



Sequential Change Detection via Denoising Score Matching

Liyan Xie

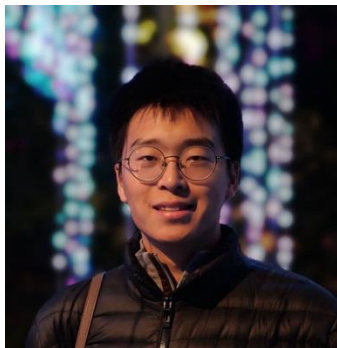
Department of Industrial and Systems Engineering

University of Minnesota

Allerton 2025



This is joint work with



Wenbin Zhou
Carnegie Mellon University



Woody (Shixiang) Zhu
Carnegie Mellon University



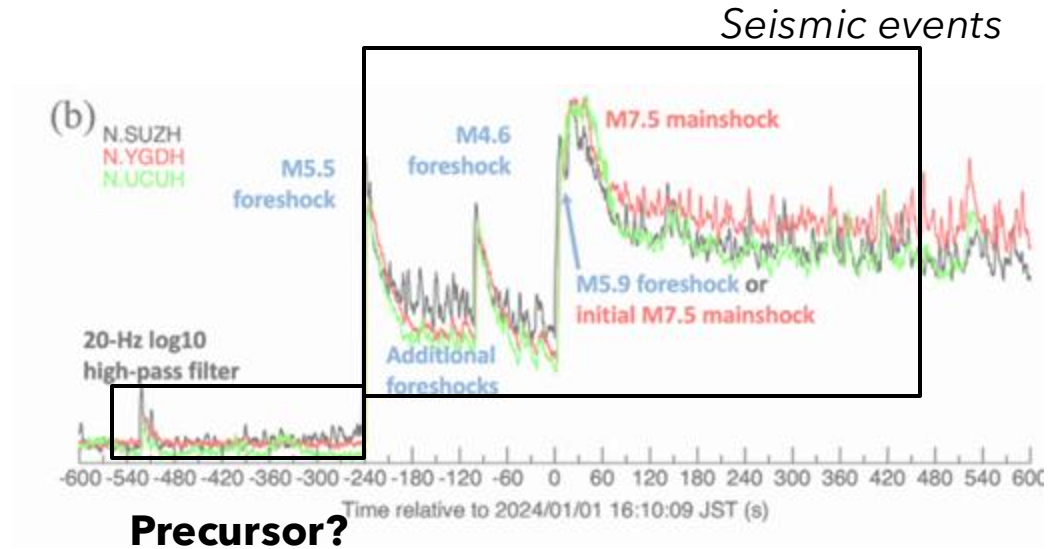
Zhigang Peng
Georgia Institute of Technology

Outline

- Motivating example and problem setup
- Proposed method: Denoising Score Matching CUSUM
- Theoretical analysis
- Simulation and real data results

Motivating example: earthquake precursor

Modern data acquisition enables high-resolution catalogues and high-dimensional geophysical monitoring signals to be collected.



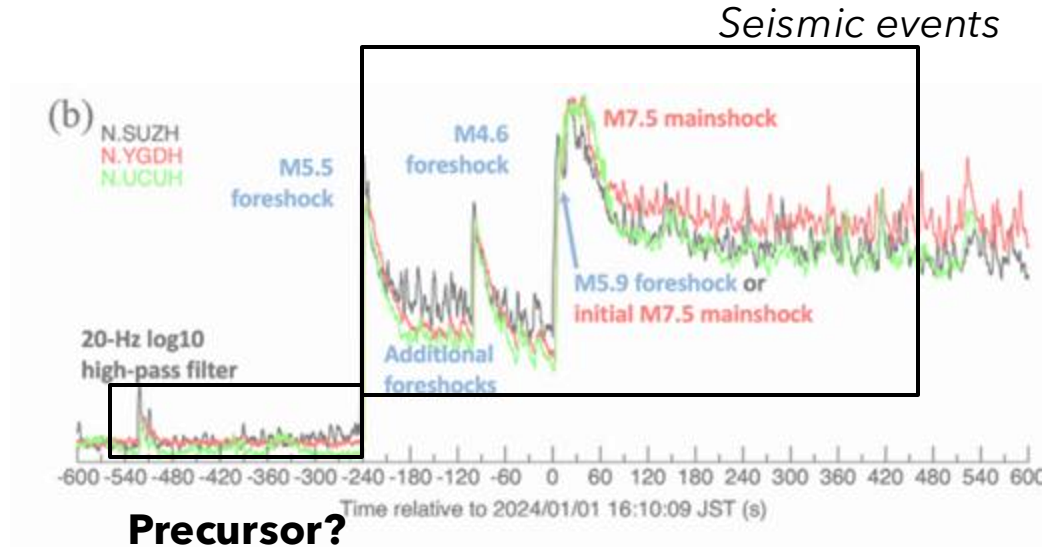
Precursor?
(signals detected prior to an earthquake)

Motivating example: earthquake precursor

Modern data acquisition enables high-resolution catalogues and high-dimensional geophysical monitoring signals to be collected.

