

# Model Averaging for Linear Regression

Erkki P. Liski

University of Tampere  
Department of Mathematics, Statistics and Philosophy

## Outline

- ▶ The Model
- ▶ Model selection
- ▶ Model average estimator (MAE)

## Outline

- ▶ The Model
- ▶ Model selection
- ▶ Model average estimator (MAE)
- ▶ Why MAE?
- ▶ General structure of MAE

## Outline

- ▶ The Model
- ▶ Model selection
- ▶ Model average estimator (MAE)
- ▶ Why MAE?
- ▶ General structure of MAE
- ▶ Selecting the model weights
- ▶ Finite Sample Performance

## Homoscedastic linear regression

**Variables** The response  $y$  and the predictors  $x_1, x_2, \dots$

## Homoscedastic linear regression

**Variables** The response  $y$  and the predictors  $x_1, x_2, \dots$

### The Model

$$y = \mu + \varepsilon, \quad \mu = \sum_{j=1}^{\infty} \beta_j x_j,$$
$$E(\varepsilon|\mathbf{x}) = 0, \quad E(\varepsilon^2|\mathbf{x}) = \sigma^2,$$

$\beta_1, \beta_2, \dots$  and  $\sigma^2$  are unknown parameters, and  $\mathbf{x} = (x_1, x_2, \dots)$ .

## Homoscedastic linear regression

**Variables** The response  $y$  and the predictors  $x_1, x_2, \dots$

### The Model

$$y = \mu + \varepsilon, \quad \mu = \sum_{j=1}^{\infty} \beta_j x_j,$$
$$E(\varepsilon|\mathbf{x}) = 0, \quad E(\varepsilon^2|\mathbf{x}) = \sigma^2,$$

$\beta_1, \beta_2, \dots$  and  $\sigma^2$  are unknown parameters, and  $\mathbf{x} = (x_1, x_2, \dots)$ .

### Further

$E(\mu^2) < \infty$  and  $\sum_{j=1}^{\infty} \beta_j x_j$  converges in mean-square.

## Model Selection

**Covariates**  $K$  potential predictors  $x_1, \dots, x_K$  available.

**Observe**  $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$ ,  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iK})$ .



## Model Selection

**Covariates**  $K$  potential predictors  $x_1, \dots, x_K$  available.

**Observe**  $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$ ,  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iK})$ .

### Approximating Linear Model

$$y_i = \sum_{j=1}^K x_{ij} \beta_j + b_i + \varepsilon_i, \quad i = 1, 2, \dots, n,$$

$$b_i = \sum_{j=K+1}^{\infty} \beta_j x_j \quad \text{is the approximation error.}$$

## Model Selection

**Covariates**  $K$  potential predictors  $x_1, \dots, x_K$  available.

**Observe**  $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$ ,  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iK})$ .

### Approximating Linear Model

$$y_i = \sum_{j=1}^K x_{ij} \beta_j + b_i + \varepsilon_i, \quad i = 1, 2, \dots, n,$$

$$b_i = \sum_{j=K+1}^{\infty} \beta_j x_j \quad \text{is the approximation error.}$$

**Multiple models** are present.

**Model**  $m$   $\{x_i \mid i \in m\} \mid i = 1, 2, \dots, K\} \subset \{1, 2, \dots, K\}$ .

## A Class of Approximating Models $A$

### The $M \times K$ Incidence Matrix

$$\mathbf{A} = \begin{pmatrix} 1 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 0 & 1 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 0 & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \\ \vdots \\ \mathbf{a}_M^T \end{pmatrix}$$

for the models in  $A$ . The 1's in row  $\mathbf{a}_m$  display the predictors in the  $m$ th model.

## A Class of Approximating Models $A$

### The $M \times K$ Incidence Matrix

$$\mathbf{A} = \begin{pmatrix} 1 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 0 & 1 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 0 & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \\ \vdots \\ \mathbf{a}_M^T \end{pmatrix}$$

for the models in  $A$ . The 1's in row  $\mathbf{a}_m$  display the predictors in the  $m$ th model.

