

A data-driven lack-of-fit test of regression functions

Erkki P. Liski

University of Tampere
Department of Mathematics and Statistics

Outline

- ▶ The Model
- ▶ Hypothesis
- ▶ Orthogonal series representation

Outline

- ▶ The Model
- ▶ Hypothesis
- ▶ Orthogonal series representation
- ▶ Distributional properties of sample coefficients
- ▶ NML model

Outline

- ▶ The Model
- ▶ Hypothesis
- ▶ Orthogonal series representation
- ▶ Distributional properties of sample coefficients
- ▶ NML model
- ▶ Stochastic complexity
- ▶ MDL criterion
- ▶ Model probabilities

Nonparametric regression model

Observe Y_1, \dots, Y_n from

Nonparametric regression model

Observe Y_1, \dots, Y_n from

the model

$$Y_i = \mu(t_i) + \varepsilon_i, \quad t_i = (i - 1/2)/n, \quad i = 1, \dots, n,$$

$\varepsilon_1, \dots, \varepsilon_n$ are i.i.d. $N(0, 1)$,

where the unknown function μ is defined on $[0, 1]$.

Hypothesis

We wish to test the **no-effect hypothesis**

$$H_0 : \mu(t) = \beta_0 \text{ for all } t \in [0, 1],$$

where $\beta_0 \in \mathbb{R}$ is an unknown constant.

Hypothesis

We wish to test the **no-effect hypothesis**

$$H_0 : \mu(t) = \beta_0 \text{ for all } t \in [0, 1],$$

where $\beta_0 \in \mathbb{R}$ is an unknown constant.

The test will rely on an orthogonal series representation for μ :

$$\mu(t) = \beta_0 + \sum_{j=1}^{\infty} \beta_j \chi_j(t), \quad x \in [0, 1], \quad (1)$$

Hypothesis

We wish to test the **no-effect hypothesis**

$$H_0 : \mu(t) = \beta_0 \text{ for all } t \in [0, 1],$$

where $\beta_0 \in \mathbb{R}$ is an unknown constant.

The test will rely on an orthogonal series representation for μ :

$$\mu(t) = \beta_0 + \sum_{j=1}^{\infty} \beta_j x_j(t), \quad x \in [0, 1], \quad (1)$$

where

$$\beta_j = \int_0^1 x_j(t) \mu(t) dt, \quad j = 0, 1, \dots$$

and

$\{1, x_1, x_2, \dots\}$ is an orthonormal basis for $\mu \in L_2[0, 1]$.

Example: An example of such a basis is

$$x_j(t) = \sqrt{2} \cos(\pi j t), \quad j = 1, 2, \dots$$

Example: An example of such a basis is

$$x_j(t) = \sqrt{2} \cos(\pi jt), \quad j = 1, 2, \dots$$

Wavelets and orthogonal polynomials are other examples.

Example: An example of such a basis is

$$x_j(t) = \sqrt{2} \cos(\pi j t), \quad j = 1, 2, \dots$$

Wavelets and orthogonal polynomials are other examples.

The basis functions are also assumed to be orthonormal with respect to the design:

$$\sum_{i=1}^n x_j(t_i) x_k(t_i) = \begin{cases} 0, & j \neq k \\ n, & j = k \end{cases}$$

for all $j, k \in \{0, 1, \dots\}$ and $x_0 \equiv 1$.

Example: An example of such a basis is

$$x_j(t) = \sqrt{2} \cos(\pi j t), \quad j = 1, 2, \dots$$

Wavelets and orthogonal polynomials are other examples.

The basis functions are also assumed to be orthonormal with respect to the design:

$$\sum_{i=1}^n x_j(t_i) x_k(t_i) = \begin{cases} 0, & j \neq k \\ n, & j = k \end{cases}$$

for all $j, k \in \{0, 1, \dots\}$ and $x_0 \equiv 1$.

Note 1 The cosine basis is orthonormal with respect to the design.

Note 2 The representation (1) has infinitely many parameters β_0, β_1, \dots

Alternative models $\mathcal{A} = \{M_1, \dots, M_K\}$ are of the form

$$\mu_j(t) = \beta_0 + \sum_{k \in \mathcal{K}_j} \beta_k x_k(t),$$

where \mathcal{K}_j is a subset of $\{1, \dots, j\}$ and $K < n$.

Alternative models $\mathcal{A} = \{M_1, \dots, M_K\}$ are of the form

$$\mu_j(t) = \beta_0 + \sum_{k \in \mathcal{K}_j} \beta_k x_k(t),$$

where \mathcal{K}_j is a subset of $\{1, \dots, j\}$ and $K < n$.