Challenges:

- Signals embedded in high-dimensional data are hard to capture
- The distribution before/after change are **unknown** and **difficult to estimate**.

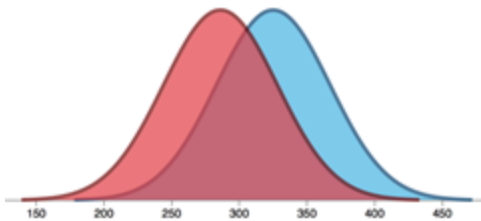


(signals detected prior to an earthquake)

Insights

Idea: capture the distributional differences **without** estimating the distribution densities

Difference of distributions



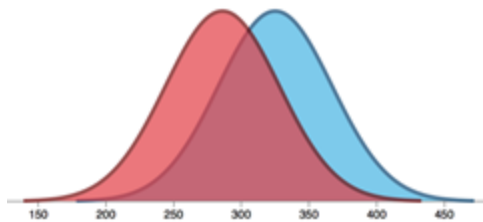
$$p_1(x) \text{ vs. } p_0(x)$$

Estimating densities p_0, p_1 may be hard

Insights

Idea: capture the distributional differences **without** estimating the distribution densities

Difference of distributions

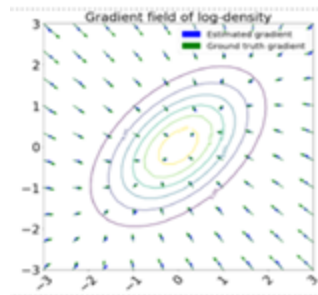


$p_1(x)$ vs. $p_0(x)$

Estimating densities p_0, p_1 may be hard



Difference of the **gradient of the log density** (score)



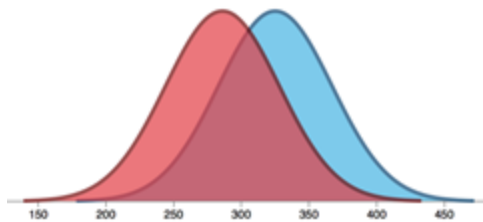
$s_1(x) = \nabla_x \log p_1(x)$ vs. $s_0(x) = \nabla_x \log p_0(x)$

Estimating scores s_0, s_1 are typically **easier**

Insights

Idea: capture the distributional differences **without** estimating the distribution densities

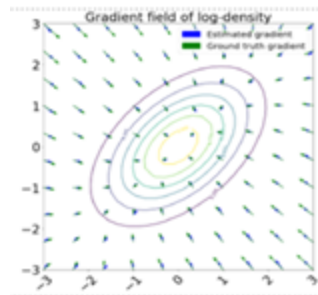
Difference of distributions



$p_1(x)$ vs. $p_0(x)$

Estimating densities p_0, p_1 may be hard

Difference of the **gradient of the log density** (score)

