

# Incremental Permutation Feature Importance (iPFI): Towards Online Explanations on Data Streams

Fabian Fumagalli<sup>1,\*</sup>, Maximilian Muschalik<sup>2,\*</sup>,  
Eyke Hüllermeier<sup>2</sup>, and Barbara Hammer<sup>1</sup>

✉ [ffumagalli@techfak.uni-bielefeld.de](mailto:ffumagalli@techfak.uni-bielefeld.de)

✉ [maximilian.muschalik@lmu.de](mailto:maximilian.muschalik@lmu.de)

<sup>1</sup> Bielefeld University, <sup>2</sup> LMU Munich, \* equal contribution



# Collaboration



Fabian <sup>1,\*</sup>  
Fumagalli



Maximilian <sup>2,\*</sup>  
Muschalik



Eyke <sup>2</sup>  
Hüllermeier

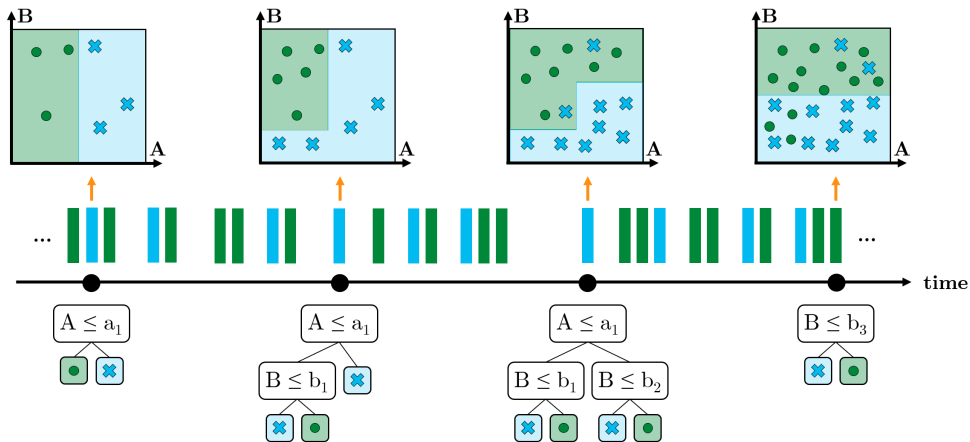


Barbara <sup>1</sup>  
Hammer



\* denotes equal contribution

# Models in Flux: Incremental Learning from Data Streams



**Various applications:** Bifet and Gavaldà (2007), Gama et al. (2014), Davari et al. (2021), etc.

# Examples of Models in Flux



Fraud  
Detection



Sensor  
Networks



Automotive  
Industry

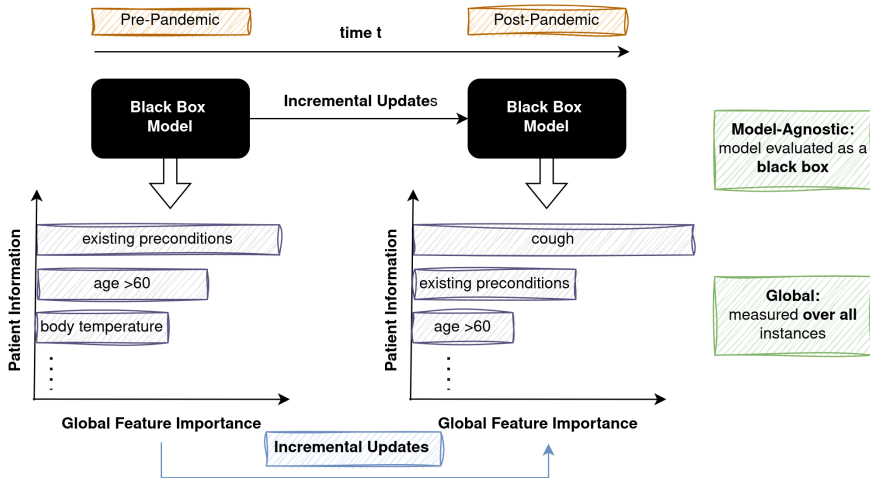


Predictive  
Maintenance

Images generated with Leonardo . ai.

