

Freie Universität Berlin  
Chair of Statistics  
Location: Berlin  
Summer Term 2017  
Lecture: Einführung in die Bayes-Statistik  
Examiner: Dr. Florian Meinfelder

## **Program leave-one-out posterior predictive checking in R**

Johannes Brinkmann (), jojo-brinkmann@gmx.de  
Carlo Michaelis (5043128), carlo.michaelis@gmail.com  
Max Reinhardt (579174), max\_reinhardt@me.com  
Adrian Rolf (), adrian.rolf@gmx.de

Master Statistics  
August 21, 2017

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Code</b>	<b>4</b>
2.1	General code . . . . .	4
2.2	Model evaluaion . . . . .	5
2.3	Cook's Distance . . . . .	6
<b>3</b>	<b>References</b>	<b>10</b>

**List of Figures**

1	Model evaluation using MSE and lppd . . . . .	6
2	Traces . . . . .	8
3	Densities . . . . .	9
4	Cook’s distance measure . . . . .	10
5	Pointwise cook’s distance after three different simulation sizes, compared to analytic cook’s distance from frequentist linear regression model using least squares. . . . .	11

**List of Tables**

# 1 Introduction

## 2 Code

### 2.1 General code

The Gibbs sampler

```
gibbsSampler <- function(X, Y, R, b = 0, initSigma = 1) {  
  B <- R + b  
  
  # Size of design matrix  
  n <- nrow(X)  
  p <- ncol(X)  
  
  # Variables to store the samples in (initialize sigma with initSigma)  
  betas <- matrix(nrow = B, ncol = p)  
  sigma <- c(initSigma, rep(NA, B-1))  
  
  # Sampling  
  for(i in 1:B){  
    # OLS of beta  
    V <- solve(t(X)%*%X)      #  $(X^T X)^{-1}$   
    beta_hat <- V%*%t(X)%*%Y #  $(X^T X)^{-1} X^T Y$   
  
    # OLS of sigma  
    sigma_hat <- t(Y-X%*%beta_hat)%*%(Y-X%*%beta_hat)/(n-p)  
  
    # Sample beta from the full conditional  
    betas[i,] <- rmvnorm(1, beta_hat, sigma[i]*V)  
  
    # Sample sigma from the full conditional  
    if(i < B) {  
      sigma[i+1] <- 1/rgamma(1, (n-p)/2, (n-p)*sigma_hat/2)  
    }  
  }  
  
  # Remove burn in  
  if(b != 0) {  
    betas <- betas[-(1:b),]  
    sigma <- sigma[-(1:b)]  
  }  
  
  return(list(betas = betas, sigma = sigma))  
}
```

Cross Validation function

```
crossValidation <- function(X, Y, R, b) {  
  # Size of design matrix  
  n <- nrow(X)  
  p <- ncol(X)  
  steps <- length(R)  
  
  # Run gibbs sampler to get sampled parameters  
  samples <- lapply(1:n, function(i) gibbsSampler(X[-i,], Y[-i], R[steps], b))  
  
  # Initialize lists
```

```

Sigma <- list()
Betas <- list()
Yhati <- list()

# Calculate sigma, betas and Yhati for every step
for(k in 1:steps) {
  Sigma[[k]] <- sapply(samples, function(sample) mean(sample$sigma[1:R[k]]))
  Betas[[k]] <- sapply(samples, function(sample) colMeans(sample$betas[1:R[k],]))
  Yhati[[k]] <- sapply(1:n, function(i) X[i,]%*%Betas[[k]][,i])
}

return(list(Sigma = Sigma, Betas = Betas, Yhati = Yhati))
}

