# iPDP: On Partial Dependence Plots in Dynamic Modeling Scenarios

**Maximilian Muschalik**[1,2,*], Fabian Fumagalli[3,*], Rohit Jagtani[1],
Barbara Hammer[3], and Eyke Hüllermeier[1,2]
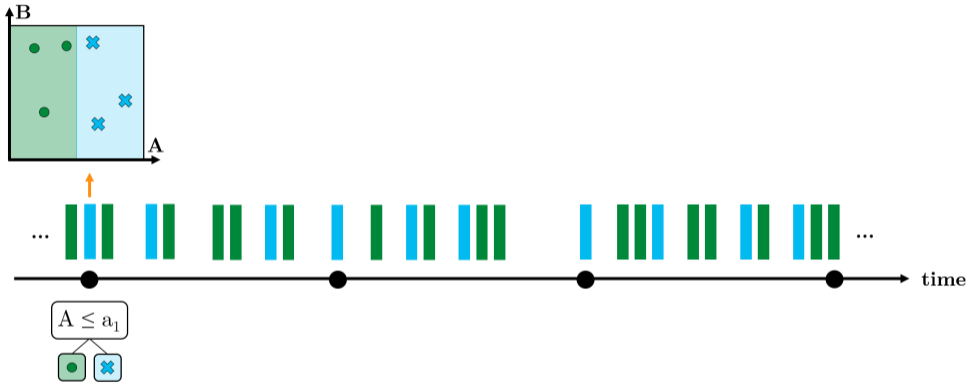
✉ maximilian.muschalik@lmu.de
✉ ffumagalli@techfak.uni-bielefeld.de
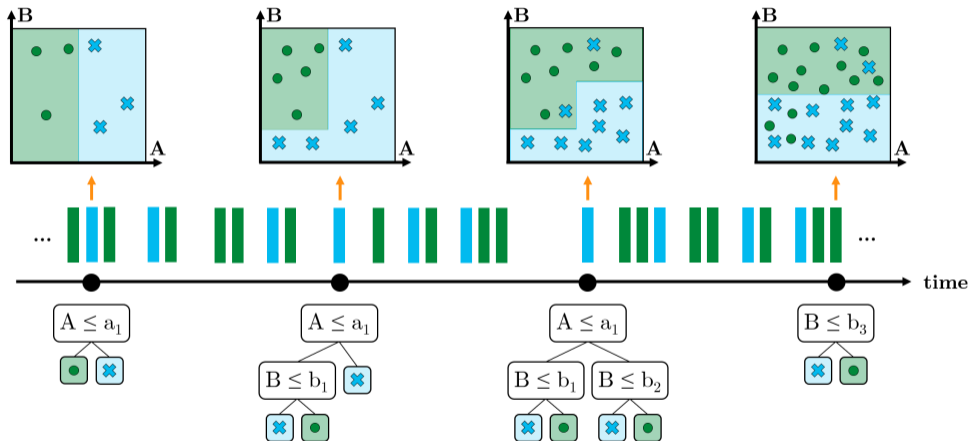[1] LMU Munich, [2] MCML Munich, [3] Bielefeld University
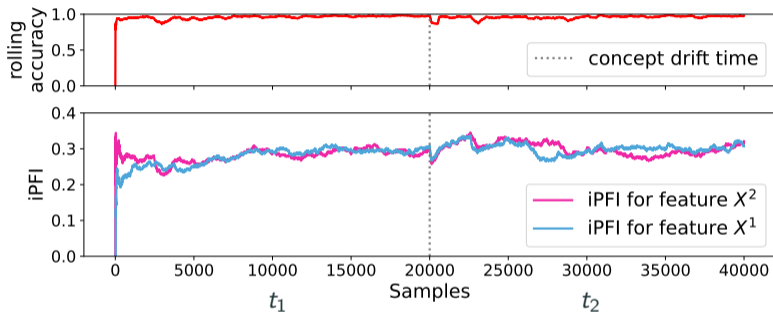* denotes equal contribution

