

Métodos Bayesianos

Estimando la proporción de obtener cara al lanzar una moneda

Resumen: Bajo la suposición que la moneda es justa, aplicaremos diferentes métodos bayesianos para comprobar esta suposición. Con una distribución prior $\text{Beta}(1,1)$ y con el método AR, vemos que nuestra moneda no es justa. Además, el Bayes Factor de la hipótesis alternativa es de 17940, por lo que hay una evidencia decisiva a favor de esta hipótesis.

Clara Albert
Enero del 2022

Índice

Introducción	2
Porque usar análisis bayesiano	3
Distribución binomial	3
Simulación con R	4
Distribución a priori.....	5
Simulación de la posterior con el método AR.....	6
Comparación de las dos distribuciones.....	8
Test de hipótesis	8
Conclusiones	11

Introducción

Este trabajo consiste en realizar un estudio bayesiano de unos datos de libre elección adecuados para realizar dicho análisis.

Estuve buscando un tiempo datos para utilizarlos para realizar este trabajo, pero no había ningún conjunto de datos que me pareciese interesante, por lo que decidí simular yo los datos.

Escogí una de las simulaciones más fáciles de realizar: tirar al aire una moneda e ir apuntando los resultados: cara o cruz. Es bien sabido que este experimento sigue una distribución binomial.

En nuestro caso, el parámetro lo sabremos para simular nuestra binomial y así obtener los parámetros de la binomial. A partir de la muestra obtenida, realizaremos todo nuestro análisis bayesiano suponiendo que no conocemos p y lo estimaremos a partir de algunos métodos empleados a lo largo del curso.

Nos ayudaremos del software estadístico R para realizar la simulación de tirar una moneda y el posterior análisis bayesiano.

Calcularemos una distribución posterior empleando dos métodos para poder compararlos. Además, realizaremos una prueba de hipótesis para testar si realmente nuestro parámetro de interés es 0.5.

Por lo tanto, el objetivo principal del trabajo es estimar a partir de diferentes métodos la proporción del número de caras cuando tiramos una moneda.

Porque usar análisis bayesiano

Hay muchas razones por el cual usar este tipo de análisis. Es más flexible que el enfoque frecuentista en:

1. Puedes incluir información adicional a los datos
2. Puedes realizar cualquier comparación entre grupos o conjuntos de datos
3. Nos deja cuantificar la incertidumbre de los datos
4. Los modelos están definidos de forma más simple y son más flexibles

Distribución binomial

Nuestro experimento consiste en tirar una moneda al aire e ir apuntando el número de caras y de cruces que sale. El resultado será un vector compuesto por 0 (si es cara) y 1 (si es cruz). Por lo tanto:

X = “número de caras obtenidas”, donde x_i es el resultado de la i -ésima tirada.

La variable se clasifica como éxito, si el evento ocurre, o fracaso, si el evento no ocurre. Por lo tanto,

$$X \sim \text{Binomial}(n, p)$$

Como X sigue una distribución binomial, ya que cuenta el número de éxitos en n tiradas de Bernoulli.

Definimos la distribución binomial:

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

Donde p está restringido al espacio $= [0,1]$ y se define como la probabilidad de éxito.

La función de verosimilitud de esta función es:

$$L(X|p) \propto p^x (1 - p)^{n-x}$$

Una vez ya sabemos de que manera se distribuye nuestro experimento, vamos a simularlo haciendo uso del software estadístico R.

Simulación con R

Hay varias opciones de simular un suceso binomial. Entre estas funciones están: `sample()` o `rbinom()`.

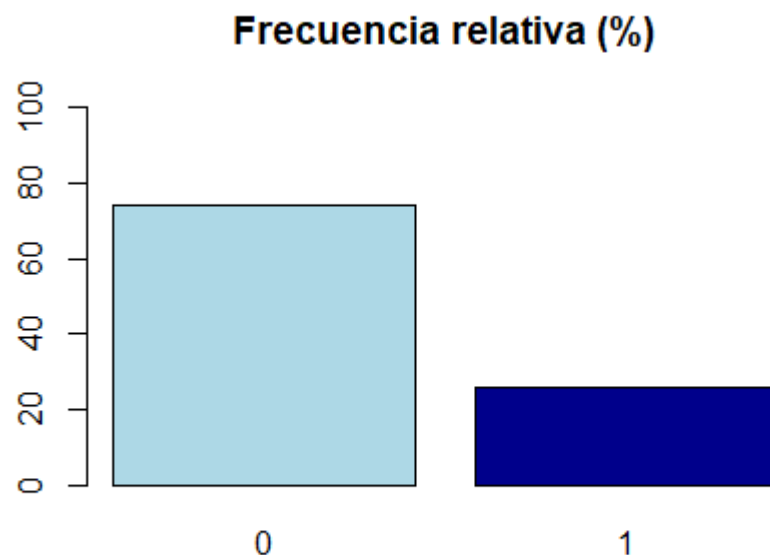
Utilizaremos la función `rbinom(100,c(0,1), prob=0.5)` ya que las tiradas de una moneda siguen esta distribución y nos dará mejores resultados.