### The Regression Matrix of the Model $m$

$$\mathbf{X}_m = \mathbf{X} \text{diag}(\mathbf{a}_m),$$

## A Class of Approximating Models $A$

### The $M \times K$ Incidence Matrix

$$\mathbf{A} = \begin{pmatrix} 1 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 0 & 1 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 0 & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \\ \vdots \\ \mathbf{a}_M^T \end{pmatrix}$$

for the models in  $A$ . The 1's in row  $\mathbf{a}_m$  display the predictors in the  $m$ th model.

### The Regression Matrix of the Model $m$

$$\mathbf{X}_m = \mathbf{X} \text{diag}(\mathbf{a}_m),$$

$\mathbf{a}_m$  is the vector diagonal entries of  $\text{diag}(\mathbf{a}_m)$ ,  
 $\mathbf{X}$  denotes the  $n \times K$  regression matrix.

## Approximating Model $m$

takes the form

$$\mathbf{y} = \mathbf{X}_m \boldsymbol{\beta}_m + \mathbf{b}_m + \boldsymbol{\varepsilon}.$$

## Approximating Model $m$

takes the form

$$\mathbf{y} = \mathbf{X}_m \boldsymbol{\beta}_m + \mathbf{b}_m + \boldsymbol{\varepsilon}.$$

The LSE of  $\boldsymbol{\beta}_m$

$$\hat{\boldsymbol{\beta}}_m = (\mathbf{X}_m^T \mathbf{X}_m)^+ \mathbf{X}_m^T \mathbf{y}$$

and of  $\boldsymbol{\mu}_m = \mathbf{X}_m \boldsymbol{\beta}_m$

$$\hat{\boldsymbol{\mu}}_m = \mathbf{H}_m \mathbf{y}$$

## Approximating Model $m$

takes the form

$$\mathbf{y} = \mathbf{X}_m \boldsymbol{\beta}_m + \mathbf{b}_m + \boldsymbol{\varepsilon}.$$

The LSE of  $\boldsymbol{\beta}_m$

$$\hat{\boldsymbol{\beta}}_m = (\mathbf{X}_m^T \mathbf{X}_m)^+ \mathbf{X}_m^T \mathbf{y}$$

and of  $\boldsymbol{\mu}_m = \mathbf{X}_m \boldsymbol{\beta}_m$

$$\hat{\boldsymbol{\mu}}_m = \mathbf{H}_m \mathbf{y}$$

under  $m \in M$ , where

$$\mathbf{H}_m = \mathbf{X}_m (\mathbf{X}_m^T \mathbf{X}_m)^+ \mathbf{X}_m^T$$

is a projector.



## Model Average Estimator (MAE)

- ▶ MAE is an alternative to model selection
- ▶ A model selection procedure can be unstable

## Model Average Estimator (MAE)

- ▶ MAE is an alternative to model selection
- ▶ A model selection procedure can be unstable
- ▶ When is combining better than selection?
- ▶ How to measure the uncertainty in selection?

## Model Average Estimator (MAE)

- ▶ MAE is an alternative to model selection
- ▶ A model selection procedure can be unstable
- ▶ When is combining better than selection?
- ▶ How to measure the uncertainty in selection?

*Draper 1995 (JRSS B), Raftery, et. al. 1997 (JASA), Burnham & Anderson 2002 (Book), Hjort & Claeskens 2003 (JASA), Hansen 2007 (Econometrica)*

## Model Average Estimator (MAE)

- ▶ MAE is an alternative to model selection
- ▶ A model selection procedure can be unstable
- ▶ When is combining better than selection?
- ▶ How to measure the uncertainty in selection?

*Draper 1995 (JRSS B), Raftery, et. al. 1997 (JASA), Burnham & Anderson 2002 (Book), Hjort & Claeskens 2003 (JASA), Hansen 2007 (Econometrica)*

## Model Average Estimator (MAE)

- ▶ MAE is an alternative to model selection
- ▶ A model selection procedure can be unstable
- ▶ When is combining better than selection?
- ▶ How to measure the uncertainty in selection?