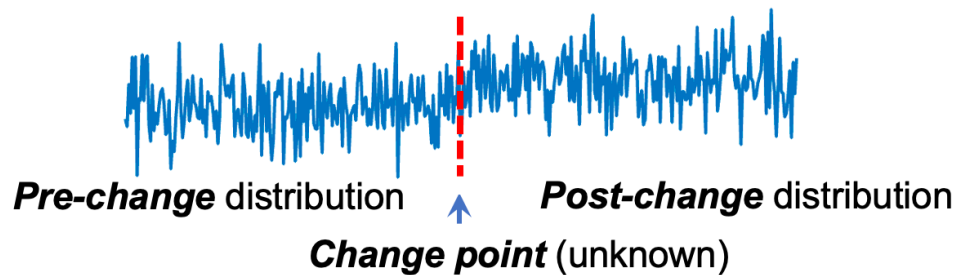


$s_1(x) = \nabla_x \log p_1(x)$ vs. $s_0(x) = \nabla_x \log p_0(x)$

Estimating scores s_0, s_1 are typically **easier**

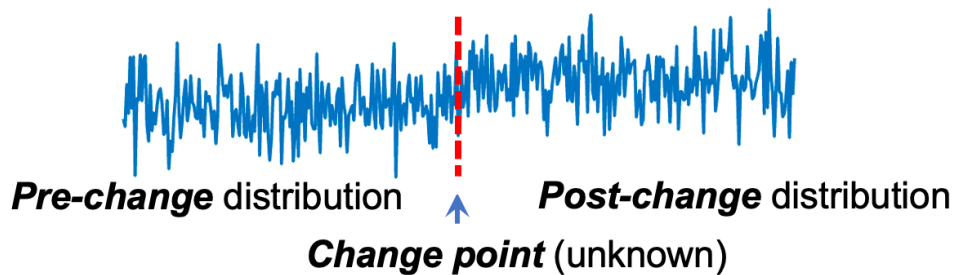
Example (unnormalized models): $p_\theta(x) = \exp(F_\theta(x))/Z_\theta \iff s_\theta(x) = F_\theta'(x)$

Problem setup



$$x_1, x_2, \dots, x_{\tau-1} \sim^{iid} p_0 \text{ (unknown)}, \quad x_{\tau}, x_{\tau+1}, \dots \sim^{iid} p_1 \text{ (unknown)}$$

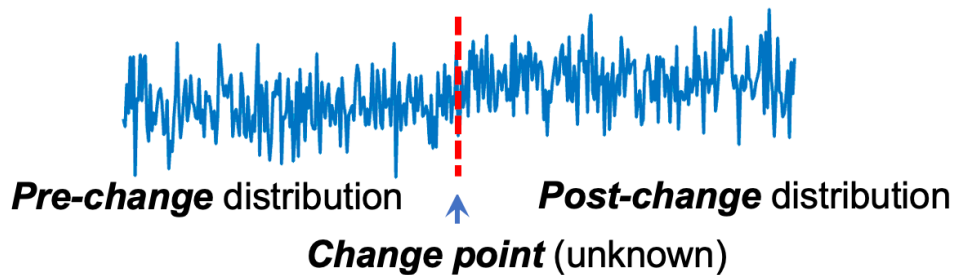
Problem setup



$$x_1, x_2, \dots, x_{\tau-1} \sim^{iid} p_0 \text{ (unknown)}, \quad x_{\tau}, x_{\tau+1}, \dots \sim^{iid} p_1 \text{ (unknown)}$$

- Assume reference datasets D_0 and D_1 sampled from p_0 and p_1 (typically available by pre-collected historical data).

Problem setup



$$x_1, x_2, \dots, x_{\tau-1} \sim^{iid} p_0 \text{ (unknown)}, \quad x_{\tau}, x_{\tau+1}, \dots \sim^{iid} p_1 \text{ (unknown)}$$

- Assume reference datasets D_0 and D_1 sampled from p_0 and p_1 (typically available by pre-collected historical data).
- **Goal:** detect the unknown change-point τ as quickly as possible.

Main component: Score-Matching

When the density p is either unavailable or difficult to estimate, we can estimate its score function instead.

- Score: $s(x) = \nabla_x \log p(x)$
- Can be reliability estimated via Score-Matching:

Main component: Score-Matching

When the density p is either unavailable or difficult to estimate, we can estimate its score function instead.

- Score: $s(x) = \nabla_x \log p(x)$
- Can be reliability estimated via Score-Matching:
 - Given training data $x_1, \dots, x_N \sim p(x)$

Main component: Score-Matching

When the density p is either unavailable or difficult to estimate, we can estimate its score function instead.

- Score: $s(x) = \nabla_x \log p(x)$
- Can be reliability estimated via Score-Matching:
 - Given training data $x_1, \dots, x_N \sim p(x)$
 - Estimate the score function by

$$\hat{s}(x) = s(x; \hat{\theta}) = \arg \min_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N [\text{tr}(\nabla_x s(x_i; \theta)) + \frac{1}{2} \|s(x_i; \theta)\|^2]$$

Main component: Score-Matching

When the density p is either unavailable or difficult to estimate, we can estimate its score function instead.

- Score: $s(x) = \nabla_x \log p(x)$
- Can be reliability estimated via Score-Matching:
 - Given training data $x_1, \dots, x_N \sim p(x)$
 - Estimate the score function by

$$\hat{s}(x) = s(x; \hat{\theta}) = \arg \min_{\theta} \hat{f}(\theta) = \frac{1}{N} \sum_{i=1}^N [\text{tr}(\nabla_x s(x_i; \theta)) + \frac{1}{2} \|s(x_i; \theta)\|^2]$$

- $s(x; \theta)$ is typically parameterized as a neural-network

Main component: Score-Matching

When the density p is either unavailable or difficult to estimate, we can estimate its score function instead.