Realizaremos 100 tiradas, suficientes para tener una muestra representativa. Además, la probabilidad de obtener cara o cruz es del 0.5, ya que partimos de la base que nuestra moneda es justa.

Es importante introducir una semilla en R antes de simular nuestra muestra, para que nuestros resultados sean reproducibles. Utilizando la semilla 836 los resultados obtenidos son:

0	1
74	26

Vemos que a simple vista los resultados no son los esperados, ya que la respuesta tendría que estar más balanceada. Representando la frecuencia relativa en % de las dos posibilidades se puede apreciar aún más la descompensación entre las dos categorías.



A partir de esta simulación, nuestra proporción de cara es:

$$\hat{p} = \frac{x}{n} = \frac{74}{100} = 0.74$$

Sin realizar ninguna prueba estadística podemos observar que se aleja bastante de la proporción introducida al realizar la simulación. Aun así, vamos a estimar ésta utilizando diferentes métodos.

Distribución a priori

Aún no hemos buscado ninguna información por internet sobre la proporción de obtener cara o cruz. Por lo tanto, nuestra distribución a priori será lo más no informativa posible.

Podemos pensar que sigue una distribución uniforme restringida en el intervalo [0,1]. Puesto que la distribución uniforme es un caso particular de la distribución beta, nuestra distribución a priori será una distribución beta.

$$\pi(p) \sim \text{Beta}(a, b) \propto x^{a-1} (1-x)^{b-1}$$

En este caso, la distribución a posterior es:

$$f(p|X) \propto L(X|p)\pi(p)$$

Si tomamos como prior una Beta(1,1) (lo menos informativa posible) tenemos que:

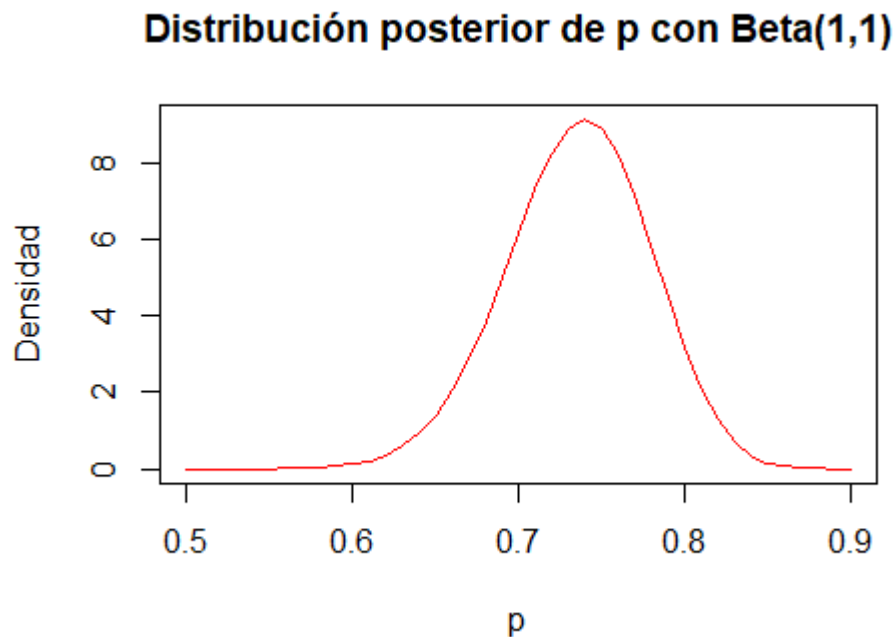
$$f(p|X) \propto p^x (1-p)^{n-x} p^{a-1} (1-p)^{b-1}$$

Tomando x=74; n=100; a=1 y b=1:

$$f(p|X) \propto p^{74} (1-p)^{100-74} p^{1-1} (1-p)^{1-1} \propto p^{74} (1-p)^{26} \sim \text{Beta}(75, 27)$$

Vemos que, en esta situación en específico, la distribución posterior es una expresión explícita y fácil de manipular, por lo que encontrar estadísticos de interés no nos supondrá ninguna complicación.

