distribution.

Consider the hedging hyperprior for the setting where the data is consistent with the prior (Table 4). The posterior mean for binomial proportion is close to the mean of the data and the posterior mean of  $\tau$  is 0.58, indicating a reduction in the weight placed on the prior distribution, but not as much of a reduction as when the data is inconsistent with the prior distribution. To understand how the hedging hyperprior works, you should look at the joint prior density of  $\pi$  and  $\tau$  (Figure 8).

There's a simple graphical illustration of how the hedging hyperprior works. Adding a



Table 3: Output from the hedging model with nasty data.

```

dat.hedge.nasty <- list(a=4,b=16,x=54,n=60)
mod.hedge.nasty <- jags.model("jhedge.txt",dat.hedge.nasty, quiet=TRUE)
update(mod.hedge.nasty,1000)
out.hedge.nasty <- coda.samples(mod.hedge.nasty, c("pi", "tau", "post.n"), 1000)
summary(out.hedge.nasty)

##
## Iterations = 2001:3000
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 1000
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##           Mean          SD Naive SE Time-series SE
## pi          0.88548 0.04341 0.001373      0.001569
## post.n 61.35855 1.02437 0.032393      0.060343
## tau         0.06793 0.05122 0.001620      0.003017
##
## 2. Quantiles for each variable:
##
##           2.5%       25%       50%       75%      97.5%
## pi          0.787043 0.85879 0.89112 0.91770 0.9535
## post.n 60.130130 60.60162 61.08463 61.86431 63.9260
## tau         0.006506 0.03008 0.05423 0.09322 0.1963

```

Table 4: Output from the hedging model with nice data.

```

dat.hedge.nice <- list(a=4,b=16,x=18,n=60)
mod.hedge.nice <- jags.model("jhedge.txt",dat.hedge.nice, quiet=TRUE)
update(mod.hedge.nice,1000)
out.hedge.nice <- coda.samples(mod.hedge.nice, c("pi", "tau", "post.n"), 1000)
summary(out.hedge.nice)

##
## Iterations = 2001:3000
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 1000
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##           Mean      SD Naive SE Time-series SE
## pi          0.2833 0.05341 0.001689      0.001689
## post.n    71.2837 5.23173 0.165442      0.244287
## tau         0.5642 0.26159 0.008272      0.012214
##
## 2. Quantiles for each variable:
##
##           2.5%      25%      50%      75%      97.5%
## pi          0.1889  0.2443  0.2816  0.3188  0.3945
## post.n    61.7040 67.1342 71.2507 75.8899 79.5801
## tau         0.0852  0.3567  0.5625  0.7945  0.9790

```

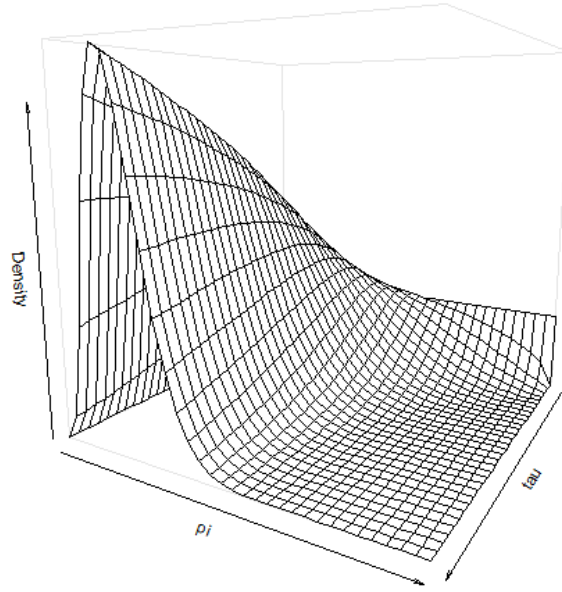


Figure 3: Prior surface for the hedging prior

hyperparameter to the prior distribution produces a prior surface, as shown in Figure 3.

Notice that the front of the surface looks like a beta distribution. In fact, it is the beta distribution that you used as your informative prior. You should think about the shapes of this surface for fixed values of  $\tau$  and for fixed values of  $\pi$  by "slicing" the surface.

For a fixed value of  $\tau$ , the slices are all beta distributions (Figure 4). For  $\tau=1$ , it is the informative prior,  $\text{beta}(4,16)$ . For smaller values of  $\tau$ , the beta distributions get spread out more, corresponding to a weaker prior. For very small values of  $\tau$ , the beta distribution converges on a uniform distribution.

