# Low-Complexity Nonparametric Bayesian Online Prediction with Universal Guarantees

#### Frédéric Cazals Alix Lhéritier



### Online prediction with side information

Features  $\in \mathbb{R}^d$   $z^n \equiv z_1, \dots, z_n$  predict:  $P(l_n | l^{n-1}, z^n)$ Labels  $\in \mathcal{L}$  (discrete)  $l^n \equiv l_1, \dots, l_n$ 

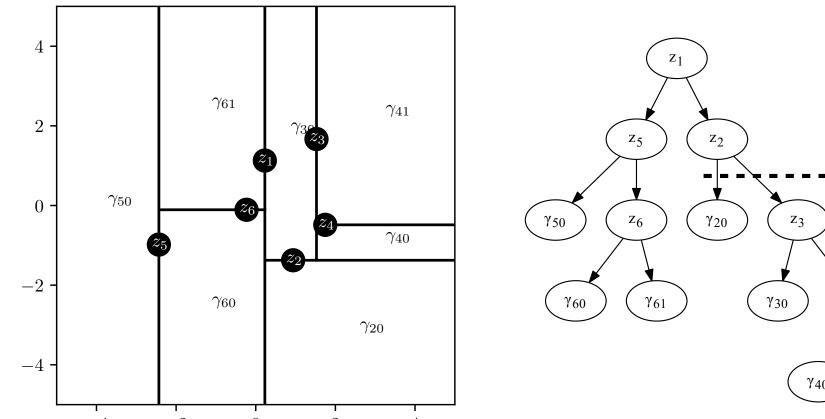
• Goal:  $-\frac{1}{\pi} \log P(l^n|z^n)$  asymptotically optimal.

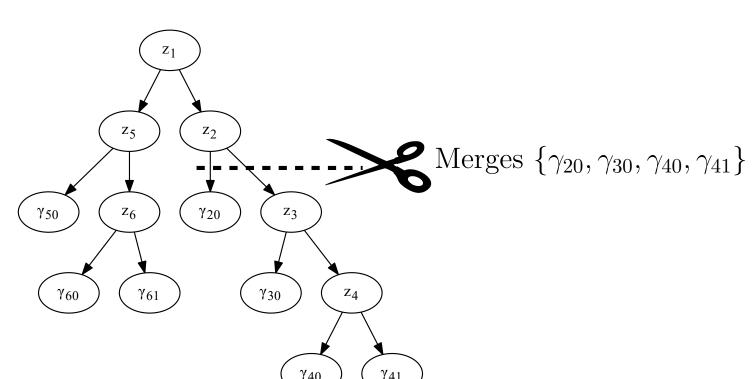
amadeus

- **Probabilistic setting:**  $(z_i, l_i)$  are i.i.d. realizations of RV (Z, L).  $\Rightarrow$  Optimum: H(L|Z).
- Previous work: scale-hyperparameter dependence and high complexity  $\Rightarrow$  k-nn: needs k(n),  $O(n^2)$ . Gaussian Processes: need kernel width,  $O(n^4)$ .

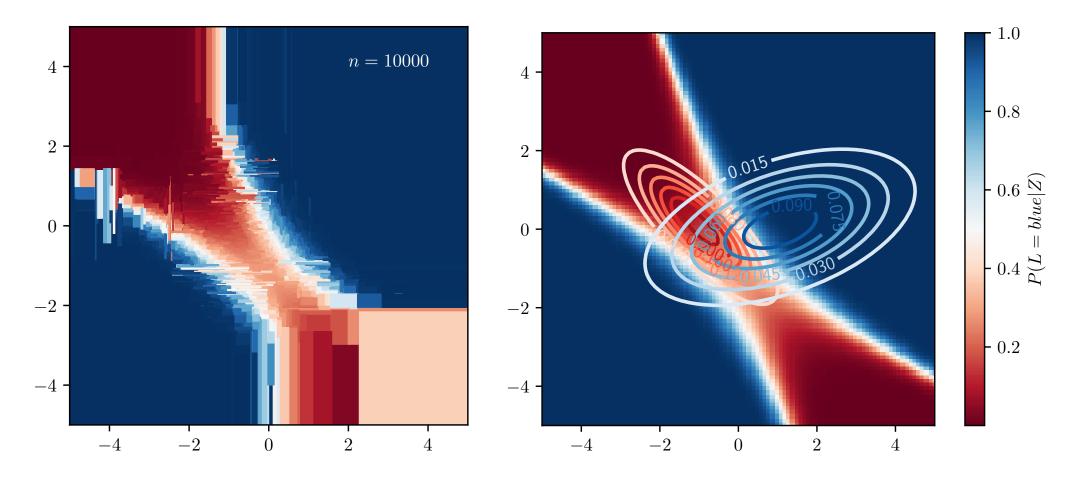
#### Online data-driven discretization of the feature space