*Draper 1995 (JRSS B), Raftery, et. al. 1997 (JASA), Burnham & Anderson 2002 (Book), Hjort & Claeskens 2003 (JASA), Hansen 2007 (Econometrica)*

**MAE** of  $\boldsymbol{\beta}$  and  $\boldsymbol{\mu}$

$$\hat{\boldsymbol{\beta}}_w = \sum_{m=1}^M w_m \hat{\boldsymbol{\beta}}_m, \text{ weights } w_i \geq 0 \text{ with } \sum_{m=1}^M w_m = 1$$

$$\hat{\boldsymbol{\mu}}_w = \mathbf{H}_w \mathbf{y}, \quad \mathbf{H}_w = \sum_{m=1}^M w_m \mathbf{H}_m \text{ is the implied hat matrix.}$$

## The Algebraic Structure of MAE

Define

$$\mathbf{AA}^T = \begin{pmatrix} k_1 & k_{12} & \dots & k_{1M} \\ k_{21} & k_2 & \dots & k_{2M} \\ \vdots & \vdots & \ddots & \\ k_{M1} & k_{M2} & \dots & k_M \end{pmatrix} = \mathbf{K}$$

## The Algebraic Structure of MAE

Define

$$\mathbf{A}\mathbf{A}^T = \begin{pmatrix} k_1 & k_{12} & \dots & k_{1M} \\ k_{21} & k_2 & \dots & k_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ k_{M1} & k_{M2} & \dots & k_M \end{pmatrix} = \mathbf{K}$$

**Properties of  $\mathbf{H}_w$ ,** model weights  $\mathbf{w}^T = (w_1, \dots, w_M)$ .

- (i)  $\text{tr}(\mathbf{H}_w) = \sum_{m=1}^M w_m k_m$ .
- (ii)  $\text{tr}(\mathbf{H}_w^2) = \mathbf{w}^T \mathbf{K} \mathbf{w}$ .
- (iii)  $\lambda_M(\mathbf{H}_w) \leq 1$ .

## The Risk under Squared-Error Loss

**Squared-Error Loss**  $L(\mathbf{w}) = \|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_{\mathbf{w}}\|^2.$



## The Risk under Squared-Error Loss

**Squared-Error Loss**  $L(\mathbf{w}) = \|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_{\mathbf{w}}\|^2$ .

**The Conditional Risk** of  $\hat{\boldsymbol{\mu}}_{\mathbf{w}}$

$$\begin{aligned} R(\mathbf{w}) &= E(L(\mathbf{w}) | \mathbf{x}_1 \dots \mathbf{x}_n) \\ &= \|(\mathbf{I} - \mathbf{H}_{\mathbf{w}})\boldsymbol{\mu}\|^2 + \sigma^2 \mathbf{w}^T \mathbf{K} \mathbf{w} \\ &= \mathbf{w}^T (\mathbf{B} + \sigma^2 \mathbf{K}) \mathbf{w}, \end{aligned}$$

## The Risk under Squared-Error Loss

**Squared-Error Loss**  $L(\mathbf{w}) = \|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_{\mathbf{w}}\|^2$ .

**The Conditional Risk** of  $\hat{\boldsymbol{\mu}}_{\mathbf{w}}$

$$\begin{aligned} R(\mathbf{w}) &= E(L(\mathbf{w}) | \mathbf{x}_1 \dots \mathbf{x}_n) \\ &= \|(\mathbf{I} - \mathbf{H}_{\mathbf{w}})\boldsymbol{\mu}\|^2 + \sigma^2 \mathbf{w}^T \mathbf{K} \mathbf{w} \\ &= \mathbf{w}^T (\mathbf{B} + \sigma^2 \mathbf{K}) \mathbf{w}, \end{aligned}$$

$$\mathbf{B} = \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1M} \\ \vdots & \vdots & \ddots & \\ b_{M1} & b_{M2} & \dots & b_{MM} \end{pmatrix}, \text{ with } b_{mk} = \mathbf{b}_m^T (\mathbf{I} - \mathbf{H}_m) (\mathbf{I} - \mathbf{H}_k) \mathbf{b}_k.$$

