

SCALABLE AND ADAPTIVE PREDICTION BANDS WITH KERNEL SUM-OF-SQUARES

ETICS 2025

Louis Allain

October 9, 2025

Thesis advisor : Sébastien Da Veiga

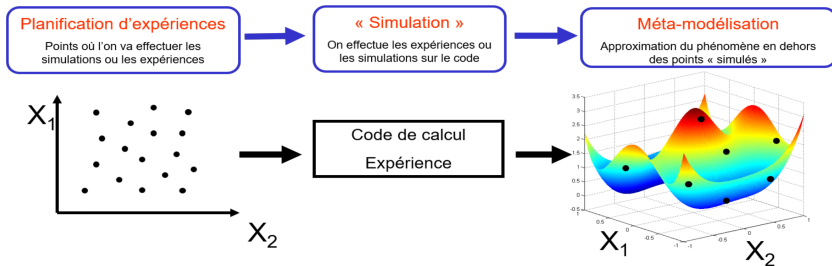
Supervisor : Brian Staber



Table of Contents

1. Introduction
2. Learning a score function
3. Experiments
4. Conclusion
5. References and appendix

- Widespread use of supervised learning for computer experiments, where expensive simulation outputs are approximated with a ML model from a DoE dataset



- Strategy adopted across multiple industries, using various ML models:
 - Linear/logistic regression
 - Random forests
 - Gaussian processes
 - Neural networks...
- Critical applications require **confidence intervals around predictions**, with guaranteed coverage:
 - Denote $\hat{C}(X)$ a confidence interval for a prediction at X , estimated from training data
 - The guarantee of **marginal coverage** at level α writes

$$\mathbb{P}(Y_{N+1} \in \hat{C}(X_{N+1})) \geq 1 - \alpha$$

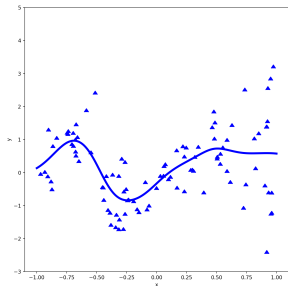
for the true unknown value of the output Y_{N+1} at an unobserved point X_{N+1}

- Limitations of traditional approaches:
 - Prediction bands are model-specific, with significant variation between models
 - Guarantees only valid as $n \rightarrow +\infty$ or under strong assumptions that cannot be verified
 - **No coverage guarantee** for practical applications
- A recent promising candidate: conformal prediction

- Conformal Prediction (CP): a rigorous method to construct prediction intervals with the following properties:
 - ✓ Coverage guarantees
 - ✓ Finite sample
 - ✓ Distribution free
 - ✓ Model agnostic
- Several variants:
 - Full CP
 - Split CP
 - Resampling strategies, e.g. jackknife+, CV+

Let us illustrate split CP, which is based on two independent datasets \mathcal{D}_n (pre-training set) and \mathcal{D}_m (calibration set)

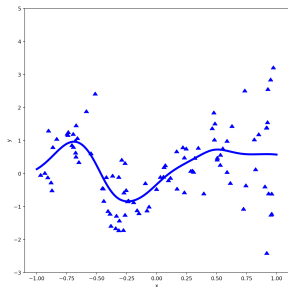
- The prediction model \hat{m} is trained on \mathcal{D}_n



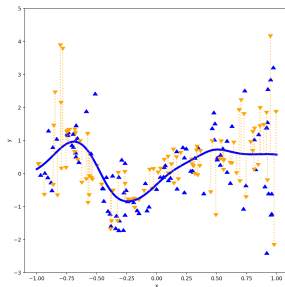
(a) Learn \hat{m} on \mathcal{D}_n

Conformal Prediction

- The prediction model \hat{m} is trained on \mathcal{D}_n
- \mathcal{D}_m is used to evaluate some prediction quality of \hat{m} , here for example the absolute residuals
- The quantile \hat{q}_α of these quality measures is computed



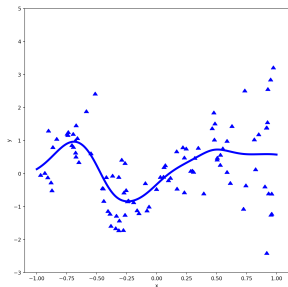
(a) Learn \hat{m} on \mathcal{D}_n



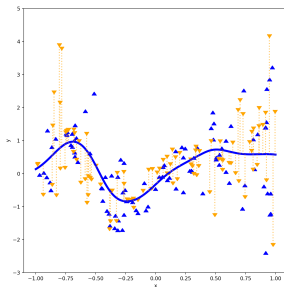
(b) Compute $|Y_i - \hat{m}(X_i)|$ on \mathcal{D}_m

