

Bayesian model averaging

Dr. Jarad Niemi

STAT 544 - Iowa State University

March 9, 2017

Outline

- Bayesian model averaging
- BIC model averaging
- Model search
- Parameter averaging
- Posterior inclusion probability
- Model selection

Bayesian Model Averaging

The posterior predictive distribution

$$p(\tilde{y}|y) = \int p(\tilde{y}|\theta)p(\theta|y)d\theta$$

assumes there is a true model $p(y|\theta)$ and accounts for the uncertainty in θ .

If you want to account for model uncertainty amongst some set of models M_1, \dots, M_h , you can use the Bayesian model averaged posterior predictive distribution

$$p(\tilde{y}|y) = \sum_{h=1}^H p(\tilde{y}|M_h, y)p(M_h|y)$$

where

- $p(M_h|y)$ is the posterior model probability and
- $p(\tilde{y}|M_h, y)$ is the predictive distribution under model M_h .

Normal example

Suppose we have two models:

$$\begin{aligned}
 Y_i | M_0 &\stackrel{\text{ind}}{\sim} N(0, 1) \\
 Y_i | M_1, \mu &\stackrel{\text{ind}}{\sim} N(\mu, 1), \mu | M_1 \sim N(0, 1)
 \end{aligned}$$

for $i = 1, \dots, n$.

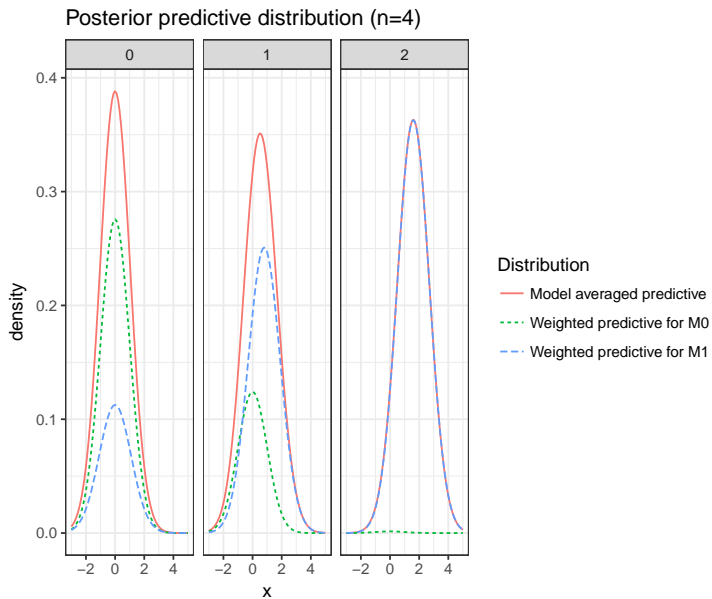
Thus, we have the following posterior predictive distributions

$$\begin{aligned}
 \tilde{y} | y, M_0 &\sim N(0, 1) \\
 \tilde{y} | y, M_1 &\sim N(n\bar{y}[n+1]^{-1}, [n+1]^{-1} + 1)
 \end{aligned}$$

for scalar \tilde{y} independent of y , but from the same distribution.
and the following posterior model probabilities:

$$\begin{aligned}
 p(M_0 | y) &\propto N(y; 0, \mathbf{I}) \\
 p(M_1 | y) &\propto N(y; 0, \mathbf{I} + \mathbf{1}\mathbf{1}^\top)
 \end{aligned}$$

where $\mathbf{1}$ is a vector of all 1s.



AIC/BIC model averaging

The generic structure for model averaging is

$$p(\tilde{y}|y) = \sum_{h=1}^H p(\tilde{y}|M_h, y)w_h$$

where w_h is the **weight** for model h .

Here are some possible weights:

- Bayesian model averaging: $w_h = p(M_h|y)$
- AIC model averaging: $w_h = e^{-\Delta_h/2}$ where $\Delta_h = AIC_h - \min AIC$
- AICc model averaging: $w_h = e^{-\Delta_h/2}$ where $\Delta_h = AICc_h - \min AICc$
- BIC model averaging: $w_h = e^{-\Delta_h/2}$ where $\Delta_h = BIC_h - \min BIC$

