

# Data-driven goodness-of-fit tests

Mikhail Langovoy

Institute for Mathematical Stochastics,  
Georg-August-University of Göttingen

July 9, 2007

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

- 1 Introduction
  - Main ideas
    - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

Constructing good tests is an essential problem of statistics.

Two approaches:

- *naive*: "distance" between theoretical and empirical distribution is proposed as statistic
- *aiming optimality*: construct tests which are asymptotically efficient (Neyman, Le Cam, Wald)

# Examples of Distance-Based Statistics

- Kolmogorov - Smirnov:

$$D_n = \sqrt{n} \|F_n - F\|_\infty$$

- Cramer - von Mises:

$$\omega_n^2 = n \int_{-\infty}^{\infty} (F_n(t) - F_0(t))^2 dF_0(t)$$

- many other

# About Distance-Based Tests

- These tests works
- Each one is asymptotically optimal only in a very few directions of alternatives

- hypothesis  $H_0 : X \sim U[0, 1]$
- $\{\phi_j\}_{j=0}^{\infty}$  orthonormal basis of  $L_2([0, 1], \lambda)$
- 

$$N_k = \sum_{j=1}^k \left\{ n^{-\frac{1}{2}} \sum_{i=1}^n \phi_j(X_i) \right\}^2$$



- hypothesis  $H_0 : X \sim U[0, 1]$
- $\{\phi_j\}_{j=0}^{\infty}$  orthonormal basis of  $L_2([0, 1], \lambda)$
- 

$$N_k = \sum_{j=1}^k \left\{ n^{-\frac{1}{2}} \sum_{i=1}^n \phi_j(X_i) \right\}^2$$

- hypothesis  $H_0 : X \sim U[0, 1]$
- $\{\phi_j\}_{j=0}^{\infty}$  orthonormal basis of  $L_2([0, 1], \lambda)$
- 

$$N_k = \sum_{j=1}^k \left\{ n^{-\frac{1}{2}} \sum_{i=1}^n \phi_j(X_i) \right\}^2$$

- 1 Introduction
  - Main ideas
  - **Data-driven tests**
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

# Idea of Selection Rule.

- model dimension  $k$  was known fixed in advance
- important to select the right model dimension
  - incorrect choice can decrease the power of a test
- **Solution:** incorporate the test statistic by some procedure which chooses the right dimension automatically by the data

# Idea of Selection Rule.

- model dimension  $k$  was known fixed in advance
- important to select the right model dimension
  - incorrect choice can decrease the power of a test
- **Solution:** incorporate the test statistic by some procedure which chooses the right dimension automatically by the data

# Idea of Selection Rule.

- model dimension  $k$  was known fixed in advance
- important to select the right model dimension
  - incorrect choice can decrease the power of a test
- **Solution:** incorporate the test statistic by some procedure which chooses the right dimension automatically by the data

Data-driven score tests are

- asymptotically optimal in an infinite number of directions
- show good overall performance in practice

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions**
  - **Framework**
  - **Selection rule**
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests



- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions**
  - Framework**
  - Selection rule
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

- $X_1, X_2, \dots$  random variables, values in a measurable space  $\mathbb{X}$
- $X_1, \dots, X_m$  have joint distribution  $P_m$  from the family of distributions  $\mathbb{P}_m$  - for every  $m$
- function  $\mathcal{F}$  acting from  $\otimes_{m=1}^{\infty} \mathbb{P}_m = (\mathbb{P}_1, \mathbb{P}_2, \dots)$  to a known set  $\Theta$
- $\mathcal{F}(P_1, P_2, \dots) = \theta$



$$H_0 : \theta \in \Theta_0 \subset \Theta$$



$$H_A : \theta \in \Theta_1 = \Theta \setminus \Theta_0$$

- observations  $Y_1, \dots, Y_n$ , values in a measurable space  $\mathbb{Y}$
- not necessarily on the basis of  $X_1, \dots, X_m$  !