# Online Models are Learning Incrementally from Data Streams

**Various applications**: Bifet and Gavaldà 2007, Gama et al. 2014, Davari et al. 2021, etc.
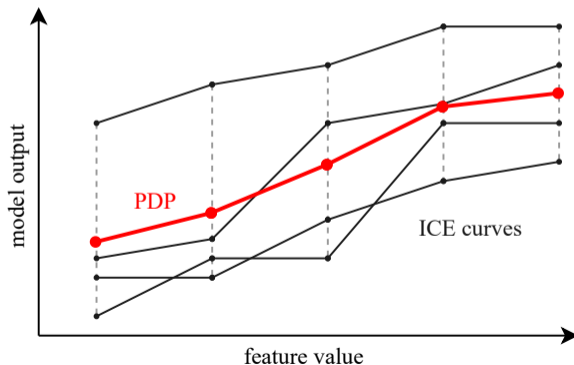
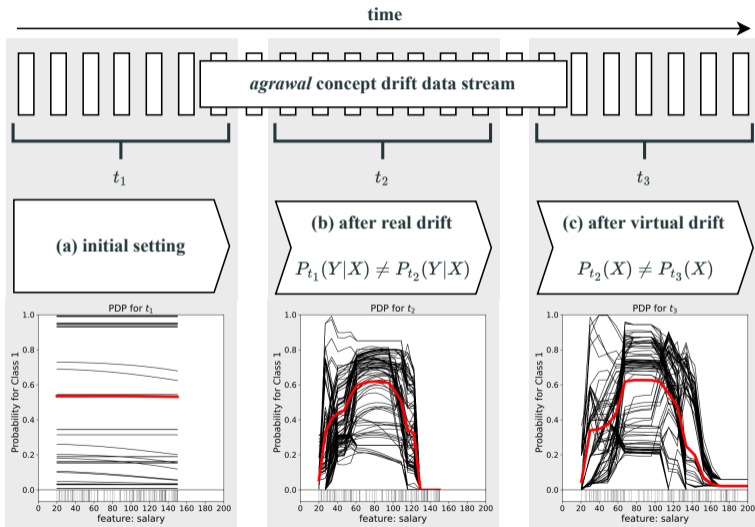**hidden concept drift**

$$P_{t_1}(Y|X) \neq P_{t_2}(Y|X)$$

# Partial Dependence Plots (PDPs) Explain Feature Effects

**Definition of PDP (Friedman 2001)**

$$f_S^{\text{PD}}(\mathbf{x}^S) = \mathbb{E}_{X^{\bar{S}}}\left[ f(\mathbf{x}^S, X^{\bar{S}}) \right] \qquad \text{in practice:} \quad \hat{f}_S^{\text{PD}}(\mathbf{x}^S) = \frac{1}{n}\sum_{i=1}^{n} f(\mathbf{x}^S, \mathbf{x}_i^{\bar{S}})$$
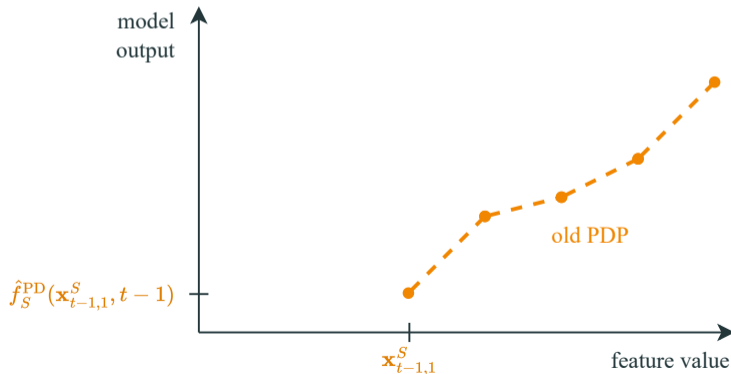
# PDP on Virtual and Real Concept Drift