- Full-fledged k-d trees: given  $z^n$ , feature space is recursively partitioned by hyperplanes passing by the projection of  $z_i$  onto a perpendicular random axis.
- Depth is  $O(\log n)$  in probability with respect to  $\mathbb{P}_{Z^n} \Rightarrow \mathbf{low\text{-}complexity}$ .
- $\mathcal{P}_n$  is the set of all the partitions obtained by pruning the partitioning tree.





• Partitioning rule  $\pi_n$  s.t.  $H(L|\pi_n(Z|Z^n)) \xrightarrow{\text{a.s.}} H(L|Z) \Rightarrow \mathbf{nonparametric}$ .



# Bayesian label prediction with context tree models

• Context Tree Weighting: a Bayesian mixture on piecewise multinomial distributions over a nested  $\mathcal{P}$ , with a natural encoding prior (P(split) = 1/2)

$$P_{\mathrm{CTW}}(l^n) \equiv \sum_{C \in \mathcal{P}} w_{\mathrm{nat}}(C) \prod_{\gamma \in C} P_{\mathrm{Jeffreys}}(\gamma(l^n))$$

where  $\gamma(l^n)$  is the subsequence falling in  $\gamma$ .

- Context Tree Switching is a mixture over sequences of these distributions.
- For any partition  $C \in \mathcal{P}$ , the context C(L) is the node  $\gamma_i$  where L falls.

$$-\lim_{n\to\infty}\frac{1}{n}\log P_{\mathrm{CT}}(L^n)\leq H\left(L|C(L)\right) \text{ a.s. }\Rightarrow \mathbf{universal}.$$

#### The kd-switch distribution

• For a given node  $\gamma$  split into  $\gamma_1$  and  $\gamma_2$  by a k-d tree on  $z^n$ ,

$$\begin{split} P_{\text{kds}}^{\gamma}(l^{n}|z^{n}) &\equiv \sum_{i^{n} \in \{a,b\}^{n}} w_{\gamma}(i^{n}) \prod_{k=1}^{n} \left[ \mathbb{1}_{\{i_{k}=a\}} \phi_{a}(l_{k}|l^{k-1}) + \mathbb{1}_{\{i_{k}=b\}} \phi_{b}^{\gamma}(l_{k}|l^{k-1},z^{k}) \right] \\ \phi_{a}(l_{k}|l^{k-1}) &\equiv P_{\text{kt}} \left( l_{k}|l^{k-1} \right) \\ \left( P_{\text{kt}} \left( l_{k}|l^{k-1} \right) \text{ if } k < \tau_{k}(\gamma) \end{split}$$

$$\phi_b^{\gamma}(l_k|l^{k-1}, z^k) \equiv \begin{cases} P_{\text{kt}}(l_k|l^{k-1}) & \text{if } k < \tau_k(\gamma) \\ P_{\text{kds}}^{\gamma_j}(\gamma_j(l^k)|\gamma_j(z^k)) \\ \hline P_{\text{kds}}^{\gamma_j}(\gamma_j(l^k)^{-1}|\gamma_j(z^k)^{-1}) & \text{with } j : z_k \in \gamma_j, \text{ otherwise} \end{cases}$$

