

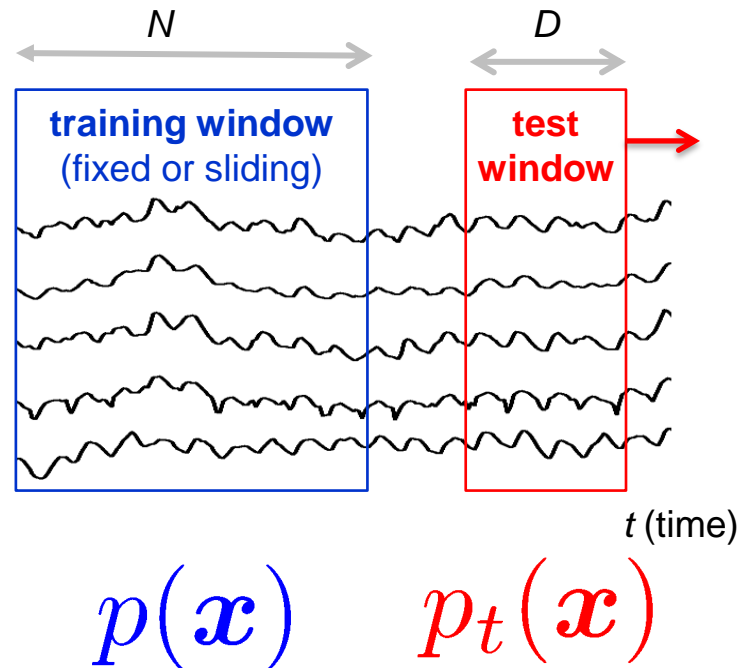
IBM Research

Change Detection Using Directional Statistics

T. Ide (“Ide-san”), D. Phan, J. Kalagnanam
IBM T. J. Watson Research Center

Problem setting: change detection from multi-variate noisy time-series data

- Change = difference between $p(\mathbf{x})$ and $p_t(\mathbf{x})$
 - \mathbf{x} : M -dimensional *i.i.d.* observation
 - $p(\mathbf{x})$: p.d.f. estimated from training window
 - $p_t(\mathbf{x})$: p.d.f. estimated from the test window at time t
- Question 1: What kind of model should we use for the pdf?
- Question 2: How can we quantify the difference between the p.d.f.'s?



We use von Mises-Fisher distribution to model $p(\mathbf{x})$ and $p_t(\mathbf{x})$

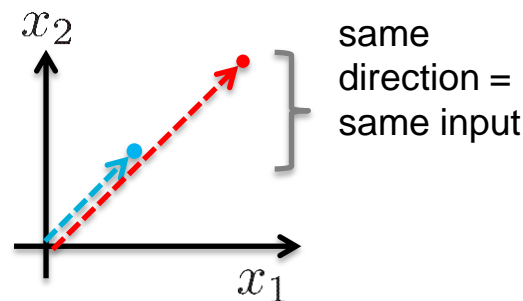
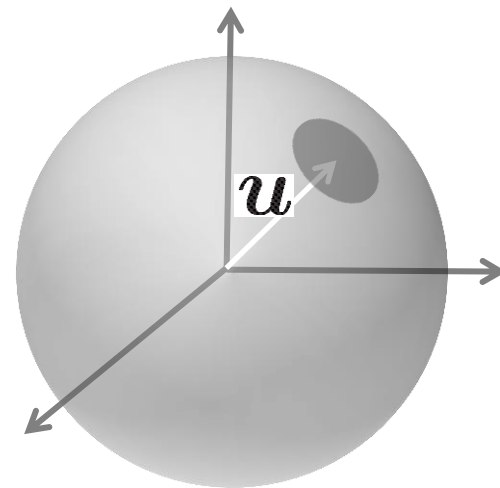
- vMF distribution: “Gaussian for unit vectors”

$$p(\mathbf{z} \mid \mathbf{u}, \kappa) = c_M(\kappa) \exp(\kappa \mathbf{u}^\top \mathbf{z})$$