Information criterion

Recall that information criteria have the form:

$$IC = -2 \log L(\hat{\theta}) + P$$

where P is a penalty. So if you take

$$w_h = e^{-\Delta_h/2} = e^{-(IC_h - \min IC)/2} \propto e^{-IC_h/2} = L_h(\hat{\theta})e^P.$$

where, if p is the number of parameters, the penalty P is

- AIC: $2p$
- AICc: $2p + 2p(p + 1)/(n - p - 1)$
- BIC: $p \log(n)$

The BIC is a large sample approximation to the marginal likelihood:

$$-2 \log p(y) \approx -2 \log p(y|\theta) + p \log(n) + C$$

Regression BMA

A common place to perform Bayesian Model Averaging is in the regression framework:

$$y \sim N(X_\gamma \beta_\gamma, \sigma_\gamma^2 \mathbf{I})$$

where γ is a vector indicator of which of the p explanatory variables are included in model γ , e.g.

$$\gamma = (1, 1, 0, \dots, 0, 1, 0)$$

indicates the first, second, \dots , and penultimate explanatory variables are included.

BIC model averaging in R

```
library(BMA)

## Loading required package: survival
## Loading required package: leaps
## Loading required package: robustbase
##
## Attaching package: 'robustbase'
## The following object is masked from 'package:survival':
##
## heart
## Loading required package: inline
## Loading required package: rrcov
## Scalable Robust Estimators with High Breakdown Point (version 1.4-3)

library(MASS)

##
## Attaching package: 'MASS'
## The following object is masked from 'package:dplyr':
##
## select

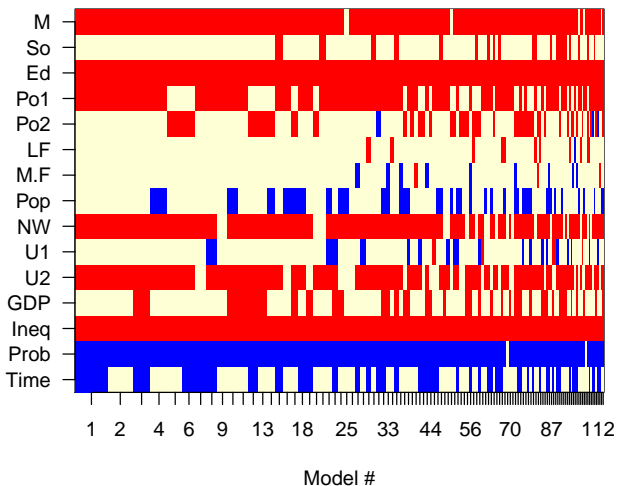
data(UScrime)
x<- UScrime[,-16]
y<- log(UScrime[,16])
x[,-2]<- log(x[,-2])
lma<- bicreg(x, y, strict = FALSE, OR = 20)
```

```
summary(lma)
```

```
##
## Call:
## bicreg(x = x, y = y, strict = FALSE, OR = 20)
##
##
## 115 models were selected
## Best 5 models (cumulative posterior probability = 0.2039 ):
##
##          p!=0   EV      SD   model 1   model 2   model 3   model 4   model 5
## Intercept 100.0 -23.45301 5.58897 -22.63715 -24.38362 -25.94554 -22.80644 -24.50477
## M          97.3  1.38103 0.53531  1.47803  1.51437  1.60455  1.26830  1.46061
## So         11.7  0.01398 0.05640  .         .         .         .         .
## Ed         100.0  2.12101 0.52527  2.22117  2.38935  1.99973  2.17788  2.39875
## Po1        72.2  0.64849 0.46544  0.85244  0.91047  0.73577  0.98597  .
## Po2        32.0  0.24735 0.43829  .         .         .         .         0.90689
## LF          6.0  0.01834 0.16242  .         .         .         .         .
## M.F         7.0 -0.06285 0.46566  .         .         .         .         .
## Pop         30.1 -0.01862 0.03626  .         .         .         -0.05685  .
## NW          88.0  0.08894 0.05089  0.10888  0.08456  0.11191  0.09745  0.08534
## U1          15.1 -0.03282 0.14586  .         .         .         .         .
## U2          80.7  0.26761 0.19882  0.28874  0.32169  0.27422  0.28054  0.32977
## GDP         31.9  0.18726 0.34986  .         .         0.54105  .         .
## Ineq       100.0  1.38180 0.33460  1.23775  1.23088  1.41942  1.32157  1.29370
## Prob        99.2 -0.24962 0.09999 -0.31040 -0.19062 -0.29989 -0.21636 -0.20614
## Time       43.7 -0.12463 0.17627 -0.28659  .         -0.29682  .         .
##
## nVar              8              7              9              8              7
## r2                0.842          0.826          0.851          0.838          0.823
## BIC               -55.91243      -55.36499      -54.69225      -54.60434      -54.40788
## post prob         0.062          0.047          0.034          0.032          0.029
```

```
imageplot.bma(lma)
```

