DETECTION OF VERY SHORT SIGNAL SEGMENTS

Ming Yuan

Morgridge Institute for Research and Department of Statistics

University of Wisconsin-Madison

myuan@stat.wisc.edu http://www.stat.wisc.edu/~myuan





(Joint work with Tony Cai)

COPY NUMBER VARIATION



[Xie and Tammi, 2009]

PEAK DETECTION



Structured signal detection

1. Introduction

ASTRONOMY



STRUCTURED SIGNAL DETECTION

$$X_i = \mu_i + \varepsilon_i, \qquad \qquad i = 1, 2, \dots, n$$

Without signal

$$\mu_1 = \mu_2 = \ldots = \mu_n = 0$$

▶ With signals located at unknown segments $S = \{S_j = (a_j, b_j]\}$

$$\mu_i = \begin{cases} f_j((i-a_j)/(b_j-a_j)) & \text{if } i \in S_j \\ 0 & \text{otherwise} \end{cases}$$

• Short signal $-|S_j| < n^{\xi}$ for some $\xi < 1$

▶ Examples

- Biology Copy number variation
- Engineering Peak detection
- Astronomy Detecting planets (see, e.g., Fabrycky et al., 2012)

•

DETECTION OF SIGNALS

• Hypothesis testing $H_0: \mu_i = 0$ vs

$$H_a: \mu_i = \begin{cases} f((i-a)/(b-a)) & \text{if } i \in S := (a,b] \\ 0 & \text{otherwise} \end{cases}$$

- ▶ Questions of interests
 - When is a signal detectable?
 - How to detect a "detectable" signal?
- ▶ Determining factors
 - Signal amplitude $A = ||f||_{L_2}$
 - Signal duration -d = b a

EFFECT OF AMPLITUDE



EFFECT OF DURATION



OPTIMAL RATES OF DETECTION



- ▶ Signal of known shape
 - $A_n^2 \asymp d_n^{-1} \log n$
- ▶ Signal of arbitrary shape

$$A_n^2 \asymp d_n^{-1} \log + \sqrt{d_n^{-1} \log n}$$

- ▶ Signal of smooth shape
 - Hölder class with $\alpha = 1$
 - $\alpha = 1/5$

SIGNALS WITH KNOWN SHAPE

$$f_n(\cdot) = A_n f_0(\cdot)$$
 where $\int_0^1 f_0^2 = 1$

Constant signals – Arias-Castro et al. (2004), Jeng et al. (2004)....
Known position and duration – Likelihood ratio type test

$$L_{(a_n, a_n + d_n]} := \left(\sum_i f_0^2(i/d_n)\right)^{-1/2} \sum_{i=1}^{d_n} X_{a_n + i} f_0\left(\frac{i}{d_n}\right)$$

▶ Unknown position and duration – scan statistic

$$L_n = \max_{a_n, d_n} L_{(a_n, a_n + d_n]} \Longrightarrow T_n := \begin{cases} \text{reject } H_0 & \text{if } L_n \ge c_n \\ \text{accept } H_0 & \text{otherwise} \end{cases}$$

▶ With critical value $c_n = c_0 \sqrt{\log n}$, T_n can detect any signal such that

 $A_n^2 \gtrsim d_n^{-1} \log(n/d_n)$

OPTIMALITY

Optimal rate of detection $-A_n^2 \simeq d_n^{-1} \log(n/d_n)$

• Detection is achieveable with $A_n^2 \gtrsim d_n^{-1} \log(n/d_n)$

$$\mathbb{P}(T_n = 1|H_0) + \sup_{H_1} \mathbb{P}(T_n = 0|H_1) \to 0.$$

▶ For signals with $A_n^2 \leq d_n^{-1} \log(n/d_n)$, any test is powerless in that

$$\inf_{\tilde{T}} \left\{ \mathbb{P}(\tilde{T}=1|H_0) + \sup_{H_1} \mathbb{P}(\tilde{T}=0|H_1) \right\} \to 1.$$

► Effect of multiplicity

- Known location $-d_n^{-1}$
- Price for not knowing the location $-\log n$

SIGNALS WITH ARBITRARY SHAPE

▶ Evidence based on quadratic statistic

$$L_{(a_n, a_n+d_n]} := 2^{-1/2} (d_n^{1/2} + \log^{1/2} n)^{-1} \sum_{i=1}^{d_n} (X_{a_n+i}^2 - 1)$$

▶ Scan over all possible segments

$$L_n = \max_{a_n, d_n} L_{(a_n, a_n + d_n]} \Longrightarrow T_n := \begin{cases} \text{reject } H_0 & \text{if } L_n \ge c_n \\ \text{accept } H_0 & \text{otherwise} \end{cases}$$

▶ T_n $(c_n = c_0 \sqrt{\log n})$ can detect any signal such that

$$A_n^2 \gtrsim \left(\frac{\log n}{d_n}\right) + \left(\frac{\log n}{d_n}\right)^{1/2}$$