Conformal Prediction

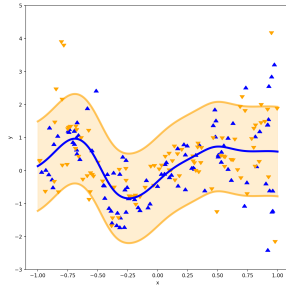
- The prediction model \hat{m} is trained on \mathcal{D}_n
- \mathcal{D}_m is used to evaluate some prediction quality of \hat{m} , here for example the absolute residuals
- The quantile \hat{q}_α of these quality measures is computed
- The prediction interval $\hat{C}(X) = [\hat{m}(X) \pm \hat{q}_\alpha]$ satisfies all the desired properties under the assumption of data exchangeability



(a) Learn \hat{m} on \mathcal{D}_n



(b) Compute $|Y_i - \hat{m}(X_i)|$ on \mathcal{D}_m

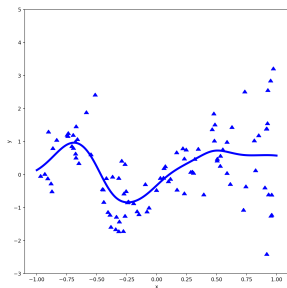


(c) $\hat{C}(X) = [\hat{m}(X) \pm \hat{q}_\alpha]$

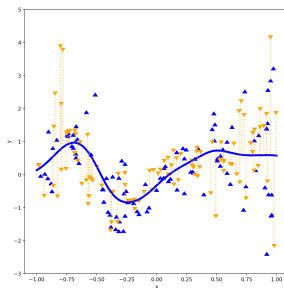
- Such evaluation of the prediction quality is performed by a **score function** s

Absolute errors

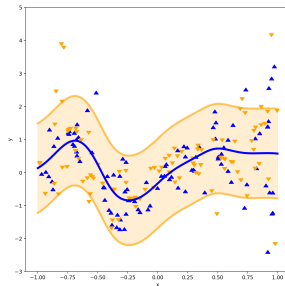
$$s(X_i, Y_i) = |Y_i - \hat{m}(X_i)|$$



(a) Learn \hat{m} on \mathcal{D}_n



(b) Compute $|Y_i - \hat{m}(X_i)|$ on \mathcal{D}_m



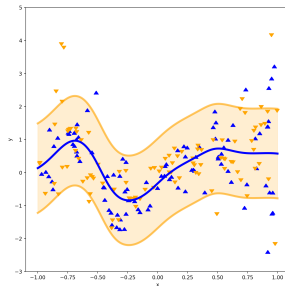
(c) $\hat{C}(X) = [\hat{m}(X) \pm \hat{q}_\alpha]$

- Such evaluation of the prediction quality is performed by a **score function** s

Absolute errors

$$s(X_i, Y_i) \quad |Y_i - \hat{m}(X_i)|$$

But you may have noticed that choosing the absolute errors leads to prediction intervals with **constant width**

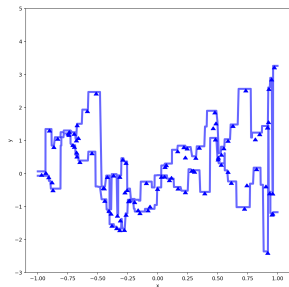


(c) $\hat{C}(X) = [\hat{m}(X) \pm \hat{q}_\alpha]$ 8 / 41

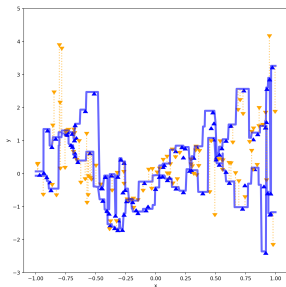
Absolute errors

Quantile regression

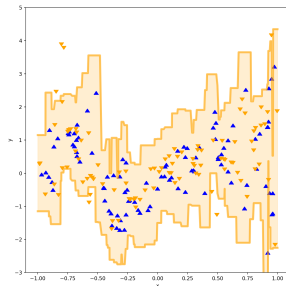
$$s(X_i, Y_i) \quad |Y_i - \hat{m}(X_i)| \quad \max(\hat{q}_l(X_i) - Y_i, Y_i - \hat{q}_u(X_i))$$