## The Risk under Squared-Error Loss

**Squared-Error Loss**  $L(\mathbf{w}) = \|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_{\mathbf{w}}\|^2$ .

**The Conditional Risk** of  $\hat{\boldsymbol{\mu}}_{\mathbf{w}}$

$$\begin{aligned} R(\mathbf{w}) &= E(L(\mathbf{w}) | \mathbf{x}_1 \dots \mathbf{x}_n) \\ &= \|(\mathbf{I} - \mathbf{H}_{\mathbf{w}})\boldsymbol{\mu}\|^2 + \sigma^2 \mathbf{w}^T \mathbf{K} \mathbf{w} \\ &= \mathbf{w}^T (\mathbf{B} + \sigma^2 \mathbf{K}) \mathbf{w}, \end{aligned}$$

$$\mathbf{B} = \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1M} \\ \vdots & \vdots & \ddots & \\ b_{M1} & b_{M2} & \dots & b_{MM} \end{pmatrix}, \text{ with } b_{mk} = \mathbf{b}_m^T (\mathbf{I} - \mathbf{H}_m) (\mathbf{I} - \mathbf{H}_k) \mathbf{b}_k.$$

At least two non-zero  $w_i$  in the optimal  $\mathbf{w}$ .

**Example:** Suppose  $M = 2$ ,  $\mathbf{w}^T = (w, 1 - w)$ . Then  $w \in (0, 1)$  unless  $b_{11} = b_{12}$  or  $b_{22} = b_{12}$ .

## Selecting the Model Weights $w_i$

**Mallows' Criterion (MMAE)** for MAE (Hansen 2007)

$$C(\mathbf{w}) = \|(\mathbf{I} - \mathbf{H}_{\mathbf{w}})\mathbf{y}\|^2 + 2\sigma^2 k_{\mathbf{w}}, \quad k_{\mathbf{w}} = \sum_{m=1}^M w_m k_m,$$

## Selecting the Model Weights $w_i$

**Mallows' Criterion (MMAE)** for MAE (Hansen 2007)

$$C(\mathbf{w}) = \|(\mathbf{I} - \mathbf{H}_{\mathbf{w}})\mathbf{y}\|^2 + 2\sigma^2 k_{\mathbf{w}}, \quad k_{\mathbf{w}} = \sum_{m=1}^M w_m k_m,$$

$\sigma^2$  is replaced with an estimate. Select  $\hat{\mathbf{w}}$  such that

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} C(\mathbf{w}).$$

## Selecting the Model Weights $w_i$

**Mallows' Criterion (MMAE)** for MAE (Hansen 2007)

$$C(\mathbf{w}) = \|(\mathbf{I} - \mathbf{H}_{\mathbf{w}})\mathbf{y}\|^2 + 2\sigma^2 k_{\mathbf{w}}, \quad k_{\mathbf{w}} = \sum_{m=1}^M w_m k_m,$$

$\sigma^2$  is replaced with an estimate. Select  $\hat{\mathbf{w}}$  such that

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} C(\mathbf{w}).$$

**Properties of  $C(\mathbf{w})$ :**  $E[C(\mathbf{w})] = E[L(\mathbf{w})] + n\sigma^2$  and

$$\frac{L(\hat{\mathbf{w}})}{\inf_{\mathbf{w}} L(\mathbf{w})} \xrightarrow{p} 1 \quad \text{as } n \rightarrow \infty.$$

## Smoothed AIC and BIC (SAIC & SBIC)

$$w_m = \exp(-\frac{1}{2}AIC_m) / \sum_{i=1}^M \exp(-\frac{1}{2}AIC_i)$$

(*Buckland 1997, Burnham & Anderson 2002*),

$$w_m = \exp(-\frac{1}{2}BIC_m) / \sum_{i=1}^M \exp(-\frac{1}{2}BIC_i),$$

## Smoothed AIC and BIC (SAIC & SBIC)

$$w_m = \exp(-\frac{1}{2}\text{AIC}_m) / \sum_{i=1}^M \exp(-\frac{1}{2}\text{AIC}_i)$$