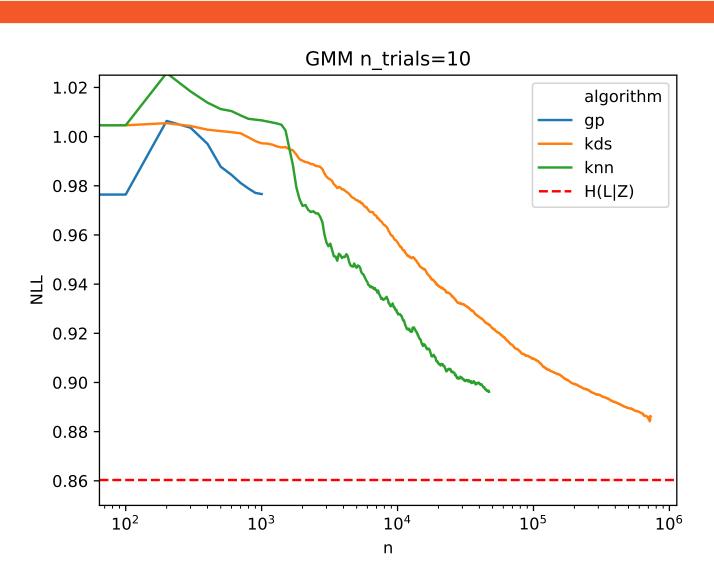
where  $w_{\gamma}(\cdot)$  is a prior over  $\{a,b\}^n$  and  $\tau_n(\gamma)$  is the splitting index in  $\gamma(z^n)$ .

#### Main result: pointwise universality

Thm: The **kd-switch** distribution is pointwise universal, i.e.

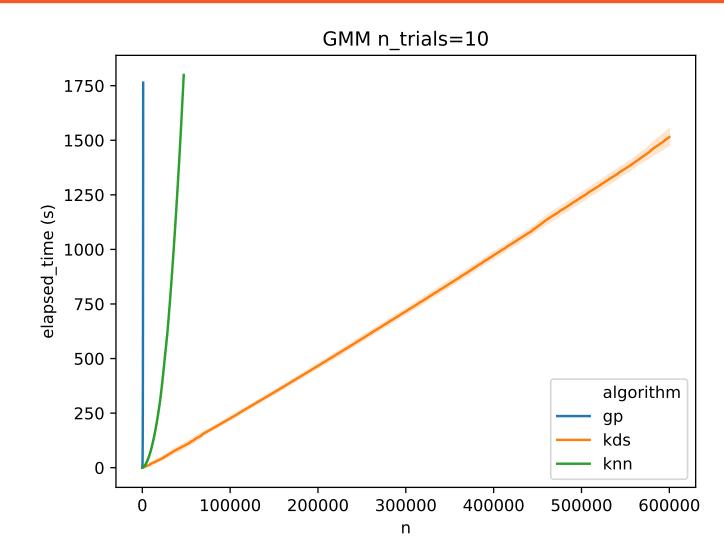
$$-\lim_{n\to\infty} \frac{1}{n} \log P_{\mathrm{kds}}(L^n|Z^n) \le H(L|Z) \text{ a.s.}$$

for any probability measure  $\mathbb{P}$  generating the samples such that  $\mathbb{P}_{Z|L}$  are absolutely continuous with respect to the Lebesgue measure.



# Low-complexity online algorithm

The **kd-switch** distribution can be computed online in  $O(n \log n)$  time.



### Application: Sequential two-sample testing [5]

• The two-sample problem:

$$\begin{cases} \mathtt{H}_0: \mathbb{P}_{Z|L=0} = \mathbb{P}_{Z|L=1} \text{ a.e.} \\ \mathtt{H}_1: \neg \mathtt{H}_0 \end{cases}.$$

• Thm: Given an online predictor P, a test rejecting  $H_0$  at **any index** n s.t.

$$\frac{\mathbb{P}(l^n)}{P(l^n|z^n)} \le c$$

has a Type I error probability less or equal than  $\alpha$ , i.e.