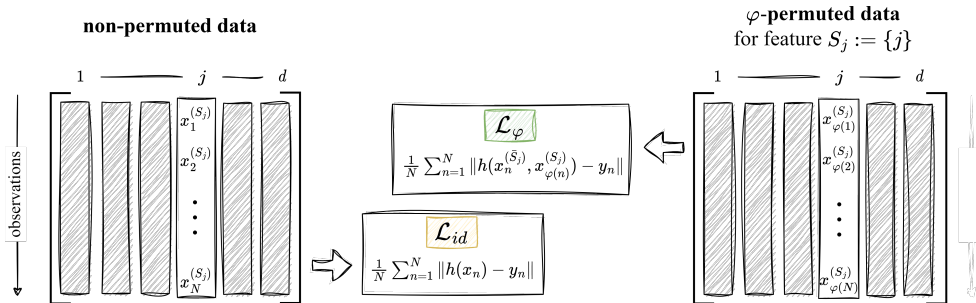
# Model-Agnostic Explanations with Global Feature Importance

## Prediction of Hospital Admission



# Permutation Feature Importance (PFI)

# Permutation Feature Importance (PFI)



## Permutation Feature Importance – (Empirical) PFI

Sample permutations  $\varphi_1, \dots, \varphi_M$  uniformly and compute loss increase  $\hat{\phi}_{\varphi}^{(S_j)} := \mathcal{L}_{\varphi} - \mathcal{L}_{id}$

$$\text{(Empirical) PFI: } \hat{\phi}^{(S_j)} := \frac{N}{N-1} \frac{1}{M} \sum_{m=1}^M \hat{\phi}_{\varphi_m}^{(S_j)}$$

# Theoretical Properties of PFI

## Global Feature Importance (Global FI) of a feature (set) $S_j$

Let  $f_{S_j}(x^{(\bar{S}_j)}, y) := \mathbb{E} \left[ \|h(x^{(\bar{S}_j)}, X^{(S_j)}) - y\| \right]$ , then global FI is defined as

$$\phi^{(S_j)}(h) := \underbrace{\mathbb{E}_{(X, Y)} \left[ f_{S_j}(X^{(\bar{S}_j)}, Y) \right]}_{\text{marginalized risk over } S_j} - \underbrace{\mathbb{E}_{(X, Y)} \left[ \|h(X) - Y\| \right]}_{\text{risk}}$$



# Theoretical Properties of PFI

## Global Feature Importance (Global FI) of a feature (set) $S_j$

Let  $f_{S_j}(x^{(\bar{S}_j)}, y) := \mathbb{E} \left[ \|h(x^{(\bar{S}_j)}, X^{(S_j)}) - y\| \right]$ , then global FI is defined as

$$\phi^{(S_j)}(h) := \underbrace{\mathbb{E}_{(X, Y)} \left[ f_{S_j}(X^{(\bar{S}_j)}, Y) \right]}_{\text{marginalized risk over } S_j} - \underbrace{\mathbb{E}_{(X, Y)} \left[ \|h(X) - Y\| \right]}_{\text{risk}}$$

## Model Reliance Fisher, Rudin, and Dominici (2019)

$$\bar{\phi}^{(S_j)} = \underbrace{\frac{1}{N(N-1)} \sum_{n=1}^N \sum_{m \neq n} \|h(x_n^{(\bar{S}_j)}, x_m^{(S_j)}) - y_n\|}_{=: \hat{e}_{\text{switch}}} - \underbrace{\frac{1}{N} \sum_{n=1}^N \|h(x_n) - y_n\|}_{=: \hat{e}_{\text{orig}}}$$

- is a U-statistic, in particular an **unbiased estimator of global FI**
- is asymptotically Normal with finite sample boundaries

# Theoretical Properties of PFI

## Theorem (PFI and Model Reliance are directly linked)

Model reliance is the expectation of PFI over uniformly drawn permutations:

$$\bar{\phi}(S_j) = \mathbb{E}_{\varphi \sim \text{unif}(\mathfrak{S}_N)}[\hat{\phi}(S_j)] = \frac{N}{N-1} \mathbb{E}_{\varphi \sim \text{unif}(\mathfrak{S}_N)} \left[ \hat{\phi}_{\varphi}(S_j) \right].$$