- For short signals $(d_n \ll \log n) A_n^2 \gtrsim d_n^{-1} \log n$
- For long signals $(d_n \gg \log n) A_n^2 \gtrsim (d_n^{-1} \log n)^{1/2}$

OPTIMALITY

Optimal rate of detection $-A_n^2 \simeq d_n^{-1} \log(n) + (d_n^{-1} \log(n))^{1/2}$

▶ Every test is powerless if

$$A_n^2 \lesssim \left(\frac{\log n}{d_n}\right) + \left(\frac{\log n}{d_n}\right)^{1/2}$$

• Adversarial case – signal is random $\pm A_n$

▶ Comparison with signals of known shape

	Known	Arbitrary
$d \lesssim \log n$	$A_n^2 \gtrsim d^{-1}\log n$	
$d\gtrsim \log n$	$A_n^2 \gtrsim d^{-1}\log n$	$A_n^2 \gtrsim (d^{-1}\log n)^{1/2}$

• Effect of multiplicity – only $\sqrt{\log n}$ when $d \gtrsim \log n$

LINEAR VS QUADRATIC SCAN





13

Smooth Signals

• f_n is α times differentiable in that

$$|f^{(\lfloor \alpha \rfloor)}(x) - f^{(\lfloor \alpha \rfloor)}(x')| \le L|x - x'|^{\alpha - \lfloor \alpha \rfloor}$$

- ▶ Optimal rate of detection
 - Determining factors amplitude, duration and degree of smoothness
- ▶ How to scan for smooth signal combining linear and quadratic statistics
 - Linear statistics most powerful if signal is almost constant
 - Quadratic statistics most powerful if signal changes rapidly

SCAN FOR SMOOTH SIGNALS



- ▶ Divide the segment (a_n, b_n] into l bins, each of size m = d_n/l
- ▶ Average within each bin

$$Y_j := m^{-1/2} \sum_{i=1}^m X_{a_n + (j-1)m+i}$$

▶ Gather evidence through

$$(l^{1/2} + \log^{1/2} n)^{-1} \sum_{j=1}^{l} (Y_j^2 - 1)$$

▶ Scan over all putative intervals – $L_n = \max_{a_n, d_n} L_{(a_n, b_n]}$

HOW MANY BINS?



If
$$d_n \le \log n$$

• if $\alpha \geq 1/4$

$$l = d_n$$

• If
$$\log n < d_n \le (\log n)^{2\alpha+1}$$

$$l = \log n$$

• If
$$d_n > (\log n)^{2\alpha+1}$$

$$l = d_n^{\frac{2}{4\alpha+1}} (\log n)^{-\frac{1}{4\alpha+1}}$$

If α < 1/4 − an extra change
 If d_n > (log n)^{1/1−4α}

 $l = d_n$

OPTIMAL RATE OF DETECTION

With critical value $c_n = c_0 \sqrt{\log n}$, detect α times differentiable signals such that

 $\begin{array}{l} \bullet \mbox{ if } \alpha \geq 1/4 \\ A_n^2 \gtrsim \left\{ \begin{array}{ll} d_n^{-1} \log n & \mbox{ if } d_n = O((\log n)^{2\alpha+1}) \\ d_n^{-\frac{4\alpha}{4\alpha+1}} (\log n)^{\frac{2\alpha}{4\alpha+1}} & \mbox{ if } d_n \gg (\log n)^{2\alpha+1} \end{array} \right. \\ \bullet \mbox{ if } \alpha < 1/4 \\ A_n^2 \gtrsim \left\{ \begin{array}{ll} d_n^{-1} \log n & \mbox{ if } d_n = O((\log n)^{2\alpha+1}) \\ d_n^{-\frac{4\alpha}{4\alpha+1}} (\log n)^{\frac{2\alpha}{4\alpha+1}} & \mbox{ if } (\log n)^{2\alpha+1} \ll d_n \ll (\log n)^{1/(1-4\alpha)} \\ (d_n^{-1} \log n)^{1/2} & \mbox{ if } d_n \gg (\log n)^{1/(1-4\alpha)} \end{array} \right.$

Compared with the case when the location of signal is known in advance (Ingster)

$$A_n^2 \gtrsim d_n^{-\frac{4\alpha}{4\alpha+1}}$$

EFFECT OF "BINNING"





18

SUMMARY

- Detection of sparse signal segments is a common problem in many high dimensional problems
- Detectability of sparse segments depends on the signal strength, duration, and shape
- ▶ Rate optimal detection can be achieved with scan statistics
 - Scan with linear statistics for signal of known shape
 - Scan with quadratic statistics for signal of arbitrary shape
 - Combing linear and quadratic statistics for signals of smoothness
- Ongoing work
 - Adaptation
 - Beyond normality