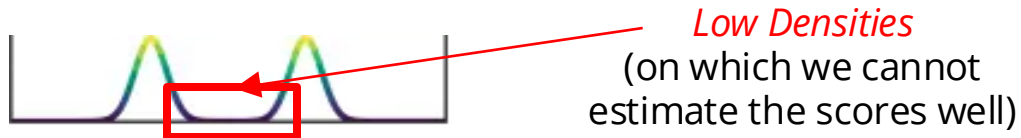
- Score: $s(x) = \nabla_x \log p(x)$
- Can be reliability estimated via Score-Matching:
 - Given training data $x_1, \dots, x_N \sim p(x)$
 - Estimate the score function by

$$\hat{s}(x) = s(x; \hat{\theta}) = \arg \min_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N [\text{tr}(\nabla_x s(x_i; \theta)) + \frac{1}{2} \|s(x_i; \theta)\|^2]$$

- $s(x; \theta)$ is typically parameterized as a neural-network
- Theoretically guaranteed to be consistent estimator

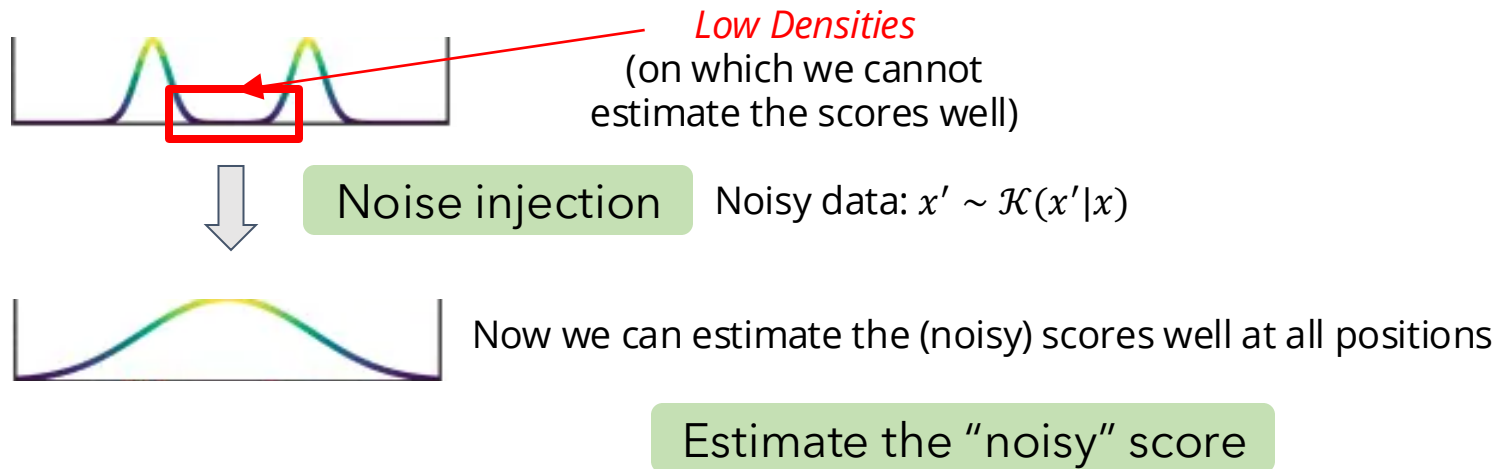
Denoising Score Matching (DSM)

Vanilla score matching is not good enough for high-dim data



Denoising Score Matching (DSM)

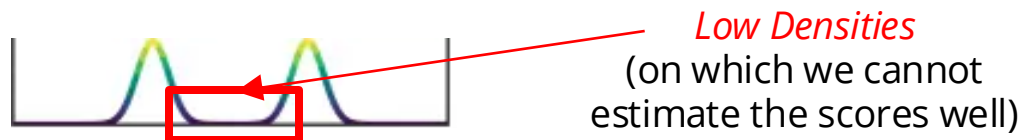
Vanilla score matching is not good enough for high-dim data



$$\mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim \mathcal{K}(\cdot|x)} \|s'(x'; \theta) - \nabla_{x'} \log \mathcal{K}(x'|x)\|_2^2. \text{ Minimize MSE}$$

Denoising Score Matching (DSM)

Vanilla score matching is not good enough for high-dim data



Noise injection

Noisy data: $x' \sim \mathcal{K}(x'|x)$



Now we can estimate the (noisy) scores well at all positions

Estimate the "noisy" score

$$\mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim \mathcal{K}(\cdot|x)} \left\| \boxed{s'(x'; \theta)} - \boxed{\nabla_{x'} \log \mathcal{K}(x'|x)} \right\|_2^2 \cdot \text{Minimize MSE}$$