- \mathbf{z} : random unit vector of $\|\mathbf{z}\| = 1$
- \mathbf{u} : mean direction
- κ : “concentration” (\sim precision in Gaussian)
- M : dimensionality
- We are concerned only with the direction of observation \mathbf{x} :

$$\mathbf{z} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$$

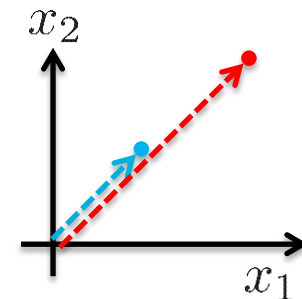
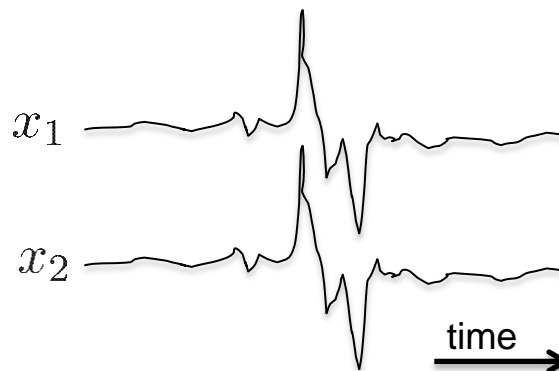
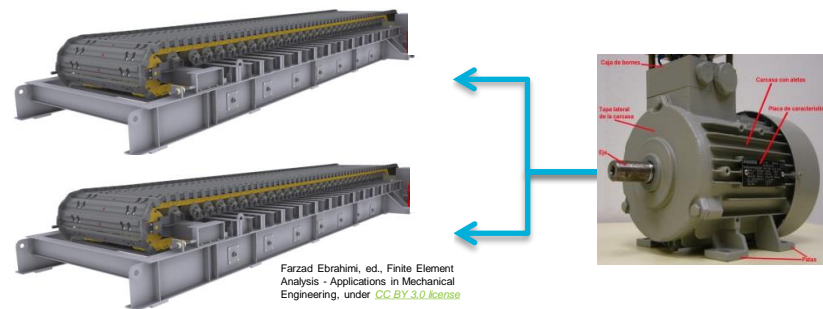
- Normalization is always made
- Do not care about the norm



Why do we enforce normalization?

Rationale from the real world

- Real mechanical systems often incur **multiplicative** noise
 - Example: two belt conveyors operated by the same motor
 - Multiplicative noise equally applied to correlated variables
- Normalization is simple but powerful method for noise reduction



Mean direction \mathbf{u} is learned via weighted maximum likelihood to down-weight contaminated samples

- Weighted likelihood function

$$L(\mathbf{u}, \kappa) = \sum_{n=1}^N w^{(n)} b^{(n)} \left\{ \ln c_M(\kappa) + \kappa \mathbf{u}^\top \mathbf{z}^{(n)} \right\}$$

$\|\mathbf{x}^{(n)}\|_2$ (normalization factor)

sample weight

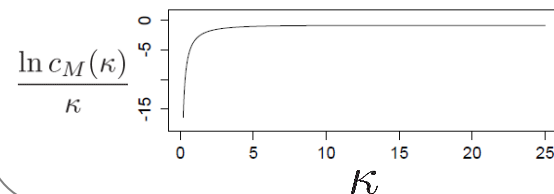
- Regularization over sample weights

$$R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \nu \|\mathbf{w}\|_1$$

- p.d.f. is learned by solving

$$(\mathbf{u}^*, \mathbf{w}^*) = \arg \max_{\mathbf{u}, \mathbf{w}} \{L(\mathbf{u}, \kappa) + \lambda R(\mathbf{w})\}$$

The term related to κ is less important. κ is treated as a given constant.