```

## 2.2 Model evaluation

```

bayesModelEvaluation <- function(models, Y, R, b) {
  # Evaluate multiple models and return results from all models
  n <- length(Y)
  k <- length(models)

  # Cross validate every model
  results <- lapply(1:k, function(i) crossValidation(models[[i]], Y, R, b))

  # Calculate Mean Squared Errors
  MSEs <- sapply(results, function(e1) {
    return((1/(n-nrow(e1$Betas[[1]])))*sum((Y-e1$Yhati[[1]])^2))
  })

  # Plot MSEs
  par(mfrow = c(1,2))
  barplot(MSEs, xlab = "Models", ylab = "MSE", names.arg = seq(1,k),
    main = "Model evaluation using MSE")

  # Calculate log posterior predictive density (log likelihood)
  #  $y \sim N(XB, s^2(X^T X)^{-1}) = N(\hat{Y}, s^2(X^T X)^{-1})$ 
  LPPDs <- sapply(results, function(eva) sum(log(dnorm(Y, eva$Yhati[[1]], eva$Sigma[[1]]))))

  # Plot LPPDs
  barplot(-2*LPPDs, xlab = "Models", ylab = "lppd", names.arg = seq(1,k),
    main = "Model evaluation using lppd")

  return(list(MSEs = MSEs, LPPDs = LPPDs))
}

```

Run model evaluation with data

```

# Swiss data
dat <- swiss

# Response variable
Y <- dat$Fertility
n <- nrow(dat)

# Design matrices
models <- list(

```

```

matrix(c(rep(1,n), dat$Education), nrow=n),
matrix(c(rep(1,n), dat$Agriculture), nrow=n),
matrix(c(rep(1,n), dat$Examination), nrow=n),
matrix(c(rep(1,n), dat$Catholic), nrow=n),
matrix(c(rep(1,n), dat$Education, dat$Agriculture), nrow=n),
matrix(c(rep(1,n), dat$Education, dat$Examination), nrow=n),
matrix(c(rep(1,n), dat$Education, dat$Catholic), nrow=n),
matrix(c(rep(1,n), dat$Education, dat$Agriculture, dat$Examination), nrow=n),
matrix(c(rep(1,n), dat$Education, dat$Agriculture, dat$Catholic), nrow=n),
matrix(c(rep(1,n), dat$Education, dat$Examination, dat$Catholic), nrow=n),
matrix(c(rep(1,n), dat$Education, dat$Agriculture, dat$Examination, dat$Catholic), nrow=n)
)

### Model evaluation ###
criteria <- bayesModelEvaluation(models, Y, R = 50, b = 10) # R = 500, b = 100