Model *Noise*
(Score network)

Our Algorithm: DSM-CUSUM

- I. Estimate score function \hat{s}_0, \hat{s}_1 via DSM using reference data D_0 (pre-change) and D_1 (post-change)

Our Algorithm: DSM-CUSUM

- I. Estimate score function \hat{s}_0, \hat{s}_1 via DSM using reference data D_0 (pre-change) and D_1 (post-change)
- II. Recursively compute the detection statistics:
 - Initialize $\mathcal{S}_0 = 0$; For $t = 1, 2, 3, \dots$

$$\Delta(x_t) = H(x_t; \hat{s}_0) - H(x_t; \hat{s}_1)$$

$$\mathcal{S}_t = \mathcal{S}_{t-1}^+ + \Delta(x_t)$$

Our Algorithm: DSM-CUSUM

- I. Estimate score function \hat{s}_0, \hat{s}_1 via DSM using reference data D_0 (pre-change) and D_1 (post-change)
- II. Recursively compute the detection statistics:

- Initialize $\mathcal{S}_0 = 0$; For $t = 1, 2, 3, \dots$

$$H(x; s) = \text{div } s(x) + \frac{1}{2} \|s(x)\|_2^2$$

Difference of (Hyvarinen) score function

$$\Delta(x_t) = H(x_t; \hat{s}_0) - H(x_t; \hat{s}_1)$$

“distributional difference”

$$\mathcal{S}_t = \mathcal{S}_{t-1}^+ + \Delta(x_t)$$

Our Algorithm: DSM-CUSUM

- I. Estimate score function \hat{s}_0, \hat{s}_1 via DSM using reference data D_0 (pre-change) and D_1 (post-change)
- II. Recursively compute the detection statistics:

- Initialize $\mathcal{S}_0 = 0$; For $t = 1, 2, 3, \dots$

$$H(x; s) = \text{div } s(x) + \frac{1}{2} \|s(x)\|_2^2$$

Difference of (Hyvarinen) score function

$$\Delta(x_t) = H(x_t; \hat{s}_0) - H(x_t; \hat{s}_1) \quad \text{"distributional difference"}$$

$$\mathcal{S}_t = \mathcal{S}_{t-1}^+ + \Delta(x_t) \quad \text{Accumulate the evidences}$$

Our Algorithm: DSM-CUSUM

I. Estimate score function \hat{s}_0, \hat{s}_1 via DSM using reference data D_0 (pre-change) and D_1 (post-change)

II. Recursively compute the detection statistics:

➤ Initialize $\mathcal{S}_0 = 0$; For $t = 1, 2, 3, \dots$

$$H(x; s) = \text{div } s(x) + \frac{1}{2} \|s(x)\|_2^2$$

Difference of (Hyvarinen) score function

$$\Delta(x_t) = H(x_t; \hat{s}_0) - H(x_t; \hat{s}_1) \quad \text{“distributional difference”}$$

$$\mathcal{S}_t = \mathcal{S}_{t-1}^+ + \Delta(x_t) \quad \text{Accumulate the evidences}$$

III. Raise alarm when \mathcal{S}_t exceeds some threshold value:

$$T = \inf \{t \in \mathbb{N} : \mathcal{S}_t \geq \tau\}$$

Extension to online score matching

- When the post-change training data is not available
- Sequentially estimate the post-change score function:

Extension to online score matching

- When the post-change training data is not available
- Sequentially estimate the post-change score function:
 - Initialize $\mathcal{S}_0 = \dots = \mathcal{S}_w = 0, \hat{\theta}_w = \theta_0;$

Extension to online score matching

- When the post-change training data is not available
- Sequentially estimate the post-change score function:
 - Initialize $\mathcal{S}_0 = \dots = \mathcal{S}_w = 0, \hat{\theta}_w = \theta_0$;
 - For $t = w+1, \dots$ update the score function parameter:

Extension to online score matching

- When the post-change training data is not available
- Sequentially estimate the post-change score function:
 - Initialize $\mathcal{S}_0 = \dots = \mathcal{S}_w = 0, \hat{\theta}_w = \theta_0;$
 - For $t = w+1, \dots$ update the score function parameter:

$$\hat{\theta}_t \leftarrow \hat{\theta}_{t-1} - \eta \nabla_{\hat{\theta}_{t-1}} \mathcal{L}_t(\hat{\theta}_{t-1})$$

“denoising score matching objective” on most recent w samples

Extension to online score matching

- When the post-change training data is not available
- Sequentially estimate the post-change score function:
 - Initialize $\mathcal{S}_0 = \dots = \mathcal{S}_w = 0$, $\hat{\theta}_w = \theta_0$;

