

Christopher Natoli

## 1 Framework of the solution

Suppose  $\{\mathbf{y}(t)\}_{t=1}^T \subset \mathbb{R}^d$  is a  $d$ -dimensional time series and let  $\mathbf{Y}(t) \in \mathbb{R}^{dk}$  be a  $k$ -long subsequence of the data beginning at time  $t$ :

$$\mathbf{Y}(t) := [\mathbf{y}(t)^\top, \mathbf{y}(t+1)^\top, \dots, \mathbf{y}(t+k-1)^\top]^\top,$$

i.e.,  $\mathbf{Y}(t)$  is the concatenation of the  $k$  vectors  $\mathbf{y}(t), \dots, \mathbf{y}(t+k-1)$ . Rather than consider the “feature space” or “state space”  $\mathbb{R}^d$ , i.e., the space in which the datapoints live, we will consider the space  $\mathbb{R}^{dk}$  in which the subsequences  $\mathbf{Y}(t)$  live, since these subsequences can describe some of the temporal structure of the time series that points in  $\mathbb{R}^d$  cannot.

This online algorithm makes use of a sliding window. Divide the sliding window into two intervals, one before the possible changepoint called the reference interval and one after the possible changepoint called the test interval. Denote the beginnings of the reference and test intervals by  $t_{\text{rf}}$  and  $t_{\text{te}}$ , respectively, and denote their lengths by  $n_{\text{rf}}$  and  $n_{\text{te}}$ . If  $t$  is the current time, i.e., the front of the sliding window, then  $t = t_{\text{te}} + n_{\text{te}}$  and  $t_{\text{te}} = t_{\text{rf}} + n_{\text{rf}}$ . Define  $\mathbf{Y}_{\text{rf}}(i) = \mathbf{Y}(t_{\text{rf}} + i - 1)$  and  $\mathbf{Y}_{\text{te}}(i) = \mathbf{Y}(t_{\text{te}} + i - 1)$ . Let  $p_{\text{rf}}$  and  $p_{\text{te}}$  be the probability densities of the reference and test intervals, respectively.

We want to test

$$\begin{aligned} H_0: & p(\mathbf{Y}(i)) = p_{\text{rf}}(\mathbf{Y}(i)) & \text{for } i = t_{\text{rf}}, \dots, t-1 \\ \text{vs} \quad H_1: & p(\mathbf{Y}(i)) = p_{\text{rf}}(\mathbf{Y}(i)) & \text{for } i = t_{\text{rf}}, \dots, t_{\text{te}}-1 \\ & p(\mathbf{Y}(i)) = p_{\text{te}}(\mathbf{Y}(i)) & \text{for } i = t_{\text{te}}, \dots, t-1. \end{aligned}$$

The likelihood ratio and log-likelihood ratio for these hypotheses are therefore

$$\begin{aligned} \Lambda &= \frac{\prod_{i=1}^{n_{\text{rf}}} p_{\text{rf}}(\mathbf{Y}_{\text{rf}}(i)) \prod_{i=1}^{n_{\text{te}}} p_{\text{te}}(\mathbf{Y}_{\text{te}}(i))}{\prod_{i=1}^{n_{\text{rf}}} p_{\text{rf}}(\mathbf{Y}_{\text{rf}}(i)) \prod_{i=1}^{n_{\text{te}}} p_{\text{rf}}(\mathbf{Y}_{\text{te}}(i))} \\ &= \frac{\prod_{i=1}^{n_{\text{te}}} p_{\text{te}}(\mathbf{Y}_{\text{te}}(i))}{\prod_{i=1}^{n_{\text{te}}} p_{\text{rf}}(\mathbf{Y}_{\text{te}}(i))} \\ S := \log \Lambda &= \sum_{i=1}^{n_{\text{te}}} \log \frac{p_{\text{te}}(\mathbf{Y}_{\text{te}}(i))}{p_{\text{rf}}(\mathbf{Y}_{\text{te}}(i))}. \end{aligned}$$

We choose the log-likelihood ratio  $S$  as our score. If the score is greater than a certain threshold, then we say that time  $t_{\text{te}}$  is a changepoint.