- $\Theta$  *any* set
- $Y_1, \dots, Y_n$  can be dependent or nonidentically distributed
- $\mathbb{Y}$  can be infinite dimensional
- additional assumptions on  $Y_i$ 's will be imposed when necessary

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - **Selection rule**
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

## Definition

- $M_k$  statistical model for specific problem
- $\Theta_k$  is the parameter set of  $M_k$
- the family of models  $M_k$  for  $k = 1, 2, \dots$  is **nested** if  $\Theta_1 \subseteq \Theta_2 \subseteq \dots$
  
- We do not require  $\Theta'_k$ s to be finite dimensional
- We do not require that all  $\Theta'_k$ s are different

$T_k$  arbitrary statistic for testing validity of  $M_k$  on the basis of  $Y_1, \dots, Y_n$

## Definition

- nested family of models  $M_k$  for  $k = 1, \dots, d(n)$
- $d(n)$  control sequence
- choose  $\pi(\cdot, \cdot) : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{R}$
- assume
  - $\pi(1, n) < \pi(2, n) < \dots < \pi(d(n), n)$  for all  $n$
  - $\pi(j, n) - \pi(1, n) \rightarrow \infty$  as  $n \rightarrow \infty$  for  $j = 2, \dots, d(n)$

Call  $\pi(j, n)$  a **penalty** attributed to  $j$ th model  $M_j$  and sample size  $n$ .



## Definition

A **selection rule**  $S$  for the sequence of statistics  $\{T_k\}$  is

$$S = \min\{k : 1 \leq k \leq d(n); T_k - \pi(k, n) \geq T_j - \pi(j, n), j = 1, \dots, d(n)\}$$

Call  $T_S$  a **data-driven** test statistic for testing validity of the initial model.

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics**
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics**
  - **Definitions**
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

## Definition

- observations  $Y_1, \dots, Y_n$ , values in a measurable space  $\mathbb{Y}$
- $k$  fixed number
- $l = (l_1, \dots, l_k)$  vector-function
- $l_i : \mathbb{Y} \rightarrow \mathbb{R}$  for  $i = 1, \dots, k$  are known Lebesgue measurable

## Definition



$$L = \{E_0[I(Y)]^T I(Y)\}^{-1}$$

- $E_0$  is with respect to  $P_0$
- $P_0$  is the d.f. of some (fixed and known) random variable  $Y$
- $Y$  has values in  $\mathbb{Y}$

### • assume

- $E_0 I(Y) = 0$
- $L$  is well defined

## Definition

$$T_k = \left\{ \frac{1}{\sqrt{n}} \sum_{j=1}^n I(Y_j) \right\} L \left\{ \frac{1}{\sqrt{n}} \sum_{j=1}^n I(Y_j) \right\}^T$$

$T_k$  - *statistic of Neyman's type* (**NT-statistic**)

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics**
  - Definitions
  - Examples**
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

$l_1, \dots, l_k$  can be

- some score functions
- any other functions, depending on the problem
  - truncated, penalized or partial likelihood
  - possible to use  $l_1, \dots, l_k$  unrelated to any likelihood



- deconvolution
- building block for many complicated statistical inverse problems

- instead of  $X_i$  one observes  $Y_i$

$$Y_i = X_i + \varepsilon_i$$

- $\varepsilon_i$ 's are i.i.d. with a known density  $h$
- $X_i$  and  $\varepsilon_i$  are independent for each  $i$
- $E \varepsilon_i = 0, 0 < E \varepsilon^2 < \infty$
- $X$  has a density

$H_0$  :  $X$  has density  $f_0$

- choose for every  $k \leq d(n)$  an auxiliary parametric family  $\{f_\theta\}$
- $\theta \in \Theta \subseteq \mathbb{R}^k$
- $f_\theta$  from this family coincides with  $f_0$  from the null hypothesis  $H_0$
- the true  $F$  possibly has no relation to the chosen  $\{f_\theta\}$

$$I(y) = \frac{\frac{\partial}{\partial \theta} \left( \int_{\mathbb{R}} f_{\theta}(s) h(y-s) ds \right) \Big|_{\theta=0}}{\int_{\mathbb{R}} f_0(s) h(y-s) ds}$$