Now let's look at slices for fixed values of  $\pi$  (Figure 5). When  $\pi$  is very large, the distribution of  $\tau$  is skewed strongly towards zero. For moderate values of  $\pi$ , the skew reverses. For very small values of  $\pi$ , it reverts again to skew towards zero.

You can also visualize the likelihood surface. The likelihood is unaffected by the hyperparameter  $\tau$ , and tends to enhance values of  $\pi$  close to the mean of the data. Figure 6 shows the likelihood surface for the case where the data is inconsistent with the prior distribution (54 successes and 6 failures). The likelihood is like a steam roller that flattens out all the small and moderate values of  $\tau$ . What is left is a small peak in the upper left

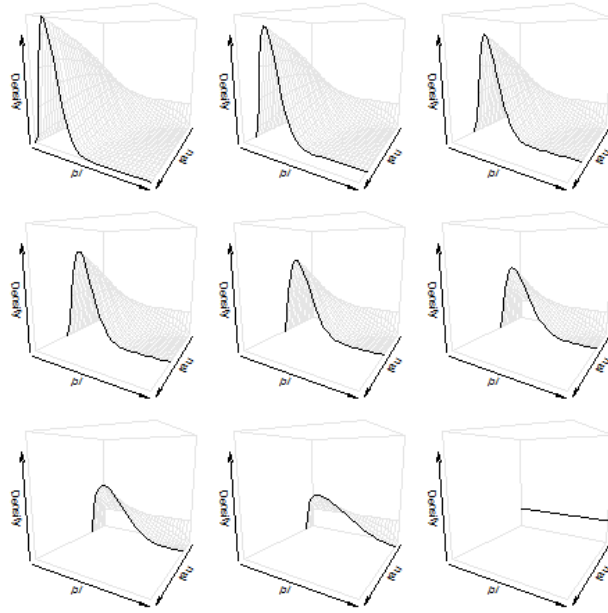


Figure 4: Prior surface with different values of tau highlighted

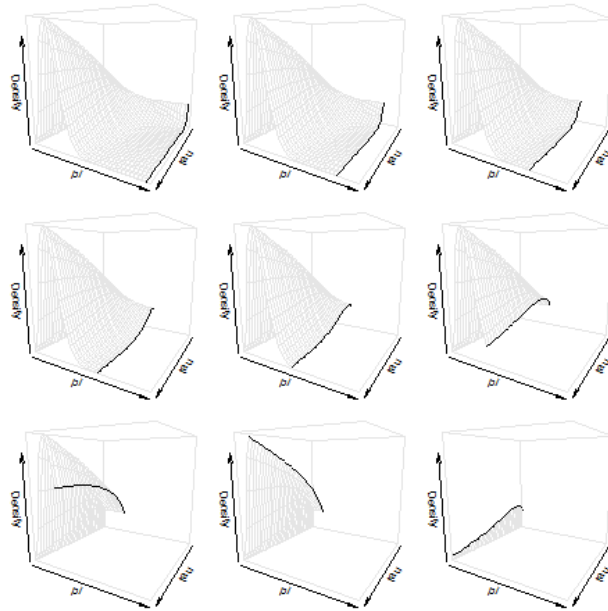


Figure 5: Prior surface with different values of pi highlighted

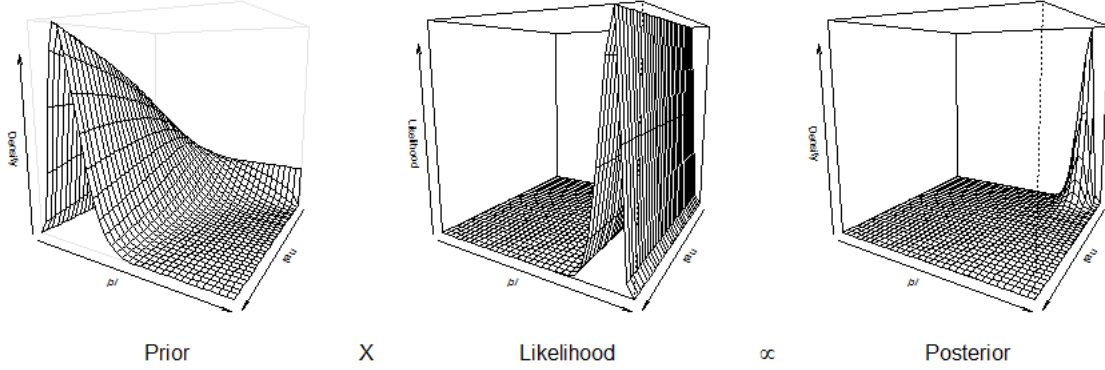


Figure 6: Posterior surface for data with 54 / 60 successes.