time

*agrawal* concept drift data stream

$t_1$      $t_2$      $t_3$

(a) initial setting

(b) after real drift
$P_{t_1}(Y|X) \neq P_{t_2}(Y|X)$

(c) after virtual drift
$P_{t_2}(X) \neq P_{t_3}(X)$

## Definition of iPDP

**iPDP:**

$$\overbrace{\hat{f}_S^{\mathrm{PD}}(\mathbf{x}_{t-1,k}^S, t - 1)}^{\text{old PDP}}$$



old PDP

$\hat{f}_S^{\mathrm{PD}}(\mathbf{x}_{t-1,1}^S, t - 1)$

model output

feature value

$\mathbf{x}_{t-1,1}^S$

# Incremental PDP (iPDP) for Moving Models and Data

## Definition of iPDP

**iPDP:**
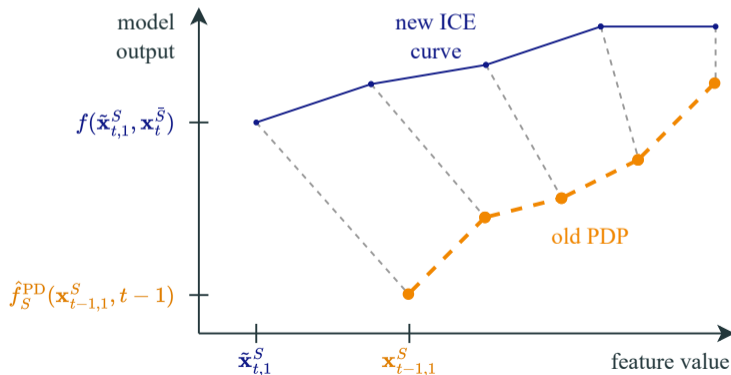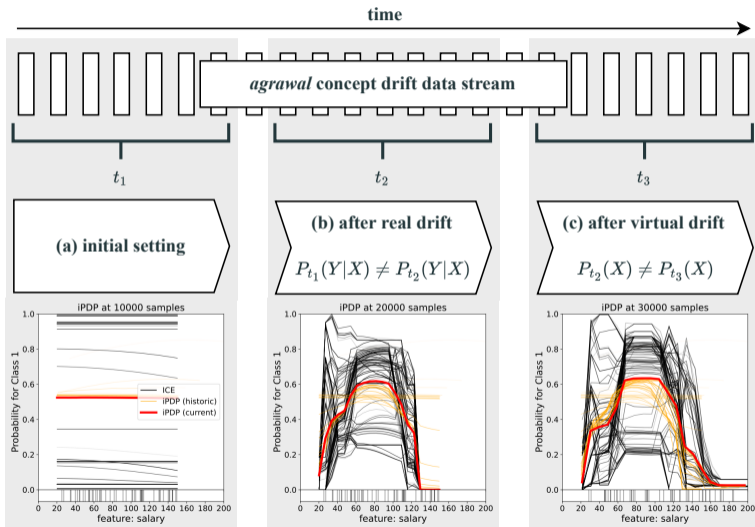$$\overbrace{\hat{f}^{\text{PD}}_S(\mathbf{x}^S_{t-1,k}, t-1)}^{\text{old PDP}} \qquad \overbrace{f_t(\tilde{\mathbf{x}}^S_{t,k}, \mathbf{x}^{\bar{S}}_t)}^{\text{new ICE}}$$

## Definition of iPDP

**iPDP:** $\overbrace{\hat{f}_S^{\text{PD}}(\mathbf{x}_{t,k}^S, t)}^{\text{new PDP}} := (1-\alpha) \cdot \overbrace{\hat{f}_S^{\text{PD}}(\mathbf{x}_{t-1,k}^S, t-1)}^{\text{old PDP}} + \alpha \cdot \overbrace{f_t(\tilde{\mathbf{x}}_{t,k}^S, \bar{\mathbf{x}}_t^{\bar{S}})}^{\text{new ICE}}$