- define  $U_k$  as above
- $U_k$  is an NT-statistic

# Rank Tests for Independence

- $(X_1, Y_1), \dots, (X_n, Y_n)$  i.i.d. random variables
- distribution function  $D$
- continuous marginal distribution functions  $F$  and  $G$

$$H_0 : D(x, y) = F(x)G(y), \quad x, y \in \mathbb{R},$$

# Rank Tests for Independence

- $b_j$  denote the  $j$ -th orthonormal Legendre polynomial  
 $b_1(x) = \sqrt{3}(2x - 1)$ ,  $b_2(x) = \sqrt{5}(6x^2 - 6x + 1)$ , ...
- $R_i$  the rank of  $X_i$  among  $X_1, \dots, X_n$
- $S_i$  the rank of  $Y_i$  among  $Y_1, \dots, Y_n$

$$T_k = \sum_{j=1}^k \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n b_j \left( \frac{R_i - 1/2}{n} \right) b_j \left( \frac{S_i - 1/2}{n} \right) \right\}^2$$

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics**
  - Definitions
  - Examples
  - Alternative distribution**
  - Null distribution
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

# Additional assumptions

- $Y_1, Y_2, \dots$  are identically distributed, with d.f.  $P$
- do not assume  $Y_1, Y_2, \dots$  are independent
- possibly,  $X_1, X_2, \dots$  are dependent and nonidentically distributed



# Additional assumptions

- $E_P I(Y)$  exists

- $$\frac{1}{n} \sum_{j=1}^n I(Y_j) \rightarrow E_P I(Y) \quad \text{in } P\text{-probability as } n \rightarrow \infty$$

- $$n^{-1/2} \sum_{j=1}^n (I(Y_j) - E_P I(Y)) \rightarrow_d \mathcal{N}(0, L^{-1})$$

⟨C⟩ there exists integer  $K = K(P) \geq 1$  such that

$$E_P I_1(Y) = 0, \dots, E_P I_{K-1}(Y) = 0, E_P I_K = C_P \neq 0$$

## Theorem

*Let the above assumptions holds and*

$$\lim_{n \rightarrow \infty} \sup_{k \leq d(n)} \frac{\pi(k, n)}{n \lambda_k^{(k)}} = 0.$$

*Then*

$$\lim_{n \rightarrow \infty} P(S \geq K) = 1.$$

# Additional assumptions

suppose there exists a sequence  $\{r_n\}_{n=1}^{\infty}$  such that

- $\lim_{n \rightarrow \infty} r_n = \infty$



$$\langle \mathbf{A} \rangle \quad P \left( \frac{1}{n} \left| \sum_{i=1}^n [l_K(Y_i) - E_P l_K(Y_i)] \right| \geq y \right) = O \left( \frac{1}{r_n} \right)$$

## Theorem

*Let the above conditions holds and*

$$d(n) = o(r_n) \quad \text{as } n \rightarrow \infty.$$

*Then  $T_S \rightarrow_P \infty$  as  $n \rightarrow \infty$ .*

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 **NT-statistics**
  - Definitions
  - Examples
  - Alternative distribution
  - **Null distribution**
- 4 GNT-statistics
  - Definitions
  - Consistency of GNT-tests

- $\{T_k\}$  sequence of NT-statistics
- $S$  selection rule
- $\pi(k, n)$  penalty

## Definition

$\pi(k, n)$  is of **proper weight**, if there exists real sequences

$\{s(k, n)\}_{k,n=1}^{\infty}$ ,  $\{t(k, n)\}_{k,n=1}^{\infty}$ , such that



$$\lim_{n \rightarrow \infty} \sup_{k \leq u_n} \frac{s(k, n)}{n \lambda_k^{(k)}} = 0,$$

$\{u_n\}_{n=1}^{\infty}$  real sequence,  $\lim_{n \rightarrow \infty} u_n = \infty$