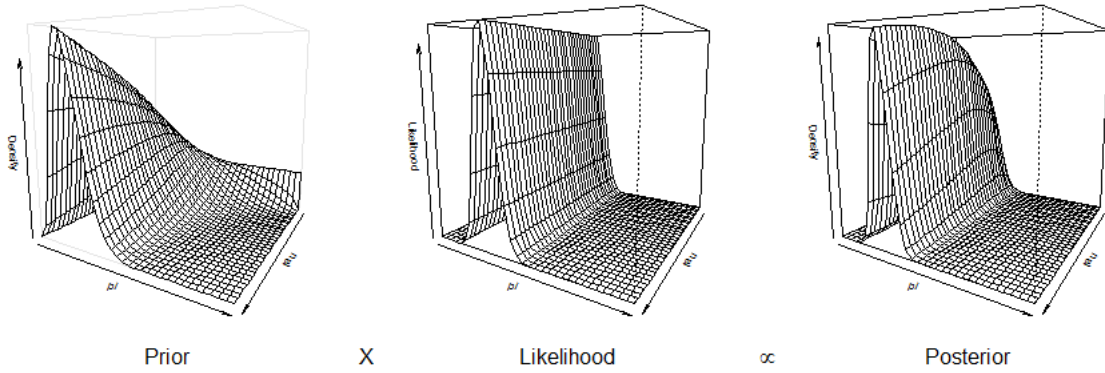


Figure 7: Posterior surface for data with 18 / 60 successes.

corner corresponding to small values of tau and large values of pi.

Figure 7 shows a setting where the data is consistent with the prior distribution (18 successes and 42 failures). The likelihood enhances the existing peak in the prior for moderate values of pi. These are also where the value of tau is large, indicating greater weight on the prior distribution.

The hedging hyperprior allows you to borrow strength from an informative prior, but only to the extent to which the informative prior agrees with the data.

Since the likelihood gets narrower and narrower as the sample size increases, the slices of the prior surface for fixed values of pi become very important. They represent the limiting case for the hedging hyperprior. In particular, the expected value of tau for a fixed value

of  $\pi$  represents how much the historic prior will end up getting discounted in the limiting case.

You can plot this mean versus  $\pi$  to get the "hedge trace," a simple two dimensional plot that shows you where serious discounting of the informative prior starts (Figure 8).

Notice that placing using a uniform distribution for the hyperprior leads to some amount of downweighting, even when the data and the prior agree perfectly. This is not too surprising, as the uniform distribution has a mean of 0.5. You can correct for this by starting with a uniform distribution on the interval 0 to 2 instead of 0 to 1. This works out well on average, and even gives you a small bonus to precision if your data is better than expected. If you want to avoid the small bonus, use a uniform distribution on the interval 0 to 1.6. This gives you an average value of  $\tau$  that maxes out fairly close to but never larger than 1 across a broad range of priors.

You are also not restricted to just a uniform distribution for the hedging hyperprior. A bit too extreme, perhaps, is the suggestion in Hobbs et al. (2011) to use a  $\text{beta}(10,1)$  if you want to strongly encourage reliance on the informative prior and a  $\text{beta}(1,10)$  if you want to strongly discourage reliance on the informative prior. But skewing slightly to the left or right from a uniform distribution does allow you to fine tune the ranges where you want to strongly downweight the prior distribution.

You can also substitute a Bernoulli distribution for the uniform hyperprior. This effectively provides a weighted average between a posterior based on the full informative prior and a posterior based on a flat prior, with the relative weights controlled by the degree to which the data disagrees with the informative prior.

## 4 Equivalence to the modified power prior

The modified power prior has been proposed as a way to downweight a bad historical prior distribution.