La densidad de esta prior gráficamente es:



A partir de su distribución podemos estimar los intervalos de credibilidad y la mediana.

$$\hat{p} \in (0.6459, 0.8158)$$

Y su mediana es que en este caso correspondería a su estimación puntual.

$$\hat{p} = 0.7368$$

Simulación de la posterior con el método AR

El método de aceptación/rechazo consiste en generar valores de una variable aleatoria X con densidad $f(x)$. Suponemos que sabemos generar fácilmente una densidad diferente $g(x)$ y que existe una constante C tal que:

$$\frac{f(x)}{cg(x)} \leq 1$$

En inferencia bayesiana, $g(x) = \pi(\theta)$.

- 1) Generamos u como una uniforme $\rightarrow u \sim Unif(0,1)$
- 2) Generamos una variable y que se distribuya como la prior $\rightarrow y \sim \pi(\theta)$
- 3) Si $u \leq \frac{L(X|Y)}{L(X|\hat{\theta})}$ devolvemos y; si no, volvemos al paso 1

Para poder aplicar este método a la simulación de una distribución posterior tenemos que calcular el MLE de \hat{p} . Calculamos la derivada del logaritmo de la función de verosimilitud y la igualamos a 0.

$$l = \log(L(X|p)) \propto x \log(p) + (n - x) \log(1 - p)$$

$$\frac{\partial l}{\partial p} = \frac{x}{p} - \frac{n - x}{1 - p} = 0 \rightarrow \hat{p} = \frac{x}{n}$$

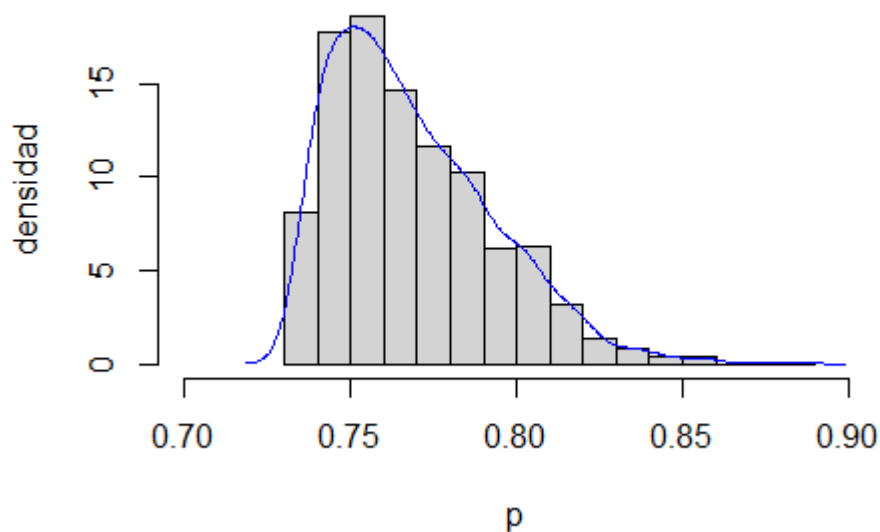
De nuestra simulación con R hemos obtenido

$$\hat{p} = \frac{x}{n} = 0.74$$

A partir del estimador, podemos encontrar el valor de MLE. Lo buscamos a partir de una función que he creado. A partir de este valor, podemos programar el método AR.

A partir del vector que guarda todos los valores que se han aceptado podemos graficar su densidad como también ver un resumen de los estadísticos básicos.

Distribución posterior de p con AR



Mínimo	0.7351
Mediana	0.7638
Media	0.7638
Máximo	0.8824
Varianza	0.0242

$$\hat{p} \in (0.7367, 0.8216)$$

Comparación de las dos distribuciones

Hemos simulado dos distribuciones: una a partir de una prior y otra con el método AR utilizando esa misma prior.

Vemos como los intervalos de confianza difieren en su precisión, haciendo el método AR mucho más preciso en este aspecto. Esto es debido a que con este método evaluamos que el valor simulado de nuestra prior esté dentro de nuestra posterior y si no es el caso descartamos este valor.