- For $t = w+1, \dots$ update the score function parameter:

$$\hat{\theta}_t \leftarrow \hat{\theta}_{t-1} - \eta \nabla_{\hat{\theta}_{t-1}} \mathcal{L}_t(\hat{\theta}_{t-1})$$

“denoising score matching objective” on most recent w samples

- Update the detection statistics \mathcal{S}_t similarly with $\Delta(x_t) = H(x_t; \hat{s}_0) - H(x_t; \hat{s}_{\hat{\theta}_t})$, and compare to detection threshold

Related work

Sequential change detection

- Seminal works [Page, 1954] [Lorden, 1971] [Pollak, 1985] [Moustakides, 1986] [Lai, 1998], etc.
- Recent books [Tartakovsky et al, 2015] [Tartakovsky, 2020], etc.
- Data-driven and non-parametric methods (based on density estimation): [Moustakides and Basioti, 2019] [Liang and Veeravalli, 2023] [Li et al, 2015], etc.

Score-method method for sequential change detection

- Hyvärinen score-based CUSUM (SCUSUM) [Wu et al, 2023] (a series of work)
- Bayesian settings [Banerjee and Tarokh, 2024]
- Robust score-based cumulative sum (RSCUSUM) algorithm [Moushegian et al, 2025]

Theoretical guarantee

Theorem (Zhou et al 2025)

Under the offline score estimation setting, for a given threshold b , and under regularity assumptions, we have with high probability,


$$\text{Detection Delay} \leq \frac{b}{D_F(p_1 \| p_0) - \epsilon_{DSM}} (1 + o(1))$$

Theoretical guarantee

Theorem (Zhou et al 2025)

Under the offline score estimation setting, for a given threshold b , and under regularity assumptions, we have with high probability,

$$\text{Detection Delay} \leq \frac{b}{D_F(p_1 \| p_0) - \epsilon_{DSM}} (1 + o(1))$$



Fisher divergence

Theoretical guarantee

Theorem (Zhou et al 2025)

Under the offline score estimation setting, for a given threshold b , and under regularity assumptions, we have with high probability,

$$\text{Detection Delay} \leq \frac{b}{D_F(p_1 \| p_0) - \epsilon_{DSM}} (1 + o(1))$$


Fisher divergence *Error term caused by DSM*

Theoretical guarantee

Theorem (Zhou et al 2025)

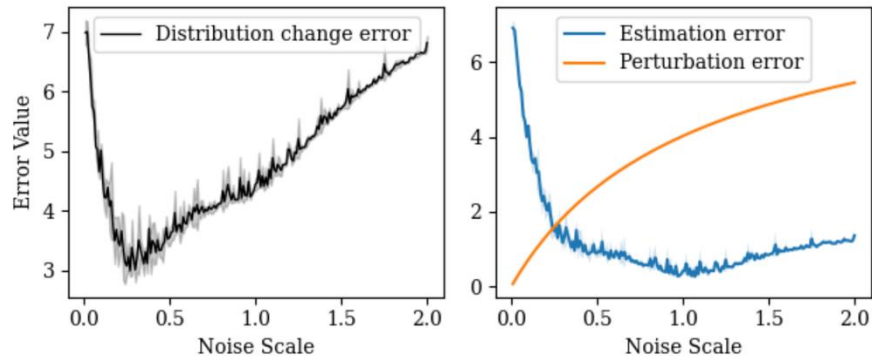
Under the offline score estimation setting, for a given threshold b , and under regularity assumptions, we have with high probability,

$$\text{Detection Delay} \leq \frac{b}{D_F(p_1 \| p_0) - \epsilon_{DSM}} (1 + o(1))$$

\swarrow Fisher divergence \searrow Error term caused by DSM

The error term ϵ_{DSM} contains two components:

- The score estimation error
- The distribution perturbation error



Theoretical guarantee

Theorem (Zhou et al 2025)

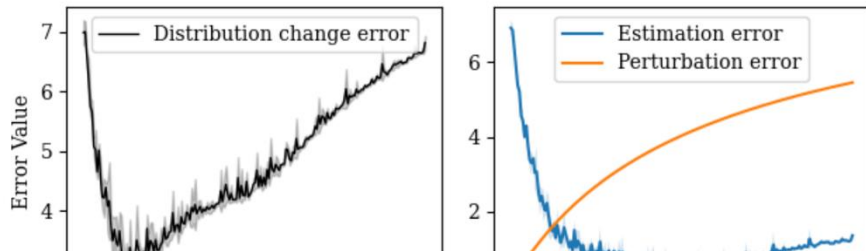
Under the offline score estimation setting, for a given threshold b , and under regularity assumptions, we have with high probability,

$$\text{Detection Delay} \leq \frac{b}{D_F(p_1 \| p_0) - \epsilon_{DSM}} (1 + o(1))$$