# Theoretical Properties of PFI

## Theorem (PFI and Model Reliance are directly linked)

Model reliance is the expectation of PFI over uniformly drawn permutations:

$$\bar{\phi}(S_j) = \mathbb{E}_{\varphi \sim \text{unif}(\mathfrak{S}_N)}[\hat{\phi}(S_j)] = \frac{N}{N-1} \mathbb{E}_{\varphi \sim \text{unif}(\mathfrak{S}_N)}[\hat{\phi}_{\varphi}(S_j)].$$

### PFI $\hat{\phi}(S_j)$ variant of Breiman (2001)

- Easy to compute in  $\mathcal{O}(N)$
- Difficult to analyze theoretically due to dependence on permutations
- **Used for computation**

### Expected PFI $\bar{\phi}(S_j) = \mathbb{E}_{\varphi}[\hat{\phi}(S_j)]$

- Hard to compute in  $\mathcal{O}(N^2)$
- U-statistic with theoretical guarantees Fisher, Rudin, and Dominici (2019)
- **Used for theoretical analysis**

# Incremental Permutation Feature Importance (iPFI)

Towards Online Explanations on Data Streams

# Incremental PFI for Online Learning

## Online Learning on Data Streams

- Unlimited data stream  $(x_0, y_0), \dots, (x_t, y_t), \dots$
- Incrementally updated model:  $h_{t+1} \leftarrow \text{incrementalUpdate}(h_t, x_t, y_t)$

### Static Permutation Tests

$$\phi_{\varphi}^{(S_j)} = \mathcal{L}_{\varphi} - L_{\text{id}} = \frac{1}{N} \sum_{n=1}^N \|h(x_n^{(\bar{S}_j)}, x_{\varphi(n)}^{(S_j)}) - y_n\| - \|h(x_n) - y_n\|$$

At time  $t$  with  $(x_t, y_t)$  and model  $h_t$

### Stochastic Sampling Strategy

$$\varphi_t : \Omega \rightarrow \{0, \dots, t-1\}$$

### Replacement with previous Observations

$$\|h_t(x_t^{(\bar{S}_j)}, x_{\varphi_t}^{(S_j)}) - y_t\| - \|h_t(x_t) - y_t\|$$

# Incremental PFI for Online Learning

## Online Learning on Data Streams

- Unlimited data stream  $(x_0, y_0), \dots, (x_t, y_t), \dots$
- Incrementally updated model:  $h_{t+1} \leftarrow \text{incrementalUpdate}(h_t, x_t, y_t)$

### Calculation at time t

$$\hat{\lambda}_t^{(S_j)}(x_t, x_{\varphi_t}, y_t) := \|h_t(x_t^{(S_j)}, x_{\varphi_t}^{(S_j)}) - y_t\| - \|h_t(x_t) - y_t\|$$

### Initial Computation

$$\hat{\phi}_{t_0-1}^{(S_j)} := 0 \text{ for } t \geq t_0 > 0$$

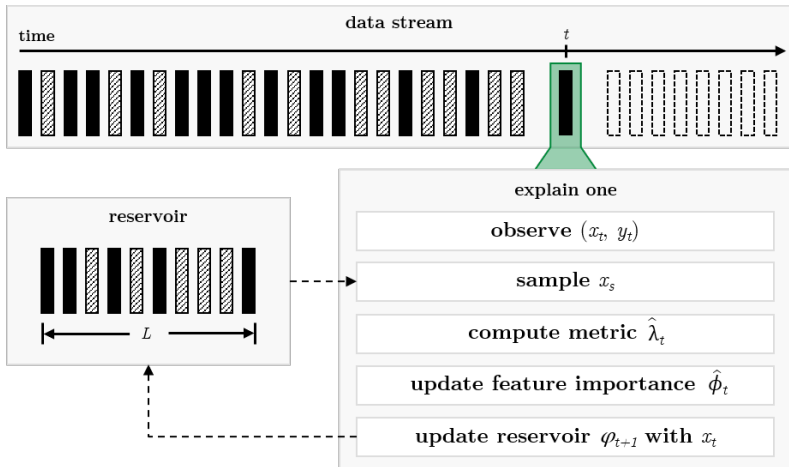
### Incremental Update of iPFI

$$\hat{\phi}_t^{(S_j)} := (1 - \alpha) \cdot \hat{\phi}_{t-1}^{(S_j)} + \alpha \cdot \hat{\lambda}_t^{(S_j)}(x_t, x_{\varphi_t}, y_t)$$