Multiple patterns (directions) can be obtained by coupling maximum likelihood equations

- Find orthogonal sequence of the mean direction $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$ by coupling the weighted regularized maximum likelihood

$$(\mathbf{u}_1^*, \mathbf{w}_1^*) = \arg \max_{\mathbf{u}_1, \mathbf{w}_1} \{L(\mathbf{u}_1, \kappa) + \lambda R(\mathbf{w}_1)\}$$

$$(\mathbf{u}_2^*, \mathbf{w}_2^*) = \arg \max_{\mathbf{u}_2, \mathbf{w}_2} \{L(\mathbf{u}_2, \kappa) + \lambda R(\mathbf{w}_2)\}$$

$$\vdots$$

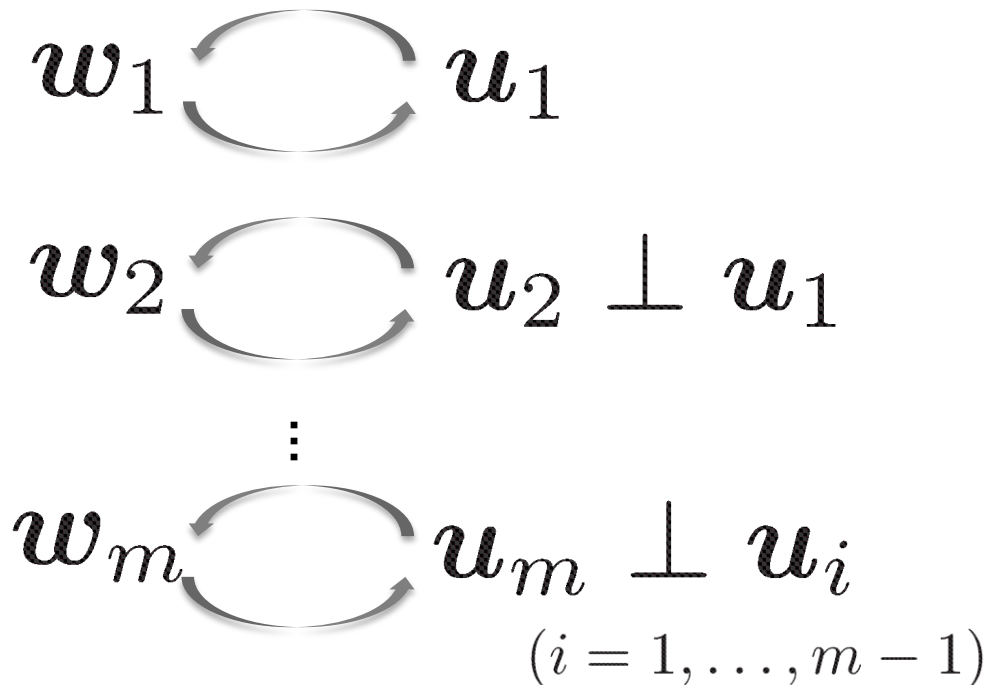
$$(\mathbf{u}_m^*, \mathbf{w}_m^*) = \arg \max_{\mathbf{u}_m, \mathbf{w}_m} \{L(\mathbf{u}_m, \kappa) + \lambda R(\mathbf{w}_m)\}$$

Orthogonality
condition

$$\mathbf{u}_i^\top \mathbf{u}_j = \delta_{i,j}$$

Kronecker
delta

Iterative sequential algorithm for the coupled maximum likelihood



- For each i , w_i and u_i are solved iteratively until convergence
- Analytic solution exists in each step
- Results in very simple fixed point equations

(For reference) Derived fixed-point iteration algorithm

▪ Example: $i=1$

Given \mathbf{w}_1 , solve

$$\max_{\mathbf{u}_1} \{ \kappa \mathbf{u}_1^\top \mathbf{X} \mathbf{w}_1 \} \quad \text{s.t.} \quad \mathbf{u}_1^\top \mathbf{u}_1 = 1$$

Given \mathbf{u}_1 , solve

$$\min_{\mathbf{w}_1} \left\{ \frac{1}{2} \|\mathbf{w}_1 - \frac{\mathbf{q}}{\lambda}\|_2^2 + \nu \|\mathbf{w}_1\|_1 \right\}$$

$$\mathbf{q} \equiv \ln c_M \mathbf{b} + \kappa \mathbf{X}^\top \mathbf{u}_1$$

This Lasso problem is solved analytically

Algorithm 1 RED algorithm.