Fisher divergence *Error term caused by DSM*

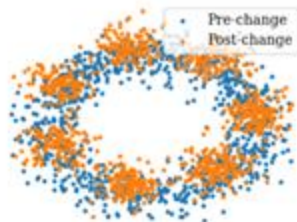
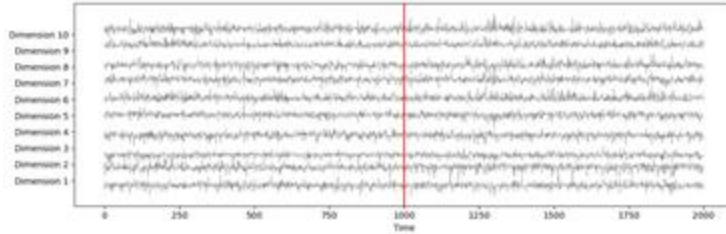
The error term ϵ_{DSM} contains two components:

- The score estimation error
- The distribution perturbation error

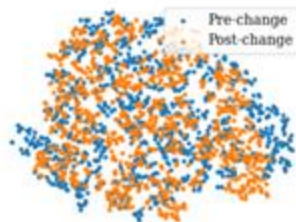


Takeaway: it's better to learn score functions more accurately when using an appropriate level of noise injection.

Simulation Experiments

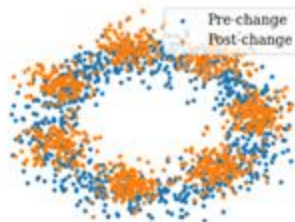
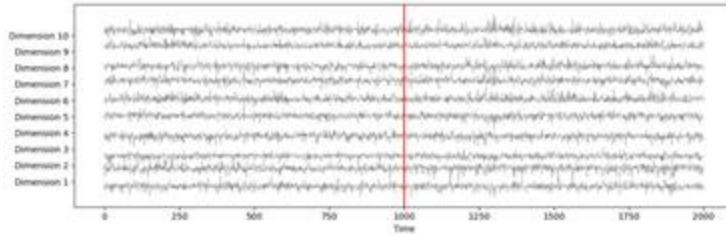


From Gaussian mixture

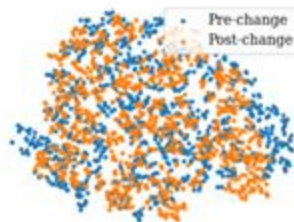


From Neural Network

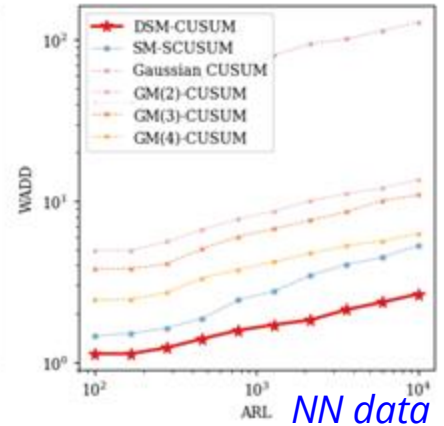
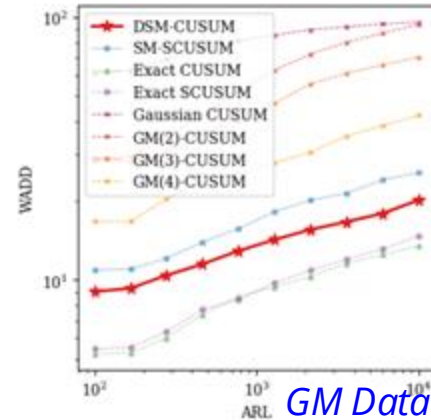
Simulation Experiments