•  $\lim_{n \rightarrow \infty} t(k, n) = \infty$  for every  $k \geq 2$

$\lim_{k \rightarrow \infty} t(k, n) = \infty$  for every fixed  $n$



## Definition

- $s(k, n) \leq \pi(k, n) - \pi(1, n) \leq t(k, n)$  for all  $k, n$



$$\lim_{n \rightarrow \infty} \sup_{k \leq m_n} \frac{\pi(k, n)}{n \lambda_k^{(k)}} = 0,$$

$\{m_n\}_{n=1}^{\infty}$  real sequence,  $\lim_{n \rightarrow \infty} m_n = \infty$ .

define for  $l = (l_1, \dots, l_k)$

- $$\bar{l}_j := \frac{1}{n} \sum_{i=1}^n l_j(Y_i)$$

- $$\bar{l} := (\bar{l}_1, \bar{l}_2, \dots, \bar{l}_k)$$

- $$Q_k(\bar{l}) = (\bar{l}_1, \bar{l}_2, \dots, \bar{l}_k) L(\bar{l}_1, \bar{l}_2, \dots, \bar{l}_k)^T$$

## Definition

- statistic  $T_k$ , selection rule  $S$ , penalty of proper weight
- Assume there exists a Lebesgue measurable  $\varphi(\cdot, \cdot) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ 
  - monotonically decreasing in the second argument
  - monotonically nondecreasing in the first argument
  - and

## Definition

- (B2) for every  $\varepsilon > 0$  there exists  $K = K_\varepsilon$  such that for  $n > n(\varepsilon)$

$$\sum_{k=K_\varepsilon}^{u_n} \varphi(k; s(k, n)) < \varepsilon$$

- (B) for all  $k \geq 1$  and  $y \in [s(k, n); t(k, n)]$

$$P_0(n Q_k(\bar{I}) \geq y) \leq \varphi(k; y)$$

We call  $\varphi$  a **proper majorant** for (large deviations of)  $T_k$ .

## Theorem

- $\{T_k\}$  sequence of NT-statistics
- $S$  selection rule, with penalty of proper weight
- large deviations of  $T_k$  are properly majorated



$$d(n) \leq \min\{u_n, m_n\}.$$

Then  $S = O_{P_0}(1)$  and  $T_S = O_{P_0}(1)$ .

## Theorem

- $\{T_k\}$  sequence of NT-statistics
- $S$  selection rule, with penalty of proper weight
- the above regularity assumptions are satisfied
- $d(n) = o(r_n)$ ,  $d(n) \leq \min\{u_n, m_n\}$ .

Then  $T_S$  is consistent against any (fixed) alternative  $P$  satisfying (C).

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 GNT-statistics**
  - **Definitions**
  - **Consistency of GNT-tests**

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 **GNT-statistics**
  - **Definitions**
  - Consistency of GNT-tests



- The notion of NT-statistics is helpful if the null hypothesis is simple
- for composite hypotheses, the concept needs to be modified

## Definition

- observations  $Y_1, \dots, Y_n$ , values in a measurable space  $\mathbb{Y}$
- for simplicity, identically distributed
- $k$  fixed number
- $l = (l_1, \dots, l_k)$  vector-function
- $l_i : \mathbb{Y} \rightarrow \mathbb{R}$  for  $i = 1, \dots, k$  (maybe **unknown**) Lebesgue measurable

$$L^{(0)} = \{E_0[l(Y)]^T l(Y)\}^{-1}$$

- $P_0$  (possibly unknown) d.f. of  $Y$  under the null hypothesis

- assume  $E_0 I(Y) = 0$
- $L^{(0)}$  is well-defined
- $L_k$  denote for every  $k$  a  $k \times k$  symmetric positive definite matrix with finite elements
- 

$$\|L_k - L^{(0)}\| = o_{P_0}(1)$$

## Definition

$I_1^*, \dots, I_n^*$  **sufficiently good** estimators of  $I(Y_1), \dots, I(Y_n)$  w.r.t.  $P_0$ , if for every  $\varepsilon > 0$

$$P_0^n \left( \frac{1}{\sqrt{n}} \left\| \sum_{i=1}^n (I_j^* - I(Y_j)) \right\| \geq \varepsilon \right) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