**time**

*agrawal* concept drift data stream

$t_1$      $t_2$      $t_3$

**(a) initial setting**

**(b) after real drift**

$P_{t_1}(Y|X) \neq P_{t_2}(Y|X)$

**(c) after virtual drift**

$P_{t_2}(X) \neq P_{t_3}(X)$

iPDP at 10000 samples

iPDP at 20000 samples

iPDP at 30000 samples

ICE
iPDP (historic)
iPDP (current)

## Theoretical Guarantees

### Theorem (Reactiveness)

*iPDP reacts to real drift and favors recent PD values, as*

$$\mathbb{E}[\hat{f}_S^{PD}(\mathbf{x}_{t,k}^S, t)] = \alpha \sum_{i=1}^{t} (1-\alpha)^{t-i} \underbrace{\mathbb{E}_{X_i^{\bar{S}}} \left[ f_i(\tilde{\mathbf{x}}_{i,k}^S, X_i^{\bar{S}}) \right]}_{PD \text{ function at time } i}, \text{ for } k = 1, \ldots, m.$$

## Theoretical Guarantees

### Theorem (Reactiveness)

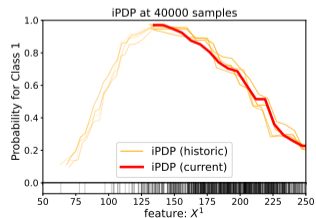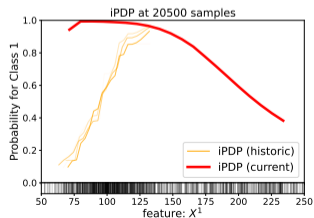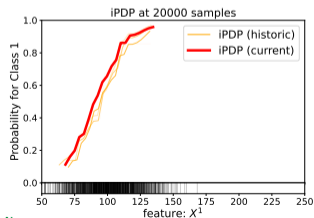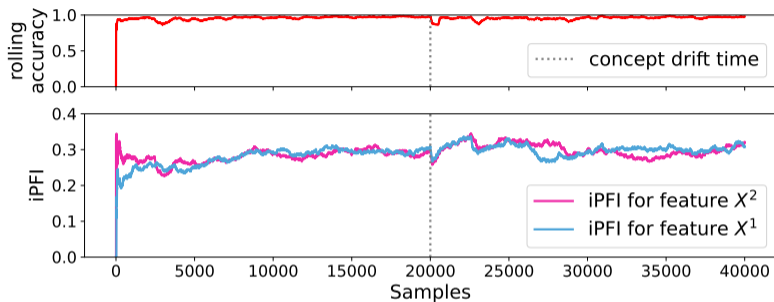iPDP reacts to real drift and favors recent PD values, as

$$\mathbb{E}[\hat{f}_S^{PD}(\mathbf{x}_{t,k}^S, t)] = \alpha \sum_{i=1}^{t} (1-\alpha)^{t-i} \underbrace{\mathbb{E}_{X_i^{\bar{S}}}\left[ f_i(\tilde{\mathbf{x}}_{i,k}^S, X_i^{\bar{S}}) \right]}_{PD \text{ function at time } i}, \text{ for } k = 1, \ldots, m.$$