### Smoothing Parameter

$$\alpha \in (0, 1)$$

# iPFI – Algorithm Illustration



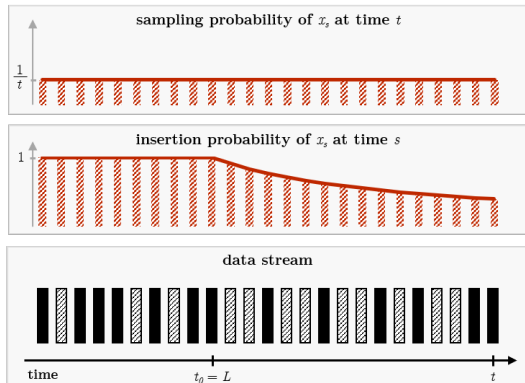
Similar Computational Complexity at time  $t$

one static PFI score

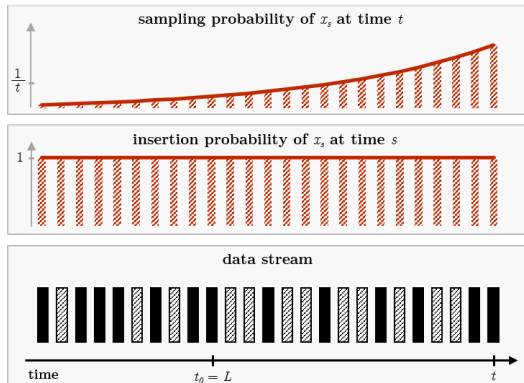
all anytime iPFI scores

# iPFI – Incremental Reservoir Sampling

uniform sampling



geometric sampling





# iPFI – Theoretical Guarantees in Static Environments

## Expected iPFI

With a (stochastic) sampling strategy  $\varphi := (\varphi_s)_{s=t_0, \dots, t}$ , we define

$$\text{Expected iPFI: } \bar{\phi}_t^{(S_j)} := \mathbb{E}_\varphi[\hat{\phi}_t^{(S_j)}].$$

## Theorem (Static Model and $(X_t, Y_t) \sim \mathbb{P}_{(X, Y)}$ )

If  $h \equiv h_t$  and  $\mathbb{V}[\|h(X_s^{(\bar{S}_j)}, X_r^{(S_j)}) - Y_s\| - \|h(X_s) - Y_s\|] < \infty$ , then

$$\phi^{(S_j)}(h) - \mathbb{E}[\bar{\phi}_t^{(S_j)}] = (1 - \alpha)^{t-t_0+1} \phi^{(S_j)}(h) \quad (\text{bias})$$

$$\mathbb{V} \left[ \lim_{t \rightarrow \infty} \bar{\phi}_t^{(S_j)} \right] = \mathcal{O}(-\alpha \log(\alpha)) \quad (\text{uniform sampling})$$

$$\mathbb{V} \left[ \lim_{t \rightarrow \infty} \bar{\phi}_t^{(S_j)} \right] = \mathcal{O}(\alpha) + \mathcal{O}(1/L) \quad (\text{geometric sampling})$$

## Controlling Change in Dynamic Environments

We define a **measure of change** between two timesteps  $t_0 \leq s \leq t$  as

$$f_S^\Delta(x^{(\bar{S}_j)}, h_s, h_t) := \mathbb{E}_{\tilde{X} \sim \mathbb{P}_S} [\|h_t(x^{(\bar{S}_j)}, \tilde{X}) - h_s(x^{(\bar{S}_j)}, \tilde{X})\|]$$
$$\Delta_S(h_s, h_t) := \mathbb{E}_X [f_S^\Delta(X, h_s, h_t)] \text{ and } \Delta(h_s, h_t) := \Delta_\emptyset(h_s, h_t).$$