(*Buckland 1997, Burnham & Anderson 2002*),

$$w_m = \exp(-\frac{1}{2}\text{BIC}_m) / \sum_{i=1}^M \exp(-\frac{1}{2}\text{BIC}_i),$$

the AIC and BIC criteria for model  $m$  are

$$\text{AIC}_m = \ln \hat{\sigma}_m^2 + 2k_m \quad \text{and} \quad \text{BIC}_m = \ln \hat{\sigma}_m^2 + k_m \ln n.$$



## Smoothed AIC and BIC (SAIC & SBIC)

$$w_m = \exp(-\frac{1}{2}\text{AIC}_m) / \sum_{i=1}^M \exp(-\frac{1}{2}\text{AIC}_i)$$

(*Bucland 1997, Burnham & Anderson 2002*),

$$w_m = \exp(-\frac{1}{2}\text{BIC}_m) / \sum_{i=1}^M \exp(-\frac{1}{2}\text{BIC}_i),$$

the AIC and BIC criteria for model  $m$  are

$$\text{AIC}_m = \ln \hat{\sigma}_m^2 + 2k_m \quad \text{and} \quad \text{BIC}_m = \ln \hat{\sigma}_m^2 + k_m \ln n.$$

## Smoothed MDL (SMDL)

$$w_m = \exp(-\text{MDL}_m) / \sum_{i=1}^M \exp(-\text{MDL}_i), \quad \text{where}$$

**Smoothed AIC and BIC (SAIC & SBIC)**

$$w_m = \exp(-\frac{1}{2}\text{AIC}_m) / \sum_{i=1}^M \exp(-\frac{1}{2}\text{AIC}_i)$$

(*Buckland 1997, Burnham & Anderson 2002*),

$$w_m = \exp(-\frac{1}{2}\text{BIC}_m) / \sum_{i=1}^M \exp(-\frac{1}{2}\text{BIC}_i),$$

the AIC and BIC criteria for model  $m$  are

$$\text{AIC}_m = \ln \hat{\sigma}_m^2 + 2k_m \quad \text{and} \quad \text{BIC}_m = \ln \hat{\sigma}_m^2 + k_m \ln n.$$

**Smoothed MDL (SMDL)**

$$w_m = \exp(-\text{MDL}_m) / \sum_{i=1}^M \exp(-\text{MDL}_i), \quad \text{where}$$

$$\text{MDL}_m = n \ln \hat{S}_m^2 + k_m \ln F_m + \ln[k_m(n-k_m)], \quad F_m = \|\hat{\boldsymbol{\mu}}_m\|^2 / k_m \hat{S}_m^2$$

(*Rissanen 2000 & 2007, Liski 2006*)

## Finite Sample Performance

**Simulation Model** is the infinite order regression

$$y_i = \sum_{j=1}^{\infty} \beta_j x_{ji} + \varepsilon_i,$$

- ▶  $x_{ji} \sim N(0, 1)$  iid ( $x_{1i} = 1$ ),  $\varepsilon_i \sim N(0, 1)$  and  $x_{ji} \perp\!\!\!\perp \varepsilon_i$ .

## Finite Sample Performance

**Simulation Model** is the infinite order regression

$$y_i = \sum_{j=1}^{\infty} \beta_j x_{ji} + \varepsilon_i,$$

- ▶  $x_{ji} \sim N(0, 1)$  iid ( $x_{1i} = 1$ ),  $\varepsilon_i \sim N(0, 1)$  and  $x_{ji} \perp\!\!\!\perp \varepsilon_i$ .
- ▶  $\beta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$  and the population  $R^2 = \frac{c^2}{1+c^2}$ .
- ▶  $50 \leq n \leq 1000$  and  $M = 3n^{1/3}$ .
- ▶  $0.5 \leq \alpha \leq 1.5$ , for larger  $\alpha$  the coefficients  $\beta_j$  decline more quickly.
- ▶  $c$  is selected such that  $0.1 \leq R^2 \leq 0.9$ .