From Gaussian mixture



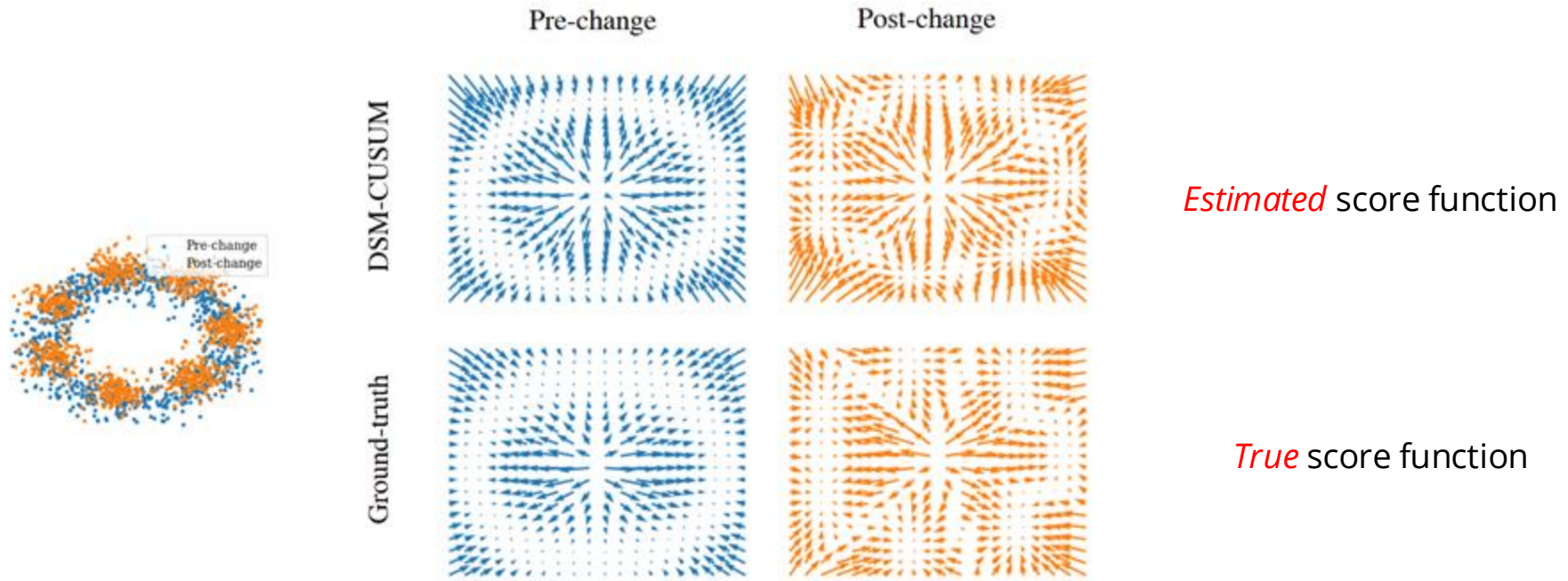
From Neural Network



Detection Delay vs. False Alarm Measures
Lower, the better

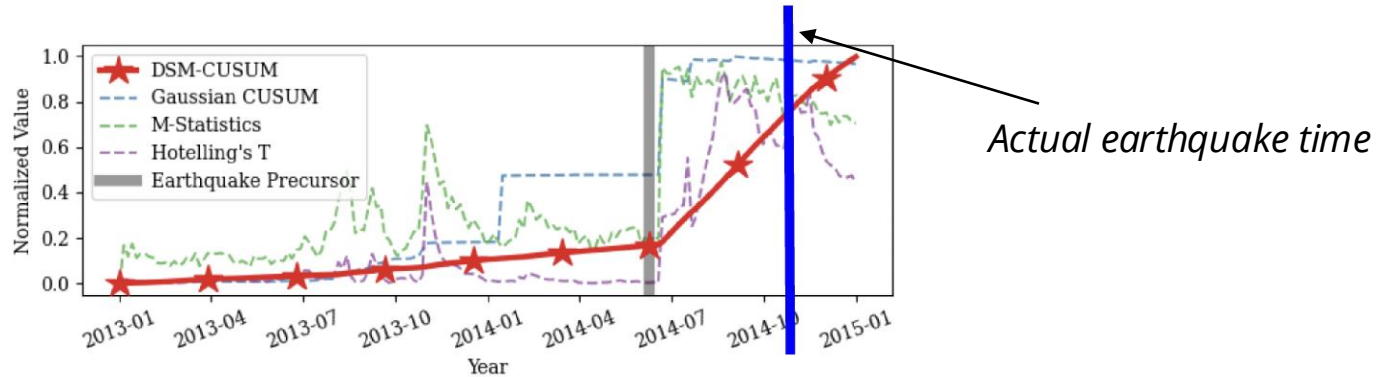
Simulation Experiments: why it works?

Our method works because it can accurately learn the score function.



A real example

- A moderate-size earthquake in 2014
- Four types of signals (water temperatures and water levels) hourly collected at three monitoring stations, over a two-year period.



Conclusion

- Proposed a denoising score matching based algorithm for sequential change detection
- Validated efficiency of our method on both synthetic and real world data
- Future directions: detection for multi-modal data, such as spatio-temporal data, graph data, etc.