## Theorem (Changing Model and $(X_t, Y_t) \sim \mathbb{P}_{(X, Y)}$ )

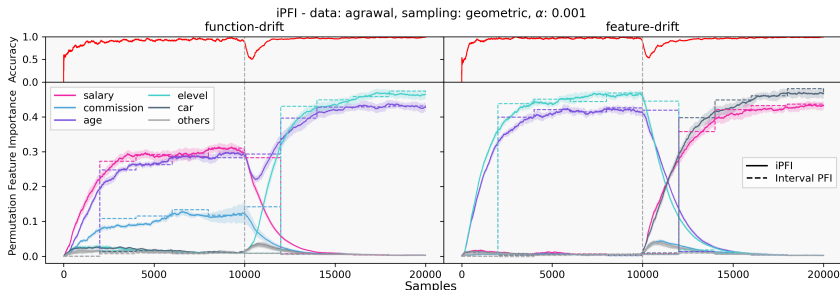
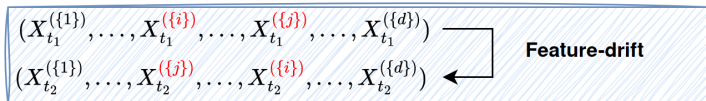
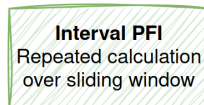
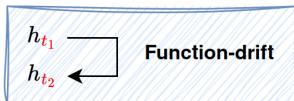
If  $\Delta(h_s, h_t) \leq \delta$  and  $\Delta_S(h_s, h_t) \leq \delta_S$  for  $t_0 \leq s \leq t$  and finite covariances, then

$$|\mathbb{E}[\bar{\phi}_t^{(S_j)}] - \phi^{(S_j)}(h_t)| \leq \delta_S + \delta + \mathcal{O}((1 - \alpha)^t) \quad (\text{bias})$$

$$\mathbb{V} \left[ \lim_{t \rightarrow \infty} \bar{\phi}_t^{(S_j)} \right] = \mathcal{O}(-\alpha \log(\alpha)) \quad (\text{uniform sampling})$$

$$\mathbb{V} \left[ \lim_{t \rightarrow \infty} \bar{\phi}_t^{(S_j)} \right] = \mathcal{O}(\alpha) + \mathcal{O}(1/L) \quad (\text{geometric sampling})$$

# iPFI vs. Interval PFI for Concept Drifts



# Conclusion & Outlook

## Permutation Feature Importance

- (Empirical) PFI as a variant of permutation test (Breiman 2001)
- Expected PFI as model reliance (Fisher, Rudin, and Dominici 2019)
- **Expected PFI is the expectation of PFI over uniformly sampled permutations**

## Incremental Permutation Feature Importance (iPFI)

- We introduce online explanations for online learning on data streams
- We propose an efficient incremental computation of PFI
- **iPFI efficiently reveals model and distribution changes over time**
- **iPFI is supported by theoretical guarantees in controlled environments**

# The Road Ahead and Open Source Implementation

## Towards Explaining Change

- iPFI is a **model-agnostic** XAI method to compute **global FI** for models **in flux**.
- Online XAI approaches include **iSAGE** (to-day at **16:30-18:30** here in room **Fucine**) and **iPDP** (xAI'23).

## Workshop Friday Afternoon Slot

- Time: **14:00-18:00**
- Room: **PoliTo Room 10i**
- Title: *Explainable Artificial Intelligence: From Static to Dynamic*



docs passing pypi v0.1.3 status alpha license MIT






### Installation

```
pip install ixai
```

### Quickstart

```
>>> for (n, (x, y)) in enumerate(stream, start=1)
...     accuracy.update(y, model.predict_one(x)) # inference
...     incremental_pfi.explain_one(x, y) # explaining
...     model.learn_one(x, y) # learning
```

# References

-  Bifet, Albert and Ricard Gavaldà (2007). “Learning from Time-Changing Data with Adaptive Windowing”. In: *Proceedings of the Seventh SIAM International Conference on Data Mining (SIAM 2007)*, pp. 443–448. DOI: 10.1137/1.9781611972771.42.
-  Breiman, Leo (2001). “Random Forests”. In: *Machine Learning* 45.1, pp. 5–32.
-  Davari, Narjes et al. (2021). “Predictive Maintenance Based on Anomaly Detection Using Deep Learning for Air Production Unit in the Railway Industry”. In: *8th IEEE International Conference on Data Science and Advanced Analytics (DSAA 2021)*. IEEE, pp. 1–10. DOI: 10.1109/DSAA53316.2021.9564181.
-  Fisher, Aaron, Cynthia Rudin, and Francesca Dominici (2019). “All Models are Wrong, but Many are Useful: Learning a Variable’s Importance by Studying an Entire Class of Prediction Models Simultaneously”. In: *Journal of Machine Learning Research* 20.177, pp. 1–81.
-  Gama, João et al. (2014). “A Survey on Concept Drift Adaptation”. In: *ACM Comput. Surv.* 46.4, 44:1–44:37. DOI: 10.1145/2523813.