## Definition

$$GT_k = \left\{ \frac{1}{\sqrt{n}} \sum_{j=1}^n I_j^* \right\} L_k \left\{ \frac{1}{\sqrt{n}} \sum_{j=1}^n I_j^* \right\}^T$$

$GT_k$  *generalized statistic of Neyman's type* (**GNT-statistic**)

$GT_S$  *data-driven GNT-statistic*

- 1 Introduction
  - Main ideas
  - Data-driven tests
- 2 General notions
  - Framework
  - Selection rule
- 3 NT-statistics
  - Definitions
  - Examples
  - Alternative distribution
  - Null distribution
- 4 **GNT-statistics**
  - Definitions
  - **Consistency of GNT-tests**

## Theorem

- $\{GT_k\}$  sequence of GNT-statistics
- $S$  selection rule, with penalty of proper weight (for  $R_k$ )
- large deviations of  $GT_k$  are properly majorated

$$d(n) \leq \min\{u_n, m_n\}.$$

Then under the null hypothesis  $S = O_{P_0}(1)$  and  $GT_S = O_{P_0}(1)$ .

# Possible additional assumptions



$$\langle \mathbf{C1} \rangle \quad \|L - L^{(0)}\| = o_P(1)$$

- $l_1^*, \dots, l_n^*$  are sufficiently good estimators of  $l(Y_1), \dots, l(Y_n)$  with respect to  $P$

- i.e., for every  $\varepsilon > 0$

$$P^n \left( \frac{1}{\sqrt{n}} \left\| \sum_{i=1}^n (l_j^* - l(Y_j)) \right\| \geq \varepsilon \right) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$



# Possible additional assumptions



$$\langle \mathbf{C1} \rangle \quad \|L - L^{(0)}\| = o_P(1)$$

- $l_1^*, \dots, l_n^*$  are sufficiently good estimators of  $l(Y_1), \dots, l(Y_n)$  with respect to  $P$

- i.e., for every  $\varepsilon > 0$

$$P^n \left( \frac{1}{\sqrt{n}} \left\| \sum_{i=1}^n (l_i^* - l(Y_i)) \right\| \geq \varepsilon \right) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

## Theorem

- $\{GT_k\}$  sequence of GNT-statistics
- $S$  selection rule
- $d(n) = o(r_n)$ .

*Under the above assumptions,  $T_S$  is consistent against any (fixed) alternative distribution  $P$  satisfying  $\langle C \rangle$  and  $\langle C1 \rangle$ .*

- **SNT-statistics**
- Selection rule: examples
- Application of NT-statistics
- GNT-statistics: examples

## Definition

Consider statistic of the form

$$T_k = \sum_{j=1}^k \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n I_j(Y_i) \right\}^2$$

- $n$  is the number of available observations  $Y_1, \dots, Y_n$
- $I_j : \mathbb{Y} \rightarrow \mathbb{R}$  are known Lebesgue measurable

We call  $T_k$  the *simplified statistic of Neyman's type* ( **SNT-statistic** ).

# Partial likelihood (Cox)

- generalizes conditional and marginal likelihood
- applications include inference in stochastic processes

- random variable  $Y$ , density  $f_Y(y; \theta)$
- let  $Y$  be transformed into the sequence

$$(X_1, S_1, X_2, S_2, \dots, X_m, S_m)$$

- components may themselves be vectors

- full likelihood

$$\prod_{j=1}^m f_{X_j|X^{(j-1)}, S^{(j-1)}}(x_j|x^{(j-1)}, s^{(j-1)}; \theta) \prod_{j=1}^m f_{S_j|X^{(j)}, S^{(j-1)}}(s_j|x^{(j)}, s^{(j-1)}; \theta)$$