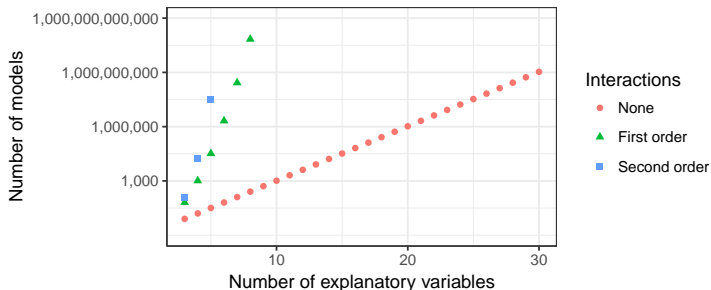
Models selected by BMA



Model space

For all subsets regression analysis with p (continuous or binary) explanatory variables, we have

- 2^p models with no interactions,
- $2^{\binom{p}{2}}$ times as many models when considering first order interactions,
- $2^{\binom{p}{3}}$ times as many models when considering second order interactions,
- etc.



Model search in R

When model enumeration isn't possible, we resort to model search. There are many ways to search the model space, but one common approach is to use Markov chain Monte Carlo.

```
library(BMS)
data(datafls)

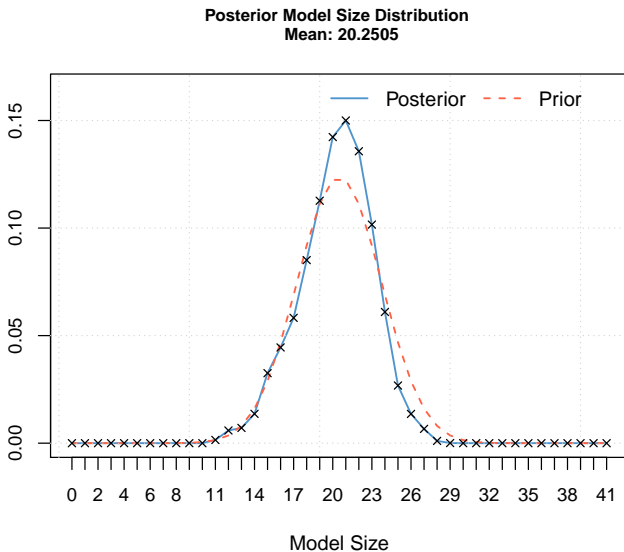
bma1 = bms(datafls,
  burn = 10000,
  iter = 20000,
  mprior = "uniform", # uniform prior over models
  user.int = FALSE)
```

If there is a uniform prior over models, what is the prior over model size (the number of explanatory variables included)?

```
summary(bma1)
```

## Mean no. regressors	Draws	Burnins	Time	No. models visited
## "20.2505"	"20000"	"10000"	"2.577701 secs"	"9083"
## Modelspace 2^K	% visited	% Topmodels	Corr PMP	No. Obs.
## "2.2e+12"	"4.1e-07"	"15"	"0.0943"	"72"
## Model Prior	g-Prior	Shrinkage-Stats		
## "uniform / 20.5"	"UIP"	"Av=0.9863"		

```
plotModelsize(bma1)
```

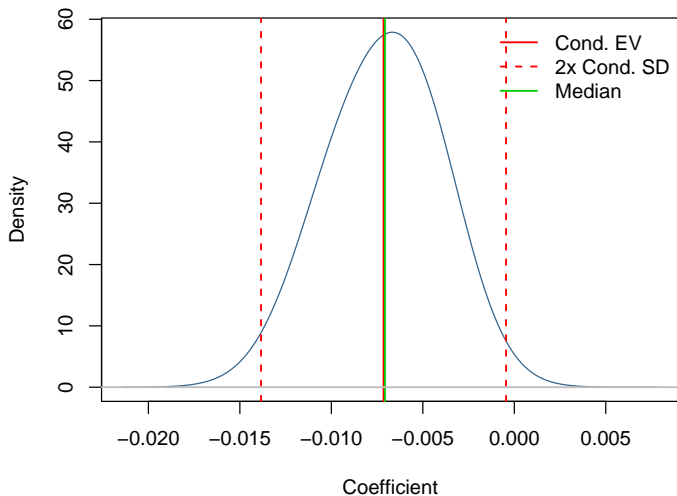


coef(bmal)