```

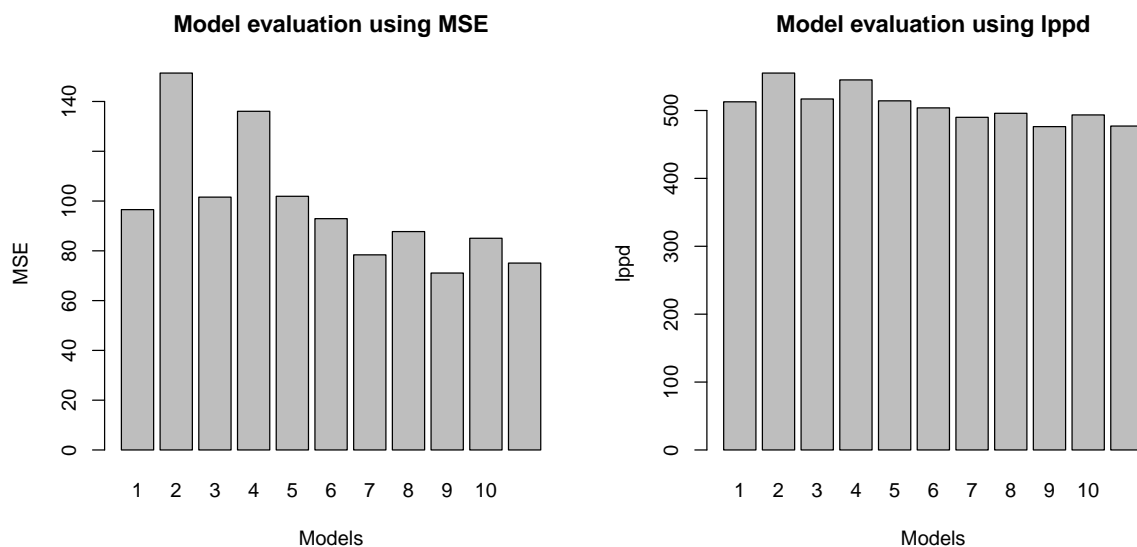


Figure 1: Model evaluation using MSE and lppd

```

# Check proportion
print(criteria$MSEs/(-2*criteria$LPPDs))

## [1] 0.1882963 0.2727617 0.1964855 0.2496199 0.1981878 0.1844312 0.1599963
## [8] 0.1769362 0.1492633 0.1723552 0.1573758

# TODO: Als Tabelle ausgeben? Explain why not fulfilled

# Choose optimal modal
optIdx <- which.min(-2*criteria$LPPDs)

# Check if MSE and LPPD would choose the same
print(paste(which.min(criteria$MSEs), "=", which.min(-2*criteria$LPPDs)))

## [1] "9 = 9"

# TODO den Vergleich einfach nur im Text verwenden?

```

## 2.3 Cook's Distance

```

cooksDistance <- function(X, Y, R, b) {
  # R is a vector

```

```

B <- R + b

# Prepare projection matrix and number of parameters
H <- X%*%solve(t(X)%*%X)%*%t(X)
p <- ncol(X)
steps <- length(R)

# Run cross validation with vector R
cv <- crossValidation(X, Y, R, b)

# Sample whole model (add R+b samples, remove b later)
sample <- gibbsSampler(X, Y, R[steps] + b)

# Cook's distance: Frequentist approach with analytic solution
betaHat <- solve(t(X)%*%X)%*%t(X)%*%Y
YhatLinReg <- X%*%betaHat
E <- Y-YhatLinReg
cooksLinReg <- (E^2/((1/(n-p))*sum(E^2)*p))*(diag(H)/(1-diag(H))^2)

cooksBayesCV <- matrix(nrow = n, ncol = steps)
for(k in 1:steps) {
  # Estimate posterior mean from betas
  betas <- colMeans(sample$betas[b:B[k],])

  # Predict values, using posterior mean
  Yhat <- X%*%betas

  # Calculate cook's distance
  dists <- apply(cv$Betas[[k]], 2, function(betas) {
    return((Yhat - X%*%betas)^2)
  })
  mse <- (1/(n-p))*sum((Y-Yhat)^2)
  cooksBayesCV[,k] <- colSums(dists)/(p*mse)
}

return(list(cooksBayesCV = cooksBayesCV, cooksLinReg = cooksLinReg, sample = sample))
}

### Cook's Distance ###

# Number of samples
b <- 100
R <- seq(100,500,50) # seq(100,5000,50)

# Get optimal model and calculate projection matrix
cooks <- cooksDistance(models[[optIdx]], Y, R, b)

```

Traces from bayes regression sampling

Posterior densities from bayes regression sampling

Apply cook's distance measure

```

# Initialize matrices
steps <- length(R)
matBayesD <- cooks$cooksBayesCV
matFreqD <- matrix(rep(cooks$cooksLinReg, steps), ncol=steps)

# Calculate differences between Cook's distance methods
cooksMeasure <- colSums((matBayesD-matFreqD)^2)