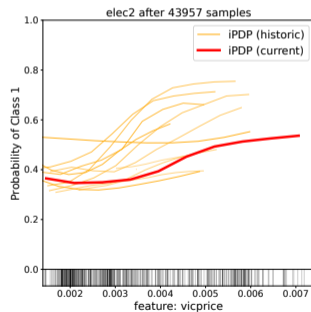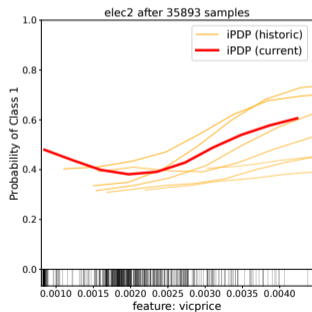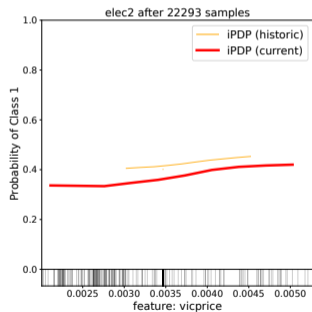### Theorem (Batch PDP Approximation in Static Settings)

Let observations $(x_0, y_0), \ldots, (x_t, y_t)$ be iid from $\mathbb{P}(X, Y)$ and $f \equiv f_t$ be a static model. If $f$ is locally linear in the range of temporary model evaluation points $\{\tilde{\mathbf{x}}_{i,k}^S\}_{i=1}^{t}$ for $k = 1, \ldots, m$, then

$$\mathbb{E}\left[ \hat{f}_S^{PD}(\mathbf{x}_{t,k}^S, t) \right] = f_S^{PD}\left( \mathbf{x}_{t,k}^S \right) \text{ and } \mathbb{E}\left[ \frac{\hat{f}_S^{PD}(\mathbf{x}_{t,k}^S, t)}{1 - (1-\alpha)^t} \right] = f_S^{PD}\left( \frac{\mathbf{x}_{t,k}^S}{1 - (1-\alpha)^t} \right).$$

elec2 after 22293 samples

elec2 after 35893 samples
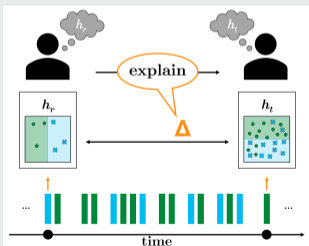
elec2 after 43957 samples

# The Road Ahead and Open Source Implementation

**Towards Explaining Change.**

- iPDP is a **model-agnostic** XAI method to capture feature effects of models **in flux**.
- iSAGE and iPFI can be used to compute global feature importance incrementally.



## *i*XAI

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.

📄 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.

📄 Friedman, Jerome H. (2001). "Greedy Function Approximation: A Gradient Boosting Machine". In: *The Annals of Statistics* 29.5, pp. 1189–1232. ISSN: 00905364. URL: http://www.jstor.org/stable/2699986.

# References

📄 Fumagalli, Fabian et al. (2022). "Incremental Permutation Feature Importance (iPFI): Towards Online Explanations on Data Streams". In: *CoRR* abs/2209.01939. DOI: 10.48550/arXiv.2209.01939. arXiv: 2209.01939. URL: https://doi.org/10.48550/arXiv.2209.01939.

📄 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.

📄 Muschalik, Maximilian et al. (2023). "iSAGE: An Incremental Version of SAGE for Online Explanation on Data Streams". In: *CoRR* abs/2303.01181. arXiv: 2303.01181. URL: https://doi.org/10.48550/arXiv.2303.01181.