## 2 Directly estimating the density ratio

Rather than estimate the density ratio  $w := \frac{p_{\text{te}}}{p_{\text{rf}}}$  by first individually estimating the densities  $p_{\text{te}}$  and  $p_{\text{rf}}$  and then dividing, the ratio itself will be directly estimated by employing the Kullback-Leibler Importance Estimation Procedure (KLIEP) developed by Sugiyama, et. al. [cite!].

Suppose the estimate  $\hat{w}$  lives in a finite-dimensional function space spanned by the basis functions  $\{\phi_\ell\}_{\ell=1}^b$ , where  $\phi_\ell: \mathbb{R}^{dk} \rightarrow \mathbb{R}^+$  for all  $\ell$ . Then  $\hat{w}(\cdot) = \sum_{\ell=1}^b \alpha_\ell \phi_\ell(\cdot)$ , and so the parameters we want to learn are the coordinates  $\{\alpha_\ell\}_{\ell=1}^b$ . We do so by minimizing the Kullback-Leibler divergence of  $\hat{p}_{\text{te}} = p_{\text{rf}} \hat{w}$  from the true distribution  $p_{\text{te}}$ , i.e.,

$$D_{KL}(p_{\text{te}} \parallel \hat{p}_{\text{te}}) = \int p_{\text{te}}(\mathbf{Y}) \log \frac{p_{\text{te}}(\mathbf{Y})}{p_{\text{rf}}(\mathbf{Y}) \hat{w}(\mathbf{Y})} d\mathbf{Y}.$$

The minimization problem is subject to two constraints. First,  $w$  is nonnegative: requiring  $\alpha_\ell \geq 0$  for all  $\ell$  implies  $\hat{w} \geq 0$ . Second,  $p_{\text{te}}$  is a probability density and therefore integrates to 1: the same should hold for  $p_{\text{rf}} \hat{w}$ . Approximation reduces this constrained minimization problem, which is convex, to the following maximization problem:

$$\begin{aligned} & \max_{\{\alpha_\ell\}_{\ell=1}^b} \left\{ \sum_{i=1}^{n_{\text{te}}} \log \sum_{\ell=1}^b \alpha_\ell \phi_\ell(\mathbf{Y}_{\text{te}}(i)) \right\} \\ & \text{subject to } \frac{1}{n_{\text{rf}}} \sum_{i=1}^{n_{\text{rf}}} \sum_{\ell=1}^b \alpha_\ell \phi_\ell(\mathbf{Y}_{\text{rf}}(i)) = 1 \\ & \alpha_\ell \geq 0, \quad \ell = 1, \dots, b. \end{aligned}$$

For this changepoint algorithm, we choose that the basis functions  $\phi_\ell(\cdot)$  are Gaussian kernels  $K_\sigma(\cdot, \mathbf{Y}_{\text{te}}(\ell))$  centered at the test subsequences, where

$$K_\sigma(\mathbf{Y}, \mathbf{Y}') = \exp \left( -\frac{\|\mathbf{Y} - \mathbf{Y}'\|_2^2}{2\sigma^2} \right).$$

It follows that  $b = n_{\text{te}}$ . The rationale for centering kernels on the test points is that it is effective to place kernels where the target function  $w$  is large. The ratio  $w$  is large when  $p_{\text{te}}$  is large and/or when  $p_{\text{rf}}$  is small. The density  $p_{\text{te}}$  of the test distribution is probably large when its argument is pulled from the test distribution, i.e., a test point  $\mathbf{Y}_{\text{te}}(i)$ . Kernels are therefore centered at the test points  $\{\mathbf{Y}_{\text{te}}(i)\}_{i=1}^{n_{\text{te}}}$ .

The kernel parameter  $\sigma$  is selected by cross-validation. The maximization problem is solved by gradient ascent.

## 3 Updating $\{\alpha_\ell\}$ online

Kawahara and Sugiyama propose a way to update  $\{\alpha_\ell\}$  in an online manner rather than repeating the above procedure with every new datapoint.