A common version of such a model is

$$\mu_j(t) = \beta_0 + \sum_{k=1}^j \beta_k x_k(t).$$

Alternative models $\mathcal{A} = \{M_1, \dots, M_K\}$ are of the form

$$\mu_j(t) = \beta_0 + \sum_{k \in \mathcal{K}_j} \beta_k x_k(t),$$

where \mathcal{K}_j is a subset of $\{1, \dots, j\}$ and $K < n$.

A common version of such a model is

$$\mu_j(t) = \beta_0 + \sum_{k=1}^j \beta_k x_k(t).$$

The MLE of β_k is

$$\hat{\beta}_k = \frac{1}{n} \sum_{i=1}^n x_k(t_i) Y_i, \quad k = 0, 1, \dots, n-1.$$

Distributional properties

Sample coefficients satisfy:

1. $\hat{\beta}_0, \dots, \hat{\beta}_{n-1}$ are mutually independent.
2. $\hat{\beta}_k \sim N\left(\frac{1}{n} \sum_{i=1}^n \mu(t_i) x_k(t_i), \frac{1}{n}\right), \quad k = 0, 1, \dots, n-1.$

Distributional properties Sample coefficients satisfy:

1. $\hat{\beta}_0, \dots, \hat{\beta}_{n-1}$ are mutually independent.
2. $\hat{\beta}_k \sim N\left(\frac{1}{n} \sum_{i=1}^n \mu(t_i) x_k(t_i), \frac{1}{n}\right), \quad k = 0, 1, \dots, n-1.$

Under H_0

$$\sqrt{n} \hat{\beta}_k \sim N(0, 1), \quad k = 1, \dots, n-1.$$

If we assume merely independence of errors
1. and 2. continue to hold asymptotically.

Normalized Maximum Likelihood (NML) model

Denote the class of normal densities for M_j , $j=0,1,\dots,K$

$$\mathcal{M}_j = \{f(\mathbf{y}; \boldsymbol{\beta}_j); \boldsymbol{\beta}_j \in \Theta_j \subset \mathbb{R}^{k_j}\}$$

where

$$\boldsymbol{\beta}_j = (\beta_0, \beta_1, \dots, \beta_{k_j})', \quad \mathbf{y} = (y_1, \dots, y_n)'$$

Normalized Maximum Likelihood (NML) model

Denote the class of normal densities for M_j , $j=0,1,\dots,K$

$$\mathcal{M}_j = \{f(\mathbf{y}; \boldsymbol{\beta}_j); \boldsymbol{\beta}_j \in \Theta_j \subset \mathbb{R}^{k_j}\}$$

where

$$\boldsymbol{\beta}_j = (\beta_0, \beta_1, \dots, \beta_{k_j})', \quad \mathbf{y} = (y_1, \dots, y_n)'$$

The shortest "codelength" of data \mathbf{y} with \mathcal{M}_j is

$$\log \frac{1}{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}, \quad \hat{\boldsymbol{\beta}}_j(\mathbf{y}) \text{ is the MLE of } \boldsymbol{\beta}_j.$$

Normalized Maximum Likelihood (NML) model

Denote the class of normal densities for M_j , $j=0,1,\dots,K$

$$\mathcal{M}_j = \{f(\mathbf{y}; \boldsymbol{\beta}_j); \boldsymbol{\beta}_j \in \Theta_j \subset \mathbb{R}^{k_j}\}$$

where

$$\boldsymbol{\beta}_j = (\beta_0, \beta_1, \dots, \beta_{k_j})', \quad \mathbf{y} = (y_1, \dots, y_n)'$$

The shortest "codelength" of data \mathbf{y} with \mathcal{M}_j is

$$\log \frac{1}{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}, \quad \hat{\boldsymbol{\beta}}_j(\mathbf{y}) \text{ is the MLE of } \boldsymbol{\beta}_j.$$

If we use $q(\mathbf{y})$, the excess code length is

$$\log \frac{1}{q(\mathbf{y})} - \log \frac{1}{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))} = \log \frac{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}{q(\mathbf{y})}.$$

Minimax solution

The **NML density function** for \mathcal{M}_j is

$$\hat{f}(\mathbf{y}; j) = \frac{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}{\int_{\hat{\boldsymbol{\beta}}_j(\mathbf{x}) \in \Theta_j} f(\mathbf{x}; \hat{\boldsymbol{\beta}}_j(\mathbf{x})) d\mathbf{x}}.$$