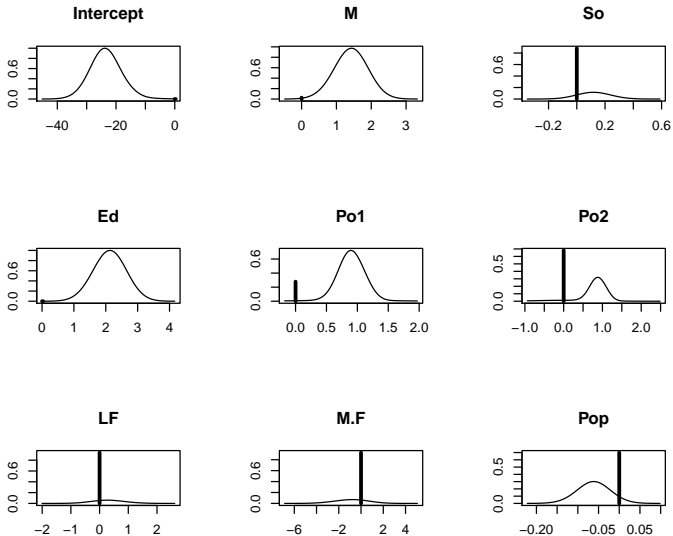
##	PIP	Post Mean	Post SD	Cond.Pos.Sign	Idx
## GDP60	1.00000	-1.684067e-02	2.767613e-03	0.00000000	12
## Confucian	1.00000	6.516425e-02	1.252213e-02	1.00000000	19
## LifeExp	0.98075	8.632707e-04	2.523443e-04	1.00000000	11
## SubSahara	0.97425	-1.948019e-02	6.303480e-03	0.00000000	7
## Hindu	0.93730	-7.671338e-02	3.554069e-02	0.00144031	21
## EquipInv	0.93100	1.298783e-01	5.497151e-02	1.00000000	38
## LabForce	0.89625	2.577815e-07	1.329617e-07	0.99709902	29
## RuleofLaw	0.87540	1.062581e-02	5.900113e-03	0.99942883	26
## Mining	0.85940	3.241153e-02	1.801450e-02	1.00000000	13
## HighEnroll	0.78600	-7.920375e-02	5.366558e-02	0.00165394	30
## EthnoL	0.78200	1.003031e-02	6.794904e-03	1.00000000	20
## NequipInv	0.69865	3.547674e-02	2.915181e-02	1.00000000	39
## LatAmerica	0.66015	-7.980742e-03	7.103132e-03	0.00098462	6
## EcoOrg	0.61870	1.205759e-03	1.191094e-03	1.00000000	14
## PrScEnroll	0.58875	1.114659e-02	1.162043e-02	0.99193206	10
## BlMktPm	0.57735	-4.387049e-03	4.509557e-03	0.00000000	41
## Spanish	0.54375	6.399302e-03	7.446361e-03	0.97765517	2
## CivlLib	0.54195	-1.179624e-03	1.460181e-03	0.02869268	34
## Protestants	0.52115	-5.070570e-03	6.144129e-03	0.00000000	25
## French	0.50410	4.733616e-03	5.722305e-03	0.99523904	3
## Muslim	0.48810	5.596900e-03	7.027168e-03	0.99539029	23
## Brit	0.46605	2.832771e-03	4.076957e-03	0.93648750	4
## English	0.40310	-3.105822e-03	4.612595e-03	0.00000000	35
## OutwarOr	0.38475	-1.271946e-03	1.982752e-03	0.00025991	8
## Buddha	0.35815	3.215589e-03	5.450239e-03	0.99804551	17
## PolRights	0.34995	-4.742097e-04	1.056684e-03	0.13730533	33
## PubEduPct	0.30710	4.961903e-02	1.017019e-01	0.96450668	31
## WarDummy	0.26605	-7.875159e-04	1.748354e-03	0.00225522	5
## Age	0.22510	-8.092115e-06	1.947451e-05	0.00244336	16
## RFEXDist	0.21665	-6.101185e-06	1.774585e-05	0.03623356	37
## Catholic	0.19810	-6.082755e-04	2.964384e-03	0.27662797	18
## UnkD	0.16990	2.149584e-04	0.704520e-03	0.20000000	00


```
density(bma1, reg="BIMktPm")
```

Marginal Density: BIMktPm (PIP 50.58 %)



```
plot(lma)
```



Model averaged parameters

Consider the following set of 4 models with $Y_i \stackrel{ind}{\sim} N(\mu_i, \sigma^2)$ where

$$M_1 : \mu_i = \beta_0$$

$$M_2 : \mu_i = \beta_0 + \beta_1 X_{i,1}$$

$$M_3 : \mu_i = \beta_0 + \beta_2 X_{i,2}$$

$$M_4 : \mu_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2}$$

It is tempting to want to obtain a model averaged posterior for the coefficients.