Input: Initialized \mathbf{w} . Regularization parameters λ, ν . Concentration parameter κ . The number of major directional patterns m .

Output: $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$ and $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_m]$.

for $j = 1, 2, \dots, m$ **do**

while no convergence **do**

$$\mathbf{u}_j \leftarrow \kappa [\mathbf{I}_M - \mathbf{U}_{j-1} \mathbf{U}_{j-1}^\top] \mathbf{X} \mathbf{w}_j \quad (17)$$

$$\mathbf{u}_j \leftarrow \text{sign}(\mathbf{u}_j^\top \mathbf{X} \mathbf{w}_j) \frac{\mathbf{u}_j}{\|\mathbf{u}_j\|_2} \quad (18)$$

$$\mathbf{q}_j \leftarrow \gamma \mathbf{b} + \kappa \mathbf{X}^\top \mathbf{u}_j \quad (19)$$

$$\mathbf{w}_j \leftarrow \text{sign}(\mathbf{q}_j) \odot \max \left\{ \frac{|\mathbf{q}_j|}{\lambda} - \nu \mathbf{1}, \mathbf{0} \right\} \quad (20)$$

end while

end for

Return \mathbf{U} and \mathbf{W} .

Theoretical property: The algorithm is reduced to the “trust-region subproblem” in $\nu \rightarrow 0$

Theorem 2. *When ν tends to 0, the nonconvex problem (5) is reduced to an optimization problem in the form of*

$$\min_{\mathbf{u}} \{ \mathbf{u}^\top \mathbf{Q} \mathbf{u} + \mathbf{c}^\top \mathbf{u} \} \quad \text{s.t.} \quad \mathbf{u}^\top \mathbf{u} = 1, \quad (23)$$

Useful to initialize
the iterative
algorithm

which has a global solution obtained in polynomial time.

Proof. The non-convex optimization problem (23) is known as the trust region subproblem. For polynomial algorithms to the global solution, see [Sorensen, 1997; Tao and An, 1998; Hager, 2001; Toint *et al.*, 2009]. Here we show how the algorithm is reduced to the trust region subproblem.

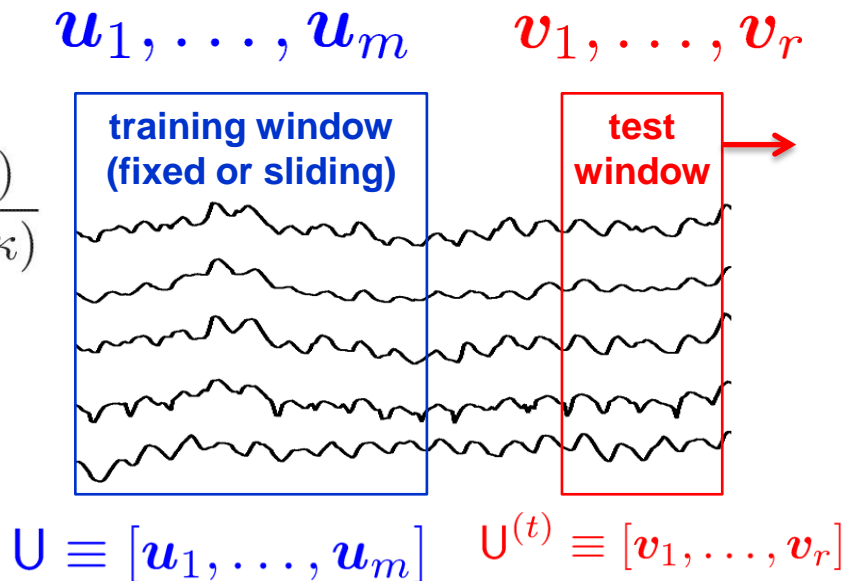
Change score as parameterized Kullback-Leibler divergence