$x^{(j)} = (x_1, \dots, x_j)$  and  $s^{(j)} = (s_1, \dots, s_j)$

- second product is *partial likelihood*

- for simplicity  $\theta$  is real
- $H_0 : \theta = \theta_0$
- define for  $j = 1, \dots, m$

$$t_j = \left. \frac{\partial \log f_{S_j | X^{(j)}, S^{(j-1)}}(s_j | x^{(j)}, s^{(j-1)}; \theta)}{\partial \theta} \right|_{\theta = \theta_0}$$

and  $\sigma_j^2 := \text{var}(t_j)$



- define  $l_j := t_j / \sigma_j$
- SNT-statistic

$$PL_m = \sum_{j=1}^m \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n l_j \right\}^2$$

- SNT-statistics
- **Selection rule: examples**
- Application of NT-statistics
- GNT-statistics: examples

- $PL_m$  depends on the number  $m$  of components
- suppose  $Y$  can be transformed into sequences  $(X_1, S_1)$ , or  $(X_1, S_1, X_2, S_2)$ , or even  $(X_1, S_1, X_2, S_2, \dots, X_m, S_m)$  for any  $m$
- which  $m$  is the best choice?
- $PL_S$  chooses the reasonable number of components automatically

- framework of Gaussian linear processes
- includes
  - Gaussian regression with fixed design
  - Gaussian sequences
  - model of Ibragimov and Has'minskii

$\mathbb{S}$  linear subspace of some Hilbert space  $\mathbb{H}$

## Definition

*Gaussian linear process* on  $\mathbb{S}$  with mean  $s \in \mathbb{H}$  and variance  $\varepsilon^2$  any process  $Y$  indexed by  $\mathbb{S}$  of the form

$$Y(t) = \langle s, t \rangle + \varepsilon Z(t),$$

for all  $t \in \mathbb{S}$

$Z$  linear isonormal process indexed by  $\mathbb{S}$  (i.e.  $Z$  is a centered and linear Gaussian process with covariance structure  $E[Z(t)Z(u)] = \langle t, u \rangle$ )

Birgé and Massart considered estimation of  $s$  in this model.

- $S$  finite dimensional subspace of  $\mathbb{S}$
- $\gamma(t) = \|t\|^2 - 2Y(t)$
- define the projection estimator on  $S$  to be the minimizer of  $\gamma(t)$  with respect to  $t \in S$

# Gaussian model selection

Given

- a finite or countable family  $\{S_m\}_{m \in \mathcal{M}}$  of finite dimensional linear subspaces of  $S$
- the corresponding family of projection estimators  $\hat{s}_m$ , built for the same realization of process  $Y$
- a nonnegative function  $pen$  defined on  $\mathcal{M}$

estimate  $s$  by  $\tilde{s} = \hat{s}_{\hat{m}}$ , where  $\hat{m}$  is any minimizer with respect to  $m \in \mathcal{M}$  of the penalized criterion

$$crit(m) = -\|\hat{s}_m\|^2 + pen(m) = \gamma(\hat{s}_m) + pen(m)$$

- $\gamma(t)$  is the least squares criterion
- $\hat{s}_m$  is the least squares estimator of  $s$ , which is in this case the maximum likelihood estimator
- $\|\hat{s}_m\|^2$  is the Neyman score for testing  $s = 0$  within this model



- $pen$  is a selection rule in our terminology
- $\|\widehat{s}_{\widehat{m}}\|^2$  is the data-driven SNT-statistics
- by the Consistency Theorem,  $\|\widehat{s}_{\widehat{m}}\|^2$  can be used for testing  $s = 0$  and has a good range of consistency

- possible to define  $PL_m$  for  $\theta$  multidimensional or even infinite-dimensional
- under additional regularity assumptions  $E(t_j) = 0$
- $PL_m$  NT-statistic, but not SNT-statistic

## Lemma

*Let  $I_K(Y_j)$ 's be bounded i.i.d. random variables with finite variance  $\sigma^2$ . Then  $\langle A \rangle$  is satisfied with  $r_n = \exp(ny^2/2\sigma)$ .*

- SNT-statistics
- Selection rule: examples
- **Application of NT-statistics**
- GNT-statistics: examples

## Theorem

- $Y_1, \dots, Y_n$  i.i.d.
- $\{T_k\}$  family of SNT-statistics,  $S$  selection rule
- $E I(Y_1) = 0$
- $(I_1(Y_i), \dots, I_k(Y_i))$  has unit covariance matrix for every  $k$
- $\|(I_1(Y_1), \dots, I_k(Y_1))\|_k \leq M(k)$  a.e.