Model averaged parameters (cont.)

Perhaps we can write a model averaged posterior for a parameter as

$$p(\beta_1|y) = \sum_{h=1}^H p(\beta_1|y, M_h)p(M_h|y)$$

But β_1 means something entirely different in these models:

- In model M_2 , β_1 is the effect of a one unit increase in $X_{i,1}$ on the expected response.
- In model M_4 , β_1 is the effect of a one unit increase in $X_{i,1}$ on the expected response **after adjusting for $X_{i,2}$** .

More accurate model

Consider the following set of 4 models with $Y_i \stackrel{ind}{\sim} N(\mu_i, \sigma^2)$ where

$$M_1 : \mu_i = \alpha_0$$

$$M_2 : \mu_i = \beta_0 + \beta_1 X_{i,1}$$

$$M_3 : \mu_i = \gamma_0 + \gamma_2 X_{i,2}$$

$$M_4 : \mu_i = \delta_0 + \delta_1 X_{i,1} + \delta_2 X_{i,2}$$

Now it seems clear that we cannot average these parameters.

Assessing explanatory variable importance

To obtain some measure of how important a particular explanatory variable is we can find its **posterior inclusion probability**, i.e. the probability it is non-zero:

$$p(\beta_j \neq 0|y) = \sum_{h:\beta_j \neq 0} p(M_h|y)$$

which is just the sum of the model probabilities for the models where β_j is not zero.

```
summary(lma)

##
## Call:
## bicreg(x = x, y = y, strict = FALSE, OR = 20)
##
##
## 115 models were selected
## Best 5 models (cumulative posterior probability = 0.2039 ):
##
##           p!=0   EV      SD   model 1   model 2   model 3   model 4   model 5
## Intercept 100.0 -23.45301 5.58897 -22.63715 -24.38362 -25.94554 -22.80644 -24.50477
## M          97.3  1.38103 0.53531  1.47803  1.51437  1.60455  1.26830  1.46061
## So         11.7  0.01398 0.05640  .         .         .         .         .
## Ed         100.0  2.12101 0.52527  2.22117  2.38935  1.99973  2.17788  2.39875
## Po1        72.2  0.64849 0.46544  0.85244  0.91047  0.73577  0.98597  .
## Po2        32.0  0.24735 0.43829  .         .         .         .         0.90689
## LF          6.0  0.01834 0.16242  .         .         .         .         .
## M.F         7.0 -0.06285 0.46566  .         .         .         .         .
## Pop         30.1 -0.01862 0.03626  .         .         .         -0.05685  .
## NW          88.0  0.08894 0.05089  0.10888  0.08456  0.11191  0.09745  0.08534
## U1          15.1 -0.03282 0.14586  .         .         .         .         .
## U2          80.7  0.26761 0.19882  0.28874  0.32169  0.27422  0.28054  0.32977
## GDP         31.9  0.18726 0.34986  .         .         0.54105  .         .
## Ineq       100.0  1.38180 0.33460  1.23775  1.23088  1.41942  1.32157  1.29370
## Prob        99.2 -0.24962 0.09999 -0.31040 -0.19062 -0.29989 -0.21636 -0.20614
## Time       43.7 -0.12463 0.17627 -0.28659  .         -0.29682  .         .
##
## nVar                8          7          9          8          7
## r2                   0.842     0.826     0.851     0.838     0.823
## BIC                  -55.91243  -55.36499  -54.69225  -54.60434  -54.40788
## post prob            0.062     0.047     0.034     0.032     0.029
```

coef(bmal)