## General Explanation Algorithm

---

**Algorithm 6** Incremental explanation procedure

---

**Require:** stream  $\{x_t, y_t\}_{t=1}^{\infty}$ , model  $f(\cdot)$ , loss function  $\mathcal{L}(\cdot)$

```
1: for all  $(x_t, y_t) \in$  stream do  
2:    $\hat{y}_t \leftarrow f_t(x_t)$   
3:    $\hat{\phi}_t \leftarrow \text{explain\_one}(x_t, y_t)$   
4:    $f_{t+1} \leftarrow \text{learn\_one}(\mathcal{L}(\hat{y}_t, y_t))$   
5: end for
```

---

- Similarly to the **prequential** training, we explain models prequentially.
- Data points are used first for explanations (model has not seen the observation, line 3) and then the model is allowed to use it for training (line 4).

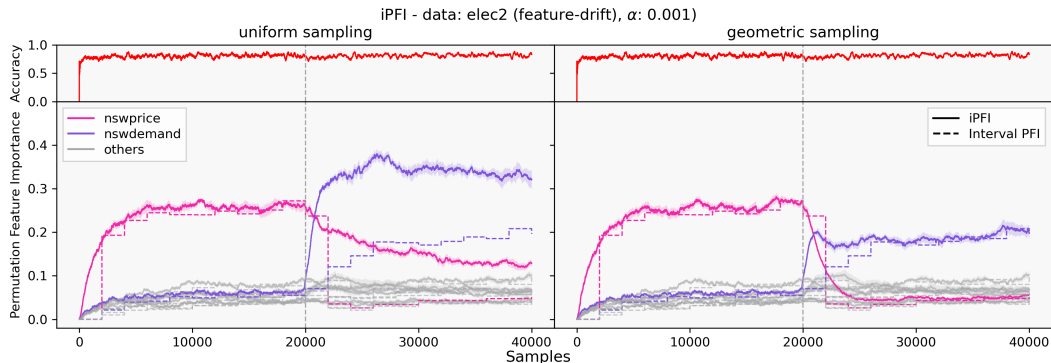


# Computational Complexity

data	<i>stagger</i>	<i>elec2</i>	<i>agrawal</i>	<i>adult</i>	<i>bank</i>	<i>insects</i>	<i>ozone</i>
feature count	3	8	9	14	16	33	72
explanation time	0.734 (.017)	1.210 (.039)	1.411 (.020)	1.976 (.118)	2.386 (.048)	5.070 (.078)	7.717 (.182)
inference time	0.959 (.001)	0.989 (.002)	0.987 (.001)	0.991 (.002)	0.991 (.001)	0.990 (.021)	0.998 (.000)

Table 1: Summary of the additional time complexity of iPFI. The additional *explanation time* is given relatively to the case where the models are trained without explaining. The *inference time* denotes the portion of the explanation time in which the models are queried. All values for each dataset are derived from ten independent runs. The run time of iPFI scales *linearly* with  $0.104 \cdot |D|$  over the number of features ( $R^2 = 0.966$ ).

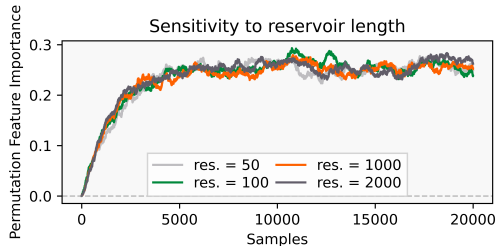
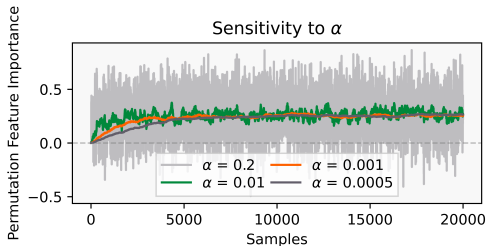
# Uniform vs. Geometric Sampling



## Geometric Sampling for Feature-Drift

If feature distributions change, then geometric sampling should be preferred.

# Parameters



## Choice of Smoothing Parameter $\alpha$

The choice depends on the application. We recommend

$\alpha = 0.001$  (conservative) and  $\alpha = 0.01$  (reactive).