Con este ejemplo tan simple podemos afirmar que el algoritmo AR es una herramienta que nos permite realizar simulaciones de la posterior más precisas que solo haciendo uso de la prior y la función de verosimilitud.

En este caso concreto, puede ser debido a que la eficiencia de nuestra simulación con AR es 0.5252, un valor muy elevado cuando se trata de simular distribuciones posteriores.

Test de hipótesis

Vamos a realizar una comparativa entre el método frecuentista y el bayesiano realizando la misma prueba de hipótesis.

H_0 : La proporción de cara es 0.5

H_1 : La proporción de cara es diferente a 0.5

Escogeremos un nivel de significación del 0.05 para todas las pruebas.

Enfoque frecuentista

Podemos utilizar dos pruebas de hipótesis que nos ofrece R: `binom.test` y `prop.test`.

`prop.test()` se utiliza cuando el tamaño de la muestra es mayor que 30 ya que la distribución muestral de la proporción muestral se puede aproximar por una normal. Aun así, podemos usar la otra función ya que conocemos la distribución exacta de nuestros datos.

Utilizaremos las dos para comparar los resultados:

binom.test(74,100,p=0.5)

$$\hat{p} \in (0.6468, 0.8226)$$

Con este test, rechazamos de que la proporción de cara es 0.5.

prop.test(74,100,p=0.5)

$$\hat{p} \in (0.6409, 0.8202)$$

También rechazamos la hipótesis nula, por lo que la proporción de cara no es 0.5.

Enfoque bayesiano

Para realizar este test de hipótesis bayesiano utilizaremos una prior impropia:

$$\pi(p) \sim 1$$

Y seremos lo más no informativos posibles:

$$P(H_0) = P(H_1) = 0.5$$

Con esta información podemos calcular $P(H_0|X)$:

$$P(H_0|X) = \frac{P(H_0) \int_{p_0} L(X|p) \pi(p) dp}{P(H_0) \int_{p_0} L(X|p) \pi(p) dp + P(H_1) \int_{p_1} L(X|p) \pi(p) dp}$$

$$P(H_0|X) = \frac{L(X|p = 0.5)}{L(X|p = 0.5) + \int_{p_1} L(X|p) \pi(p) dp} = \frac{7.88 \cdot 10^{-31}}{7.88 \cdot 10^{-31} + 1.41 \cdot 10^{-26}}$$

$$= 5.57 \cdot 10^{-5}$$

$$L(X|p = 0.5) = 0.5^{74} (1 - 0.5)^{26} = 7.88 \cdot 10^{-31}$$

$$\int_{p_1} \text{Beta}(75,27) = \int_0^1 p^{75-1} (1-p)^{27-1} = \frac{\Gamma(75)\Gamma(27)}{\Gamma(75+27)} = 1.41 \cdot 10^{-26}$$

Y la $P(H_1|X)$ es:

$$P(H_1|X) = 1 - P(H_0|X) = 0.9999$$

El factor de Bayes de H_0 respecto H_1 es:

$$\frac{P(H_0|X)}{1 - P(H_0|X)} = \frac{5.57 \cdot 10^{-5}}{0.9999} = 5.573 \cdot 10^{-3}$$

El factor de Bayes de H_1 respecto H_0 es:

$$\frac{P(H_1|X)}{1 - P(H_1|X)} = \frac{0.999}{5.57 \cdot 10^{-5}} = 17940.89$$

Según la escala de Jeffreys la evidencia a favor de H_1 es decisiva y según la escala de Kass y Raftery la evidencia es muy fuerte.

Por lo tanto, tenemos suficiente evidencia para concluir que nuestra proporción sea igual a 0.5.

Conclusiones

A partir de todos los métodos bayesianos utilizados a lo largo del trabajo podemos afirmar que nuestra simulación realizada con R no sigue la binomial esperada.

Cuando simulé los datos, supuse que la moneda era justa y con estadística descriptiva ya se podía intuir que en verdad no lo era.

Cuando tomamos una distribución a priori Beta(1,1) se podía ver claramente que la distribución en 0.5 era prácticamente 0, haciendo muy poco probable nuestra suposición de moneda justa. Además, el intervalo de credibilidad nos respaldaba esta conclusión ya que con un 95% de confianza el intervalo no contenía el parámetro deseado.

Cuando simulamos la distribución con el método de aceptación/rechazo, el intervalo de credibilidad era mucho menos ancho, por lo que era mejor distribución, ya que la varianza era mucho menor. Además, no obtuvimos ni una vez el parámetro deseado, ya que el mínimo fue 0.7351.