```

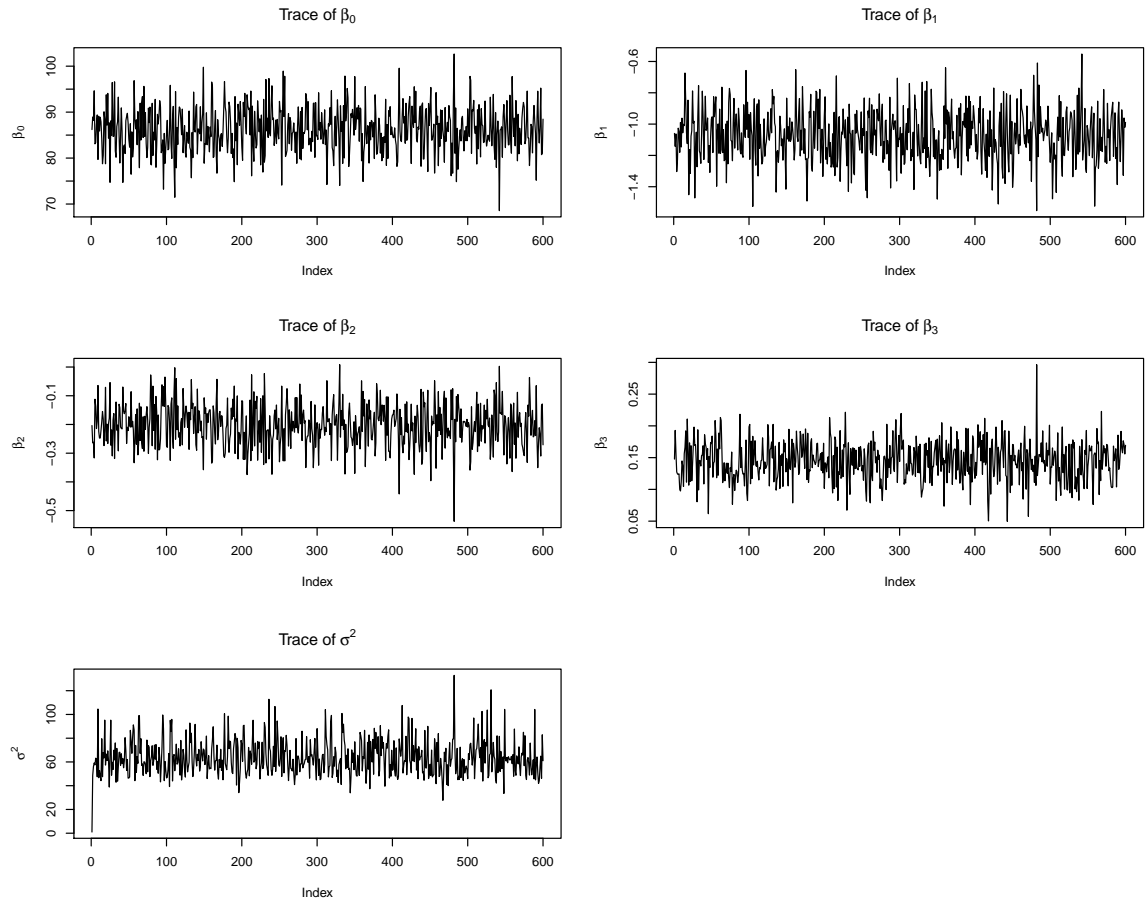


Figure 2: Traces



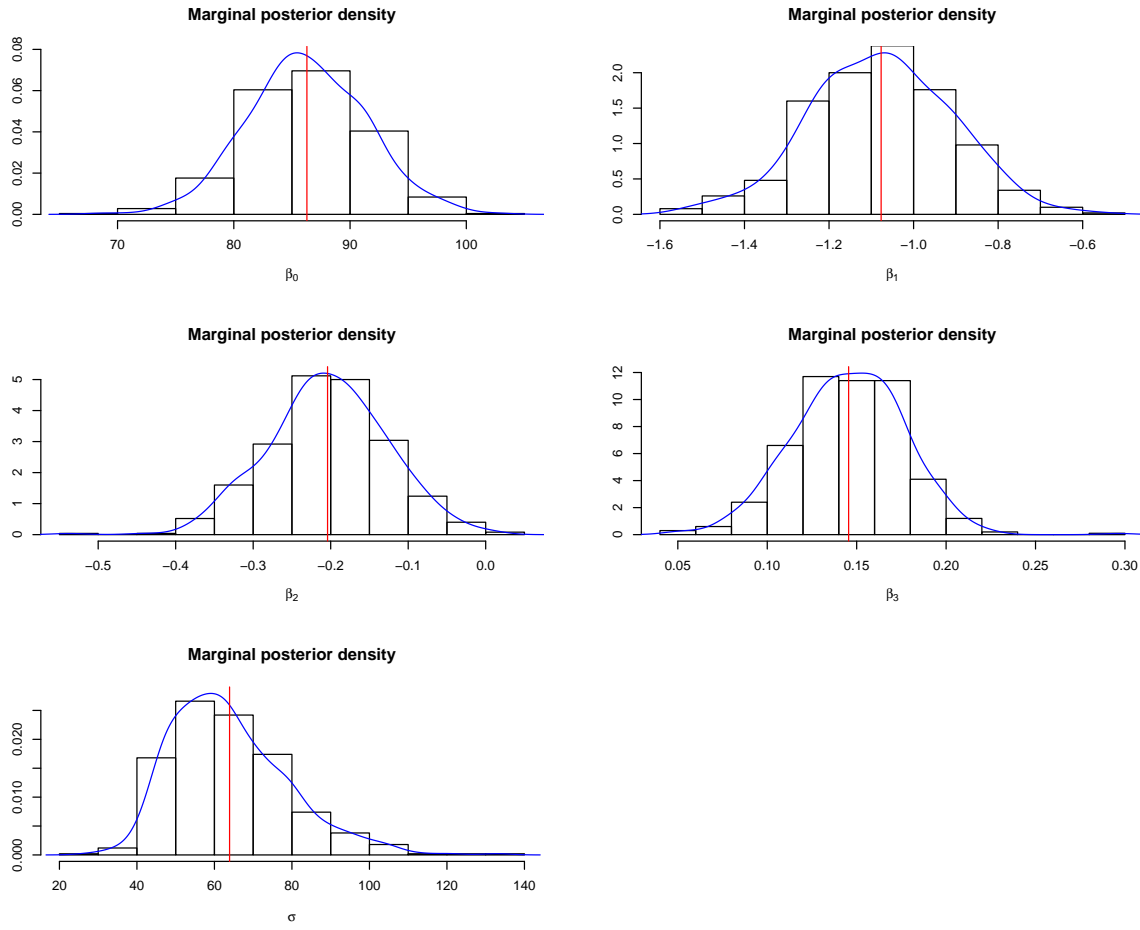


Figure 3: Densities

Cook's distance measure