- With extracted directions, define the change score at time t as

$$a^{(t)} = \min_{\mathbf{f}, \mathbf{g}} \int dx \overset{\text{VMF dist.}}{\mathcal{M}(\mathbf{x} | \mathbf{U} \mathbf{f}, \kappa)} \ln \frac{\overset{\text{VMF dist.}}{\mathcal{M}(\mathbf{x} | \mathbf{U} \mathbf{f}, \kappa)}}{\overset{\text{VMF dist.}}{\mathcal{M}(\mathbf{x} | \mathbf{U}^{(t)} \mathbf{g}, \kappa)}}$$

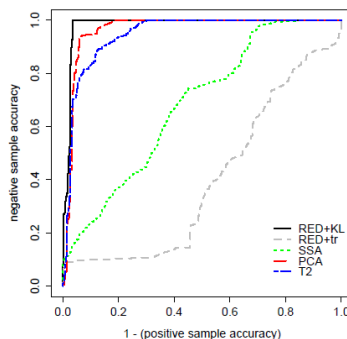
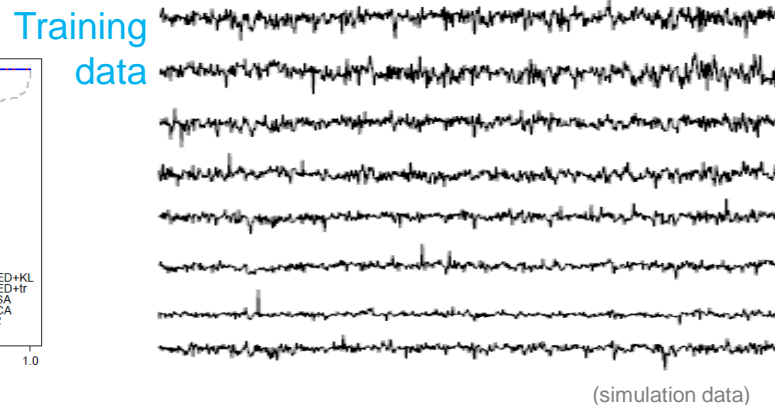
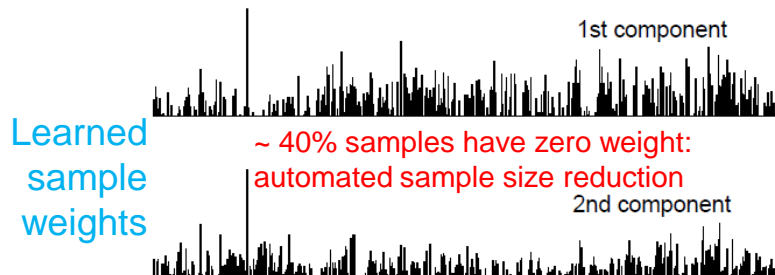
$$\mathbf{f}^\top \mathbf{f} = 1, \mathbf{g}^\top \mathbf{g} = 1$$

- Concisely represented by the top singular value of $\mathbf{U}^\top \mathbf{U}^{(t)}$



Experiment: Failure detection of ore belt conveyors

- vMF formulation successfully suppressed very noisy non-Gaussian noise of multiplicative nature
- ~40% of samples were automatically excluded from the model
- Better than alternatives
 - PCA, Hotelling T²
 - Stationary subspace analysis [Blythe et al., 2012]



Summary – Thank you for your attention!

- Proposed a new change detection algorithm featuring
 - (1) New feature extraction method based on weighted max. likelihood of the vMF distribution
 - (2) Change score based on parameterized KL divergence
- Showed the linkage with the trust region sub-problem