Minimax solution

The NML density function for \mathcal{M}_j is

$$\hat{f}(\mathbf{y}; j) = \frac{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}{\int_{\hat{\boldsymbol{\beta}}_j(\mathbf{x}) \in \Theta_j} f(\mathbf{x}; \hat{\boldsymbol{\beta}}_j(\mathbf{x})) d\mathbf{x}}.$$

It solves the minimax problem

$$\min_q \max_{\mathbf{y}} \log \frac{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}{q(\mathbf{y})}$$

with the unique solution $\hat{q} = \hat{f}(\cdot; j)$.
(Shtarkov 1987, Rissanen 1996).

Maxmin solution

The NML density $\hat{f}(\cdot; j)$ is also the solution of the maxmin problem

$$\max_g \min_q E_g \log \frac{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}{q(\mathbf{y})}$$

(Rissanen 2001 and 2007).

NML Density for \mathcal{M}_j

$$\hat{f}(\mathbf{y}; j) = \frac{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}{C(j)},$$

where the numerator is the LF evaluated at the MLE:

$$f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y})) = (2\pi)^{-n/2} \exp\left[-\frac{1}{2} \sum_{i=1}^n (y_i - \hat{\beta}_0)^2 + \frac{n}{2} \sum_{k \in \mathcal{K}_j} \hat{\beta}_k^2\right]$$

NML Density for \mathcal{M}_j

$$\hat{f}(\mathbf{y}; j) = \frac{f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y}))}{C(j)},$$

where the numerator is the LF evaluated at the MLE:

$$f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y})) = (2\pi)^{-n/2} \exp\left[-\frac{1}{2} \sum_{i=1}^n (y_i - \hat{\beta}_0)^2 + \frac{n}{2} \sum_{k \in \mathcal{K}_j} \hat{\beta}_k^2\right]$$

and the denominator

$$C(j) = \int_{\hat{\boldsymbol{\beta}}_j(\mathbf{x}) \in \Theta_j} f(\mathbf{x}; \hat{\boldsymbol{\beta}}_j(\mathbf{x})) \, d\mathbf{x}$$

is the parametric complexity of the model \mathcal{M}_j .

The MDL codelength for \mathbf{y} with \mathcal{M}_j :

$$\log \frac{1}{\hat{f}(\mathbf{y}; j)} = -\log \hat{f}(\mathbf{y}; j) = -\log f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y})) + \log C(\mathcal{M}_j).$$

The MDL codelength for \mathbf{y} with \mathcal{M}_j :

$$\log \frac{1}{\hat{f}(\mathbf{y}; j)} = -\log \hat{f}(\mathbf{y}; j) = -\log f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y})) + \log C(\mathcal{M}_j).$$

is the *stochastic complexity* of data \mathbf{y} relative to \mathcal{M}_j .

The MDL codelength for \mathbf{y} with \mathcal{M}_j :

$$\log \frac{1}{\hat{f}(\mathbf{y}; j)} = -\log \hat{f}(\mathbf{y}; j) = -\log f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y})) + \log C(\mathcal{M}_j).$$

is the *stochastic complexity* of data \mathbf{y} relative to \mathcal{M}_j .

It can be shown that

$$-2 \log \hat{f}(\mathbf{y}; j) = k_j \left(\log \frac{n \|\hat{\boldsymbol{\beta}}_j\|^2}{k_j \sigma_n^2} - \frac{n \|\hat{\boldsymbol{\beta}}_j\|^2}{k_j \sigma_n^2} + 1 \right) + \log k_j + a(\mathbf{y}),$$

where $\sigma_n^2 = \frac{1}{n}$ and $a(\mathbf{y})$ is common to all \mathcal{M}_j .

Denote

$$MDL_j = -2 \log \hat{f}(\mathbf{y}; j).$$

The MDL codelength for \mathbf{y} with \mathcal{M}_j :

$$\log \frac{1}{\hat{f}(\mathbf{y}; j)} = -\log \hat{f}(\mathbf{y}; j) = -\log f(\mathbf{y}; \hat{\boldsymbol{\beta}}_j(\mathbf{y})) + \log C(\mathcal{M}_j).$$

is the *stochastic complexity* of data \mathbf{y} relative to \mathcal{M}_j .