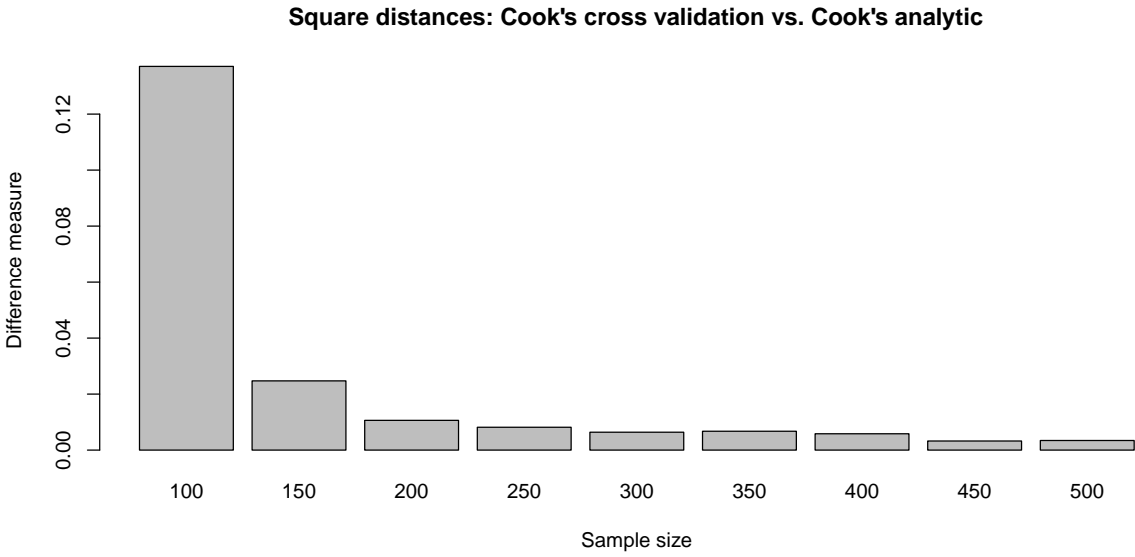


Figure 4: Cook's distance measure

Cook's distance, comarison between bayes and frequentist

### 3 References

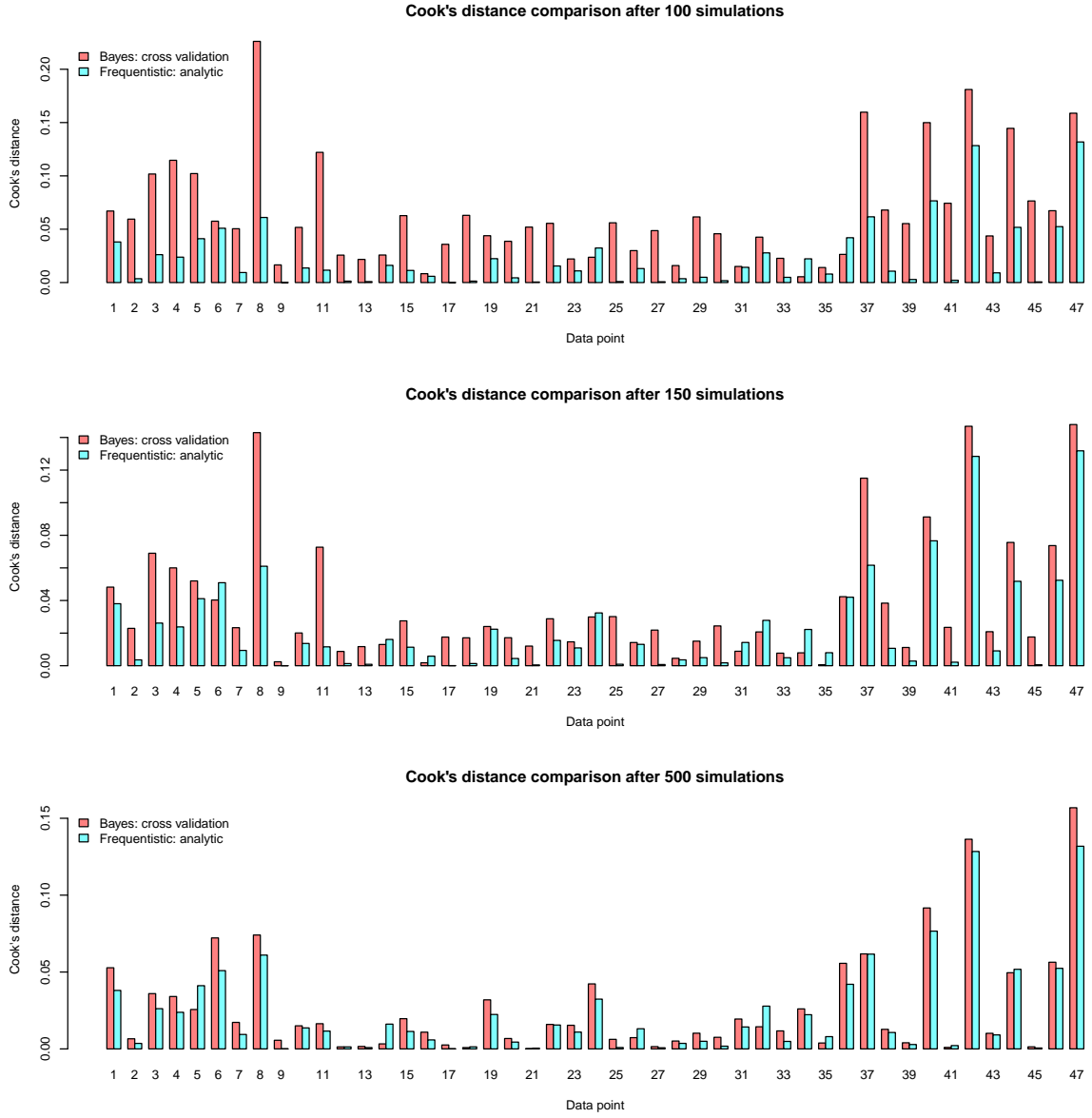


Figure 5: Pointwise cook's distance after three different simulation sizes, compared to analytic cook's distance from frequentist linear regression model using least squares.