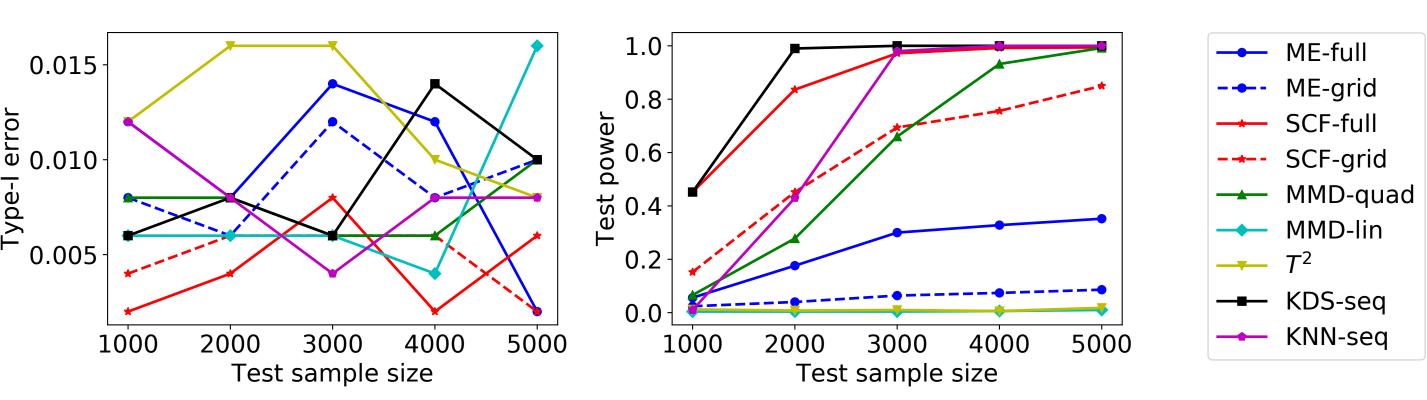
$$\mathbb{P}_{\mathtt{H}_0}\left(\exists n: rac{\mathbb{P}\left(L^n
ight)}{P(L^n|Z^n)} \leq lpha
ight) \leq lpha.$$

• Thm: A **pointwise universal** P yields a **consistent** two-sample test:

$$\mathbb{P}\left(\text{Type II error}\right) \to 0.$$

#### Sequential vs training-set-optimized two-sample tests

- KNN-seq: k-nn based sequential test [5].
- SCF-full and ME-full [1, 3]:
- Kernel based consistent two-sample tests.
- Train/test paradigm to optimize revealing locations and kernel width.
- Valid if *n* is fixed in advance.



#### References

- [1] K. P. Chwialkowski, A. Ramdas, D. Sejdinovic, and A. Gretton.
- Fast two-sample testing with analytic representations of probability measures.
- In Advances in Neural Information Processing Systems, pages 1972–1980, 2015.
- [2] L. Devroye, L. Györfi, and G. Lugosi.
- A probabilistic theory of pattern recognition, volume 31. Springer Verlag, 1996.
- [3] W. Jitkrittum, Z. Szabó, K. P. Chwialkowski, and A. Gretton. Interpretable distribution features with maximum testing power.
- In Advances in Neural Information Processing Systems, pages 181–189, 2016.
- [4] R. Krichevsky and V. Trofimov. The performance of universal encoding.
- Information Theory, IEEE Transactions on, 27(2):199–207, 1981.
- [5] A. Lhéritier and F. Cazals.
- A sequential non-parametric multivariate two-sample test. IEEE Transactions on Information Theory, 64(5):3361–3370, 2018.
- [6] J. Veness, K. S. Ng, M. Hutter, and M. Bowling.
- Context tree switching.
- In Data Compression Conference (DCC), 2012, pages 327–336. IEEE, 2012.
- [7] F. M. J. Willems, Y. M. Shtarkov, and T. J. Tjalkens.
  - The context-tree weighting method: Basic properties.
  - Information Theory, IEEE Transactions on, 41(3):653–664, 1995.