Aún ya habiendo comprobado que la proporción de nuestra moneda no era 0.5, la prueba de hipótesis respaldó aún más esa conclusión.

Realizamos la prueba con un enfoque frecuentista y obtuvimos unos intervalos de confianza mucho más anchos, pero que seguían sin comprender el valor esperado. Por lo tanto, con un 95% de confianza pude confirmar que el parámetro no había caído en el intervalo.

Tomando una prior impropia y siendo lo menos informativos posibles, el bayes factor de H_1 respecto H_0 nos concluyó que la evidencia en contra de H_0 era decisiva.

Con todos estos resultados, la conclusión es que nuestra moneda no es justa.

Bibliografía

<https://www.rstudio.com/wp-content/uploads/2015/02/rmarkdown-cheatsheet.pdf>

<https://www.uv.es/mlejarza/actuarios/iibayes.pdf>

<https://www.rdocumentation.org/packages/BayesianFirstAid/versions/0.1/topics/bayes.prop.test>

<https://www.sumsar.net/blog/2014/06/bayesian-first-aid-prop-test/>

<https://www.r-bloggers.com/2019/04/understanding-bayesian-inference-with-a-simple-example-in-r/>

<https://methods.sagepub.com/dataset/howtoguide/bayesian-influence-ncmp-2017-r>

<https://a-little-book-of-r-for-bayesian-statistics.readthedocs.io/en/latest/src/bayesianstats.html>

<https://www.statology.org/bayes-factor/>

https://es.wikipedia.org/wiki/Distribuci%C3%B3n_binomial

<https://sites.warnercnr.colostate.edu/gwhite/wp-content/uploads/sites/73/2017/04/BinomialLikelihood.pdf>

Anexo

```
set.seed(836)
moneda=rbinom(100, c(0,1),prob=0.5)
moneda=as.factor(moneda)

summary(moneda)
barplot(prop.table(table(moneda)) * 100, col=c("lightblue","darkblue"),
, main="Frecuencia relativa (%)", ylim=c(0,100))

# Estimación de La propotción
x=sum(moneda == 0); x
n=length(moneda); n
n-x

# gráfico de La Likelihood
calcLikelihoodForProportion <- function(successes, total)
{
  curve(dbinom(successes,total,x)); abline(v=successes/total, col =
"blue") # plot the Likelihood
}
calcLikelihoodForProportion(x,n)

# Podemos observar que el pico de La función de verosimilitud está apr
oximadamente en 0.74 que equivale a La media de La muestra (74/100 = 0.
74). En otras palabras, La proporción más probable, dado esta muestra,
es 0.74.

# Prior Beta(1,1)
p <- seq(0.5, 0.9, by = 0.01)
plot(p, dbeta(p, 75, 27), type =
"l",col="red",main="Distribución posterior de p con Beta(1,1)", ylab="
Densidad")
qbeta(0.025, 75,27); qbeta(0.5, 75,27); qbeta(0.975, 75,27)

# Método AR
log_likelihoood = function(p){
  ll=sum(dbinom(x, size=round(p*100), prob=p, log=T))
  return(ll)
}

llmax=log_likelihoood(x/n)

nsim=1000 # Number of simulations
co=0;tot=0
prob_ar=numeric(nsim)

set.seed(836)
while(co<nsim){
  lu=log(runif(1))
  prob=rbeta(1,75,27)
  if(lu<=log_likelihoood(prob)-llmax){
    co=co+1
    prob_ar[co]=prob}
}
```

```

    tot=tot+1
  }
  efficiency <- nsim/tot; efficiency

  hist(prob_ar, probability = TRUE, main = "Distribución posterior de p
  con AR", ylab="densidad", xlab="p", xlim=c(0.7,0.9))
  lines(density(prob_ar), col='blue')

  summary(prob_ar)
  quantile(prob_ar,c(0.025,0.975))
  sd(prob_ar)

# Test de hipotesis

# Enfoque frecuentista
  binom.test(74,100,p=0.5)
  prop.test(74,100,p=0.5)

# Enfoque bayesiano
  h = 0.5^(74) * 0.5^(26)
  g = (gamma(75)*gamma(27))/gamma(75+27)
  H0=h/(h+g)
  H1=g/(g+h)
  1-h/(h+g)
  H0/H1
  H1/H0

```