(a) Learn \hat{q}_l, \hat{q}_u on \mathcal{D}_n



(b) Compute $\max(\hat{q}_l(X_i) - Y_i, Y_i - \hat{q}_u(X_i))$
on \mathcal{D}_m



(c) $\hat{C} = [\hat{q}_l(X) - \hat{q}_\alpha, \hat{q}_u(X) + \hat{q}_\alpha]$

Absolute errors

Quantile regression

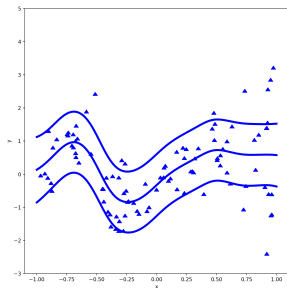
Normalization

$$s(X_i, Y_i)$$

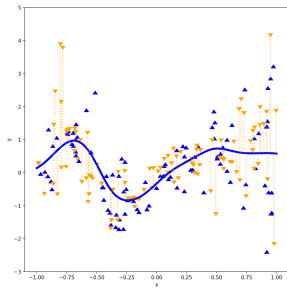
$$|Y_i - \hat{m}(X_i)|$$

$$\max(\hat{q}_l(X_i) - Y_i, Y_i - \hat{q}_u(X_i))$$

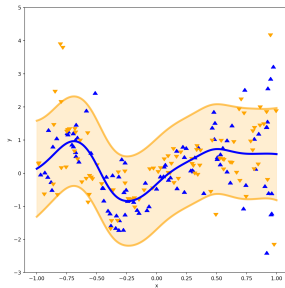
$$\frac{(Y_i - \hat{m}(X_i))^2}{\hat{f}(X_i)}$$



(a) Learn \hat{m}, \hat{f} on \mathcal{D}_n



(b) Compute $\frac{(Y_i - \hat{m}(X_i))^2}{\hat{f}(X_i)}$ on \mathcal{D}_m



(c) $\hat{C} = [\hat{m}(X) \pm \sqrt{\hat{q}_\alpha \hat{f}(X)}]$

Absolute errors		Quantile regression ¹	Normalization ²
$s(X_i, Y_i)$	$ Y_i - \hat{m}(X_i) $	$\max(\hat{q}_l(X_i) - Y_i, Y_i - \hat{q}_u(X_i))$	$\frac{(Y_i - \hat{m}(X_i))^2}{\hat{f}(X_i)}$

- \hat{f} is **any** estimate of the errors of \hat{m} (e.g. other ML models trained on the absolute residuals, resampling procedure, Bayesian approach such as GPs, ...)
- Key fact: this estimation is made without any consideration for coverage nor adaptivity

We then propose to *learn* the score function in a way that targets both adaptivity and coverage

¹[Romano et al. 2019]

²[Lei et al. 2014; Johansson et al. 2014; Papadopoulos 2024; Jaber et al. 2024]



Table of Contents

1. Introduction
2. Learning a score function
3. Experiments
4. Conclusion
5. References and appendix

Learning problem for a score function

- We consider a normalized score: $\frac{(Y - m(X))^2}{f(X)}$, with $f \geq 0$
- As for all learning problems, we must:
 - Specify the criterion to minimize, to be discussed later
 - Choose a search space for our functions, here we rely on **kernel methods**
 - m lives in the Reproducible Kernel Hilbert Space (RKHS) \mathcal{H}^m with kernel k^m and lengthscales θ^m
 - f is a *kernel sum-of-squares* function, in order to impose its **positivity**

Kernel sum-of-squares (kSoS)

- Consider a RKHS \mathcal{H}^f with a feature map $\phi: \mathcal{X} \rightarrow \mathcal{H}^f$, a kernel SoS function is defined as

$$f(X) = \phi(X)^\top \mathcal{A} \phi(X), \quad \text{with } \mathcal{A} \in \mathcal{S}_+(\mathcal{H}^f)$$

- f can be written as

$$f_{\mathcal{A}}(X) = \sum_{l \geq 0} \lambda_l u_l(X) u_l(X)^\top$$

for functions $u_l \in \mathcal{H}^f$ with λ_l the eigenvalues of the operator \mathcal{A} , hence the **sum-of-squares** name

Learning the score function amounts to simultaneously learning

$$m \in \mathcal{H}^m, f \in \mathcal{SoS}(\mathcal{H}^f) \quad \Leftrightarrow \quad m \in \mathcal{H}^m, \mathcal{A} \in \mathcal{S}_+(\mathcal{H}^f)$$