## Finite Sample Performance

**Simulation Model** is the infinite order regression

$$y_i = \sum_{j=1}^{\infty} \beta_j x_{ji} + \varepsilon_i,$$

- ▶  $x_{ji} \sim N(0, 1)$  iid ( $x_{1i} = 1$ ),  $\varepsilon_i \sim N(0, 1)$  and  $x_{ji} \perp\!\!\!\perp \varepsilon_i$ .
- ▶  $\beta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$  and the population  $R^2 = \frac{c^2}{1+c^2}$ .
- ▶  $50 \leq n \leq 1000$  and  $M = 3n^{1/3}$ .
- ▶  $0.5 \leq \alpha \leq 1.5$ , for larger  $\alpha$  the coefficients  $\beta_j$  decline more quickly.
- ▶  $c$  is selected such that  $0.1 \leq R^2 \leq 0.9$ .

Mean of predictive loss  $\|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_w\|^2$  over simulations.

# Simulation Results

## Comments

- ▶ AIC and Mallows' C similar, MMAE better than SAIC.

# Simulation Results

## Comments

- ▶ AIC and Mallows' C similar, MMAE better than SAIC.
- ▶ SAIC has lower risk than AIC

# Simulation Results

## Comments

- ▶ AIC and Mallows' C similar, MMAE better than SAIC.
- ▶ SAIC has lower risk than AIC
- ▶ MMAE better than SBIC in most cases.

| Method | Performs well for              |
|--------|--------------------------------|
| SBIC   | $n$ and $R^2$ small, $a$ large |
| BIC    | 'small' models                 |
| AIC    | 'large' models                 |



# Simulation Results

## Comments

- ▶ AIC and Mallows' C similar, MMAE better than SAIC.
- ▶ SAIC has lower risk than AIC
- ▶ MMAE better than SBIC in most cases.

| Method | Performs well for              |
|--------|--------------------------------|
| SBIC   | $n$ and $R^2$ small, $a$ large |
| BIC    | 'small' models                 |
| AIC    | 'large' models                 |

- ▶ SMDL better than MDL.

# Simulation Results

## Comments

- ▶ AIC and Mallows' C similar, MMAE better than SAIC.
- ▶ SAIC has lower risk than AIC
- ▶ MMAE better than SBIC in most cases.

| Method | Performs well for              |
|--------|--------------------------------|
| SBIC   | $n$ and $R^2$ small, $a$ large |
| BIC    | 'small' models                 |
| AIC    | 'large' models                 |

- ▶ SMDL better than MDL.
- ▶ SMDL emulates the best performance criterion.

# Simulation Results






## Comments





- ▶ AIC and Mallows' C similar, MMAE better than SAIC.
- ▶ SAIC has lower risk than AIC
- ▶ MMAE better than SBIC in most cases.

| Method | Performs well for              |
|--------|--------------------------------|
| SBIC   | $n$ and $R^2$ small, $a$ large |
| BIC    | 'small' models                 |
| AIC    | 'large' models                 |

- ▶ SMDL better than MDL.
- ▶ SMDL emulates the best performance criterion.
- ▶ SMDL has the best overall performance.

## References

-  Buckland, S. T., Burnham, K. P. and Augustin, N. H. (1997), Model Selection: An Integral Part of Inference. *Biometrics*, 53, 603–618.
-  Burnham & Anderson (2002), *Model Selection and Multi-model Inference*, Springer
-  Draper, D. (1995), Assessment and Propagation of Model Uncertainty. *Journal of the Royal Statistical Society B*, 57, 45–70.
-  Hansen, B. E. (2007), Least Squares Model Averaging. *Econometrica*, Forthcoming.
-  Hjort, L. H. and Claeskens, G. (2003), Frequentist Model Average Estimators. *Journal of the American Statistical Association*, 98, 879–899.

-  [Liski, E. P. \(2006\)](#), Normalized ML and the MDL Principle for Variable Selection in Linear Regression In: *Festschrift for Tarmo Pukkila on His 60th Birthday*, 159-172.
-  [Raftery, A. E., Madigan, D. and Hoeting, J. A. \(1997\)](#), Bayesian Model Averaging for Regression Models. *Journal of the American Statistical Association*, 92, 179–191.
-  [Rissanen, J. \(2000\)](#). MDL Denoising. *IEEE Trans. Information Theory*, IT-46, pp. 2537–2543.
-  [Rissanen, J. \(2007\)](#), *Information and Complexity in Statistical Modeling*, Springer