Let  $\mathcal{Y}_{\text{te}} = \{\mathbf{Y}_{\text{te}}(i)\}_{i=1}^{n_{\text{te}}}$  be the set of test subsequences. Since the Gaussian kernel  $K_\sigma: \mathcal{Y}_{\text{te}} \times \mathcal{Y}_{\text{te}} \rightarrow \mathbb{R}^+$  is a positive definite kernel, we have by the Moore-Aronszajn theorem that there

exists a unique Hilbert space  $\mathcal{H}$  of functions on  $\mathcal{Y}_{\text{te}}$  with  $K_\sigma$  as the reproducing kernel. Moreover, the span of  $\{K_\sigma(\cdot, \mathbf{Y})\}_{\mathbf{Y} \in \mathcal{Y}_{\text{te}}}$  is dense in  $\mathcal{H}$ , so if  $w \in \mathcal{H}$ , then  $\hat{w}$  can approximate  $w$  arbitrarily well. This somewhat justifies our previous assumption that  $\hat{w}$  lives in the span of  $\{K_\sigma(\cdot, \mathbf{Y})\}_{\mathbf{Y} \in \mathcal{Y}_{\text{te}}}$ .

We choose to minimize the regularized risk functional

$$E_i(\hat{w}) = -\log \hat{w}(\mathbf{Y}_{\text{te}}(i)) + \frac{\lambda}{2} \|\hat{w}\|_{\mathcal{H}}^2$$

over the span of  $\{K_\sigma(\cdot, \mathbf{Y})\}_{\mathbf{Y} \in \mathcal{Y}_{\text{te}}}$ , where  $-\log \hat{w}(\mathbf{Y}_{\text{te}}(i))$  is the loss function,  $\frac{\lambda}{2} \|\hat{w}\|_{\mathcal{H}}^2$  is the regularization term, and  $\|\cdot\|_{\mathcal{H}}$  is the norm in  $\mathcal{H}$ . Rather than minimize the average regularized risk  $\frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} E_i(\hat{w})$ , we follow the method of stochastic gradient descent [cite!] and minimize the regularized risk  $E_{n_{\text{te}}+1}$  at only the new training example.

Before differentiating the regularized risk with respect to the function  $\hat{w}$ , we note a few useful derivatives. Let  $e_{\mathbf{Y}}(\hat{w}) := \hat{w}(\mathbf{Y})$  be the evaluation functional at  $\mathbf{Y}$ . Since  $K_\sigma$  is a reproducing kernel, we have by definition that

$$\langle \hat{w}(\cdot), K_\sigma(\cdot, \mathbf{Y}) \rangle = \hat{w}(\mathbf{Y}) = e_{\mathbf{Y}}(\hat{w}).$$

Therefore,

$$\partial_{\hat{w}} e_{\mathbf{Y}}(\hat{w}) = \partial_{\hat{w}} \langle \hat{w}(\cdot), K_\sigma(\cdot, \mathbf{Y}) \rangle = K_\sigma(\cdot, \mathbf{Y})$$

and so

$$\partial_{\hat{w}} (-\log \hat{w}(\mathbf{Y})) = \partial_{\hat{w}} (-\log e_{\mathbf{Y}}(\hat{w})) = -\frac{1}{\log e_{\mathbf{Y}}(\hat{w})} \partial_{\hat{w}} e_{\mathbf{Y}}(\hat{w}) = -\frac{K_\sigma(\cdot, \mathbf{Y})}{\log \hat{w}(\mathbf{Y})}.$$

We also note that

$$\partial_{\hat{w}} \left( \frac{1}{2} \|\hat{w}\|^2 \right) = \hat{w}.$$

Putting these together, we find that

$$\partial_{\hat{w}} E_i(\hat{w}) = -\frac{K_\sigma(\cdot, \mathbf{Y}_{\text{te}}(i))}{\log \hat{w}(\mathbf{Y}_{\text{te}}(i))} + \lambda \hat{w}.$$