# Efficient Access to Feature Distribution over Time

## Maximum Value Storage
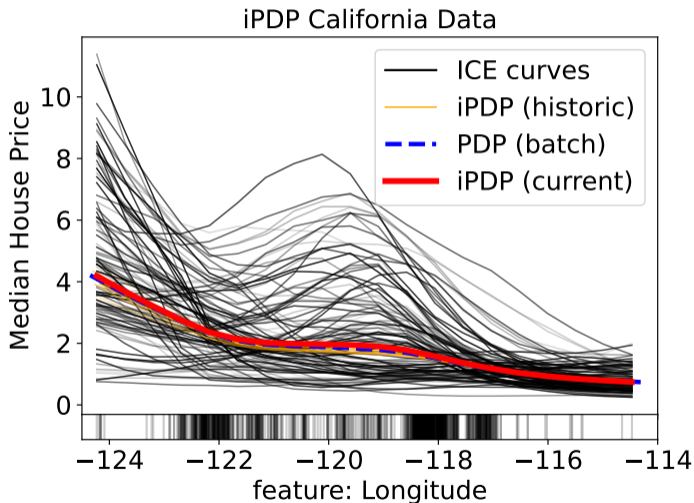


**(a) expected case with iid data**



**(b) worst case scenario**

## Removal Stratgies

- **Interventional** removal (or categorical features) can be stored in **Geometric Reservoirs** (Fumagalli et al. 2022)

- **Observational** removal can be stored in **Incremental Subgroups** (Muschalik et al. 2023)

iPDP California Data

time

**agrawal concept drift data stream**

$t_1$

$X_{t_1}^{\text{salary}} \sim \text{Unif}([20, 150])$

$Y_{t_1} = f_1(X_{t_1}^{\text{age}}, X_{t_1}^{\text{education}})$

**(a) initial setting**

$t_2$

$X_{t_2}^{\text{salary}} \sim \text{Unif}([20, 150])$

$Y_{t_2} = f_2(X_{t_2}^{\text{age}}, X_{t_2}^{\text{salary}})$

**(b) after real drift**
$P_{t_1}(Y|X) \neq P_{t_2}(Y|X)$

$t_3$

$X_{t_3}^{\text{salary}} \sim \text{Unif}([20, 200])$

$Y_{t_3} = f_2(X_{t_3}^{\text{age}}, X_{t_3}^{\text{salary}})$

**(c) after virtual drift**
$P_{t_2}(X) \neq P_{t_3}(X)$

PDP for $t_1$

PDP for $t_2$

PDP for $t_3$

**Algorithm 1** iPDP Explanation Procedure

**Require:** stream $\{\mathbf{x}_t, y_t\}_{t=1}^{\infty}$, model $f_t(.)$, feature set of interest $S$, smoothing parameter $0 < \alpha \leq 1$, number of grid points $m$, and storage object $R_t$

1: initialize $\hat{f}_S^{\text{PD}}(\mathbf{x}_{0,k}^S, 1) \leftarrow 0$
2: **for all** $(\mathbf{x}_t, y_t) \in$ stream **do**
3: $\quad \{\tilde{\mathbf{x}}_{t,k}^S\}_{k=1}^m \leftarrow \text{GetGridPoints}(R_t, m)$ {e.g., equidistant points, quantiles, etc.}
4: $\quad$ **for** $k = 1, \ldots, m$ **do**
5: $\quad\quad \mathbf{x}_{t,k}^S \leftarrow (1 - \alpha) \cdot \mathbf{x}_{t-1,k}^S + \alpha \cdot \tilde{\mathbf{x}}_{t,k}^S$ {update grid point}
6: $\quad\quad \hat{y}_k \leftarrow f_t\left(\tilde{\mathbf{x}}_{t,k}^S, \mathbf{x}_t^{\bar{S}}\right)$ {evaluate on model evaluation point}
7: $\quad\quad \hat{f}_S^{\text{PD}}(\mathbf{x}_{t,k}^S, t) \leftarrow (1-\alpha) \cdot \hat{f}_S^{\text{PD}}(\mathbf{x}_{t-1,k}^S, t-1) + \alpha \cdot \hat{y}_k$ {update point-wise estimates}
8: $\quad$ **end for**
9: $\quad R_t \leftarrow \text{UpdateStorage}(R_{t-1}, x_t^S)$ {add $x_t^S$ to the storage object}
10: $\quad$ **Output:** $\frac{\hat{f}_S^{\text{PD}}(\mathbf{x}_{t,k}^S, t)}{1-(1-\alpha)^t}, \frac{\mathbf{x}_{t,k}^S}{1-(1-\alpha)^t}$ {debiasing of estimates and grid points}
11: **end for**