## Theorem

- $\pi(k, n) - \pi(1, n) \geq 2k$  for all  $k \geq 2$



$$\lim_{n \rightarrow \infty} \frac{M(d(n)) \pi(d(n), n)}{\sqrt{n}} = 0.$$

Then  $S = O_{P_0}(1)$  and  $T_S = O_{P_0}(1)$ .

## Theorem

- $T_S$  Neyman's smooth data-driven test statistic
- $\pi(k, n) - \pi(1, n) \geq 2k$  for all  $k \geq 2$  and for all  $k \leq d(n)$

- 

$$\lim_{n \rightarrow \infty} \frac{d(n)\pi(d(n), n)}{\sqrt{n}} = 0.$$

Then  $S = O_{P_0}(1)$  and  $T_S = O_{P_0}(1)$ .

## Theorem

- $d(n) = o\left(\left\{\frac{n}{\log n}\right\}^{1/10}\right)$
- $\mathbb{P}$  *alternative*
- $F$  and  $G$  the marginal distribution functions of  $X$  and  $Y$  under  $\mathbb{P}$
- for some  $j$

$$E_{\mathbb{P}} b_j(F(X))b_j(G(Y)) \neq 0.$$

If  $d(n) \rightarrow \infty$ , then  $T_S \rightarrow \infty$  as  $n \rightarrow \infty$ .



# Rank tests for independence

- condition  $\langle C \rangle$  : for some  $j$

$$E_{\mathbb{P}} b_j\left(\frac{R_1 - 1/2}{n}\right) b_j\left(\frac{S_1 - 1/2}{n}\right) \neq 0$$

- for continuous  $F$  and  $G$  this is asymptotically equivalent to

$$E_{\mathbb{P}} b_j(F(X)) b_j(G(Y)) \neq 0.$$

- SNT-statistics
- Selection rule: examples
- Application of NT-statistics
- **GNT-statistics: examples**

# Smooth test for composite hypotheses

- $X_1, \dots, X_n$  i.i.d. with density  $f(\mathbf{x})$
- $H_0: f(\mathbf{x}) \in \{f(\mathbf{x}; \beta), \beta \in \mathcal{B}\}$
- $\mathcal{B} \subset \mathbb{R}^q$
- $\{f(\mathbf{x}; \beta), \beta \in \mathcal{B}\}$  given family of densities

# Smooth test for composite hypotheses

- $F$  distribution function corresponding to  $f$



$$Y_n(\beta) = n^{-1} \sum_{i=1}^n (\phi_1(F(X_i; \beta)), \dots, \phi_j(F(X_i; \beta)))^T$$

- $j$  depends on the context
- $I$  the  $k \times k$  identity matrix

# Smooth test for composite hypotheses

- $$I_{\beta} = \left\{ -E_{\beta} \frac{\partial}{\partial \beta_t} \phi_j(F(X_j; \beta)) \right\}_{t=1, \dots, q; j=1, \dots, k}$$

- $$I_{\beta\beta} = \left\{ -E_{\beta} \frac{\partial^2}{\partial \beta_t \partial \beta_u} \log f(X; \beta) \right\}_{t=1, \dots, q; u=1, \dots, q}$$

- $$R(\beta) = I_{\beta}^T (I_{\beta\beta} - I_{\beta} I_{\beta}^T) I_{\beta}$$

# Smooth test for composite hypotheses

- $\hat{\beta}$  maximum likelihood estimator of  $\beta$  under  $H_0$
- the score statistic

$$W_k(\hat{\beta}) = n Y_n^T(\hat{\beta}) \{I + R(\hat{\beta})\} Y_n(\hat{\beta})$$

- in a regular situation (Cox and Hinkley),  $W_k(\hat{\beta})$  is a GNT-statistic