Consider the historic prior and the data as two sets of data,  $D_0$  and  $D$ . If, for example, your informative prior were  $\text{beta}(4, 16)$ , then you would create a pseudo data set with  $(4-1)$  successes and  $(16-1)$  failures. Apply this pseudo data to a flat prior ( $\pi_0(\theta)$ ), get a posterior

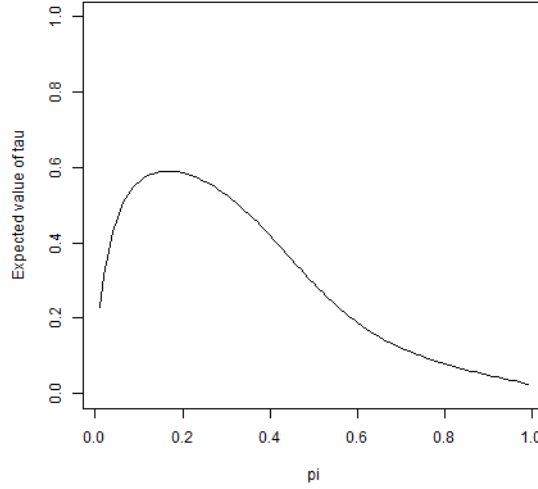


Figure 8: Hedge trace.

distribution the resembles your informative prior.

$$\pi(\theta|D_0) \propto L(\theta|D_0)\pi_0(\theta)$$

Then use  $\pi(\theta|D_0)$  use that as your prior distribution for the new data.

$$\pi(\theta|D, D_0) \propto L(\theta|D)\pi(\theta|D_0)$$

What Ibrahim et al. (2003) suggested is that you could downweight the historical prior by raising the historical likelihood to the power.

$$\pi(\theta|D_0, \tau) \propto L(\theta|D_0)^\tau \pi_0(\theta)$$

With  $\tau=1$ , you get the full strength of the historical prior and with  $\tau=0$  you wipe out the historic data and are left with just a flat prior. This formulation does not satisfy the likelihood principle, but you can include a fudge factor  $(\int L(\theta|D_0)^\tau \pi_0(\theta) d\theta)$  to make things work nicely.

It is fairly easy to show that hedging hyperprior in the beta-binomial model is equivalent to the modified power prior. Notice that

$$L(\theta|D_0)^\tau = \left( \frac{(\alpha+\beta-2)!}{(\alpha-1)!(\beta-1)!} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right)^\tau$$

is proportional to

$$\theta^{(\alpha-1)\tau} (1-\theta)^{(\beta-1)\tau}$$

If you multiply this by a uniform prior for theta and apply the fudge factor, you get a

$\text{beta}(1 + (\alpha - 1)\tau, 1 + (\beta - 1)\tau)$  distribution for  $L(\theta|D_0)$ .

The advantage to formulating the problem as the hedging hyperprior is that it makes it easy to see how to use standard Bayesian software like BUGS, JAGS, or STAN. Also, the hedging hyperprior is easier to visualize, allowing you to readily discern its behavior.

## 5 Conclusion

The hedging hyperprior is a simple modification that you can use whenever you have an informative prior distribution and the distribution can be reparameterized to include a measure of the prior sample size. It produces a posterior estimate that uses some of the precision of the prior distribution when the data is reasonably consistent with the data. When the data is inconsistent with the data, however, the hedging hyperprior sharply downweights the prior distribution.

The hedge trace, a graph of how the mean of the hyperparameter changes as the main parameter varies, helps you understand how much disagreement between the prior is needed before you see significant downweighting of the prior distribution.

The  $\text{uniform}(0,1)$  hyperprior is perhaps not the best choice for a hyperprior. It still significantly downweights the prior distribution even when the data is consistent with the prior. Two ways to improve the performance of the hedging hyperprior is to use a wider interval, like a  $\text{uniform}(0,2)$  or select a beta distribution which puts greater weight on the larger values of the hedging hyperprior, like a  $\text{beta}(2,1)$ . Another approach worth considering is replacing the continuous uniform hyperprior with a discrete 0-1 Bernoulli distribution.

The beta-binomial model is an ideal setting for understanding how the hedging hyperprior works because the parameter space is only two dimensional and is bounded in all directions. Further research is needed to explore the role of the hedging hyperparameter in more complex models and to better understand the performance of various alternative hyperprior distributions.



## 6 Bibliography

### References

- Hobbs, B. P., B. P. Carlin, S. J. Mandrekar, and D. J. Sargent (2011, September). Hierarchical Commensurate and Power Prior Models for Adaptive Incorporation of Historical Information in Clinical Trials. *Biometrics* 67(3), 1047–1056.
- Ibrahim, J. G., M.-H. Chen, and D. Sinha (2003). On optimality properties of the power prior. *Journal of the American Statistical Association* 98(461), 204–213.
- Jiang, Y., S. Simon, M. S. Mayo, and B. J. Gajewski (2015, February). Modeling and Validating Bayesian Accrual Models on Clinical Data and Simulations Using Adaptive Priors. *Statistics in medicine* 34(4), 613–629.