It can be shown that

$$-2 \log \hat{f}(\mathbf{y}; j) = k_j \left(\log \frac{n \|\hat{\boldsymbol{\beta}}_j\|^2}{k_j \sigma_n^2} - \frac{n \|\hat{\boldsymbol{\beta}}_j\|^2}{k_j \sigma_n^2} + 1 \right) + \log k_j + a(\mathbf{y}),$$

where $\sigma_n^2 = \frac{1}{n}$ and $a(\mathbf{y})$ is common to all \mathcal{M}_j .

Denote

$$MDL_j = -2 \log \hat{f}(\mathbf{y}; j).$$

Model probabilities

Given the data \mathbf{y} , $\hat{f}(\mathbf{y}; j)$ can be interpreted as the likelihood of the model \mathcal{M}_j , $j = 1, 2, \dots, K$

Model probabilities

Given the data \mathbf{y} , $\hat{f}(\mathbf{y}; j)$ can be interpreted as the likelihood of the model \mathcal{M}_j , $j = 1, 2, \dots, K$

This leads to the *NML distribution* for models

$$\hat{p}(j; \mathbf{y}) = \frac{\hat{f}(\mathbf{y}; j)}{\sum_{i=1}^K \hat{f}(\mathbf{y}; i)} = \frac{\exp(-MDL_j/2)}{\sum_{i=1}^K \exp(-MDL_i/2)}.$$

Model probabilities

Given the data \mathbf{y} , $\hat{f}(\mathbf{y}; j)$ can be interpreted as the likelihood of the model \mathcal{M}_j , $j = 1, 2, \dots, K$

This leads to the *NML distribution* for models

$$\hat{p}(j; \mathbf{y}) = \frac{\hat{f}(\mathbf{y}; j)}{\sum_{i=1}^K \hat{f}(\mathbf{y}; i)} = \frac{\exp(-MDL_j/2)}{\sum_{i=1}^K \exp(-MDL_i/2)}.$$

Given the data \mathbf{y} , we may compute the probability

$$p_0(\mathbf{y}) = \hat{p}(\mathcal{M}_0; \mathbf{y}).$$

Model probabilities

Given the data \mathbf{y} , $\hat{f}(\mathbf{y}; j)$ can be interpreted as the likelihood of the model \mathcal{M}_j , $j = 1, 2, \dots, K$

This leads to the *NML distribution* for models

$$\hat{p}(j; \mathbf{y}) = \frac{\hat{f}(\mathbf{y}; j)}{\sum_{i=1}^K \hat{f}(\mathbf{y}; i)} = \frac{\exp(-MDL_j/2)}{\sum_{i=1}^K \exp(-MDL_i/2)}.$$

Given the data \mathbf{y} , we may compute the probability

$$p_0(\mathbf{y}) = \hat{p}(\mathcal{M}_0; \mathbf{y}).$$

The idea of the test:

reject \mathcal{M}_0 at level α if $p_0(\mathbf{y})$ less than the α -quantile of $p_0(\mathbf{y})$'s null distribution.

Distribution of $p_0(\mathbf{y})$

A related test based on BIC. (Aerts & Claeskens 2004).

Results on the asymptotic distribution of $p_0(\mathbf{y})$ under

- ▶ the null hypothesis and
- ▶ local alternatives.

Distribution of $p_0(\mathbf{y})$

A related test based on BIC. (Aerts & Claeskens 2004).





Results on the asymptotic distribution of $p_0(\mathbf{y})$ under





- ▶ the null hypothesis and
- ▶ local alternatives.

Here the distribution of $p_0(\mathbf{y})$ seems to be more complicated.

A topic for further research.

References

-  [Aerts, Claeskens & Hart \(2004\)](#), *Bayesian-motivated tests of function fit and their asymptotic frequentist properties*, *Annals of Statistics* 32, pp. 2580–2615.
-  [Burnham & Anderson \(2002\)](#), *Model Selection and Multi-model Inference*, Springer
-  [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.
-  [Rissanen, J. \(1996\)](#). Fisher Information and Stochastic Complexity. *IEEE Trans. Information Theory*, IT-42, pp. 40–47.

-  [Rissanen, J. \(2000\)](#). MDL Denoising. *IEEE Trans. Information Theory*, IT-46, pp. 2537–2543.
-  [Rissanen, J. \(2001\)](#). Strong Optimality of the NML Models as Universal Codes and Information in Data. *IEEE Trans. Information Theory*, IT-47, pp. 1712–1717.
-  [Rissanen, J. \(2007\)](#), *Information and Complexity in Statistical Modeling*, Springer
-  [Shtarkov, Yu. M. \(1987\)](#). Universal Sequential Coding of Single Messages. *Problems of Information Transmission*, 23, pp. 3–17.