##	PIP	Post Mean	Post SD	Cond.Pos.Sign	Idx
## GDP60	1.00000	-1.684067e-02	2.767613e-03	0.00000000	12
## Confucian	1.00000	6.516425e-02	1.252213e-02	1.00000000	19
## LifeExp	0.98075	8.632707e-04	2.523443e-04	1.00000000	11
## SubSahara	0.97425	-1.948019e-02	6.303480e-03	0.00000000	7
## Hindu	0.93730	-7.671338e-02	3.554069e-02	0.00144031	21
## EquipInv	0.93100	1.298783e-01	5.497151e-02	1.00000000	38
## LabForce	0.89625	2.577815e-07	1.329617e-07	0.99709902	29
## RuleofLaw	0.87540	1.062581e-02	5.900113e-03	0.99942883	26
## Mining	0.85940	3.241153e-02	1.801450e-02	1.00000000	13
## HighEnroll	0.78600	-7.920375e-02	5.366558e-02	0.00165394	30
## EthnoL	0.78200	1.003031e-02	6.794904e-03	1.00000000	20
## NequipInv	0.69865	3.547674e-02	2.915181e-02	1.00000000	39
## LatAmerica	0.66015	-7.980742e-03	7.103132e-03	0.00098462	6
## EcoOrg	0.61870	1.205759e-03	1.191094e-03	1.00000000	14
## PrScEnroll	0.58875	1.114659e-02	1.162043e-02	0.99193206	10
## BlMktPm	0.57735	-4.387049e-03	4.509557e-03	0.00000000	41
## Spanish	0.54375	6.399302e-03	7.446361e-03	0.97765517	2
## CivlLib	0.54195	-1.179624e-03	1.460181e-03	0.02869268	34
## Protestants	0.52115	-5.070570e-03	6.144129e-03	0.00000000	25
## French	0.50410	4.733616e-03	5.722305e-03	0.99523904	3
## Muslim	0.48810	5.596900e-03	7.027168e-03	0.99539029	23
## Brit	0.46605	2.832771e-03	4.076957e-03	0.93648750	4
## English	0.40310	-3.105822e-03	4.612595e-03	0.00000000	35
## OutwarOr	0.38475	-1.271946e-03	1.982752e-03	0.00025991	8
## Buddha	0.35815	3.215589e-03	5.450239e-03	0.99804551	17
## PolRights	0.34995	-4.742097e-04	1.056684e-03	0.13730533	33
## PubEduPct	0.30710	4.961903e-02	1.017019e-01	0.96450668	31
## WarDummy	0.26605	-7.875159e-04	1.748354e-03	0.00225522	5
## Age	0.22510	-8.092115e-06	1.947451e-05	0.00244336	16
## RFEXDist	0.21665	-6.101185e-06	1.774585e-05	0.03623356	37
## Catholic	0.19810	-6.082755e-04	2.964384e-03	0.27662797	18
## UrbanPop	0.16990	2.149584e-04	2.704520e-03	0.20000000	00

Multiple posterior inclusion probability

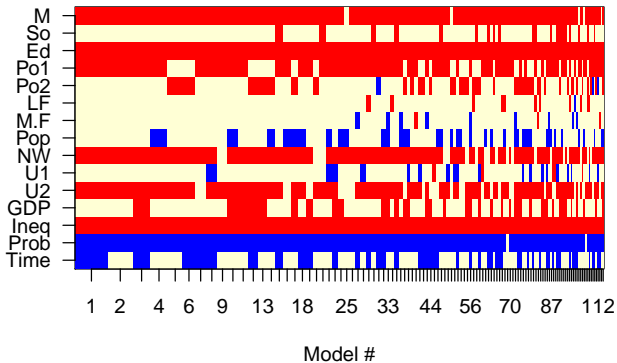
If explanatory variables are correlated, then it is possible to have low posterior inclusion probability for the correlated explanatory variable, but the probability of at least one of the explanatory variables being included is high.

For example,

$$P(\beta_i \neq 0 \text{ or } \beta_j \neq 0 | y) = \sum_{h: \beta_i \neq 0 \text{ or } \beta_j \neq 0} p(M_h | y)$$

```
imageplot.bma(lma)
```

Models selected by BMA



```
cor(UScrime$Po1, UScrime$Po2)
```

```
## [1] 0.9935865
```

Model selection

Sometimes, we will want to select a model. Selecting model M_h is clearly justified if $p(M_h|y) \approx 1$.

If forced to choose a model, it might seem that choosing the model with the highest $p(M_h|y)$ would be the way to go, but Barbieri and Berger (2004) show that if prediction is the goal, then the **median probability model** is better. The **median probability model** is the model that includes all explanatory variables whose posterior inclusion probability is greater than $1/2$.