Learning problem for a score function

$$\inf_{m \in \mathcal{H}^m, \mathcal{A} \in \mathcal{S}_+(\mathcal{H}^f)} \frac{a}{n} \sum_{i=1}^n (Y_i - m(X_i))^2 + \frac{b}{n} \sum_{i=1}^n f_{\mathcal{A}}(X_i) + \lambda_1 \|\mathcal{A}\|_{\star} + \lambda_2 \|\mathcal{A}\|_F^2 \quad (1)$$

$$\text{s.t.} \quad f_{\mathcal{A}}(X_i) \geq (Y_i - m(X_i))^2, \quad i \in [n], \quad (2)$$

$$\|m\|_{\mathcal{H}^m}^2 \leq s \quad (3)$$

Learning problem for a score function

$$\inf_{m \in \mathcal{H}^m, \mathcal{A} \in \mathcal{S}_+(\mathcal{H}^f)} \frac{a}{n} \sum_{i=1}^n (Y_i - m(X_i))^2 + \frac{b}{n} \sum_{i=1}^n f_{\mathcal{A}}(X_i) + \lambda_1 \|\mathcal{A}\|_{\star} + \lambda_2 \|\mathcal{A}\|_F^2 \quad (1)$$

$$\text{s.t.} \quad f_{\mathcal{A}}(X_i) \geq (Y_i - m(X_i))^2, \quad i \in [n], \quad (2)$$

$$\|m\|_{\mathcal{H}^m}^2 \leq s \quad (3)$$

i) Faithful estimation of the mean function

Learning problem for a score function

$$\inf_{m \in \mathcal{H}^m, \mathcal{A} \in \mathcal{S}_+(\mathcal{H}^f)} \quad \frac{a}{n} \sum_{i=1}^n (Y_i - m(X_i))^2 + \frac{b}{n} \sum_{i=1}^n f_{\mathcal{A}}(X_i) + \lambda_1 \|\mathcal{A}\|_{\star} + \lambda_2 \|\mathcal{A}\|_F^2 \quad (1)$$

$$\text{s.t.} \quad f_{\mathcal{A}}(X_i) \geq (Y_i - m(X_i))^2, \quad i \in [n], \quad (2)$$

$$\|m\|_{\mathcal{H}^m}^2 \leq s \quad (3)$$

- i) Faithful estimation of the mean function
- ii) 100% coverage on the training sample - **convex** constraint (later adjusted with split CP)

Learning problem for a score function

$$\inf_{m \in \mathcal{H}^m, \mathcal{A} \in \mathcal{S}_+(\mathcal{H}^f)} \quad \frac{a}{n} \sum_{i=1}^n (Y_i - m(X_i))^2 + \frac{b}{n} \sum_{i=1}^n f_{\mathcal{A}}(X_i) + \lambda_1 \|\mathcal{A}\|_{\star} + \lambda_2 \|\mathcal{A}\|_F^2 \quad (1)$$

$$\text{s.t.} \quad f_{\mathcal{A}}(X_i) \geq (Y_i - m(X_i))^2, \quad i \in [n], \quad (2)$$

$$\|m\|_{\mathcal{H}^m}^2 \leq s \quad (3)$$

- i) Faithful estimation of the mean function
- ii) 100% coverage on the training sample - **convex** constraint (later adjusted with split CP)
- iii) Minimization of the interval mean width

Learning problem for a score function

$$\inf_{m \in \mathcal{H}^m, \mathcal{A} \in \mathcal{S}_+(\mathcal{H}^f)} \quad \frac{a}{n} \sum_{i=1}^n (Y_i - m(X_i))^2 + \frac{b}{n} \sum_{i=1}^n f_{\mathcal{A}}(X_i) + \lambda_1 \|\mathcal{A}\|_{\star} + \lambda_2 \|\mathcal{A}\|_F^2 \quad (1)$$

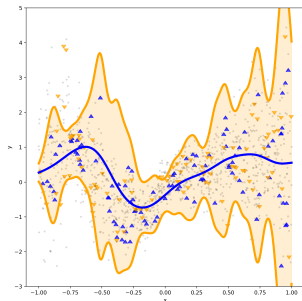
$$\text{s.t.} \quad f_{\mathcal{A}}(X_i) \geq (Y_i - m(X_i))^2, \quad i \in [n], \quad (2)$$

$$\|m\|_{\mathcal{H}^m}^2 \leq s \quad (3)$$

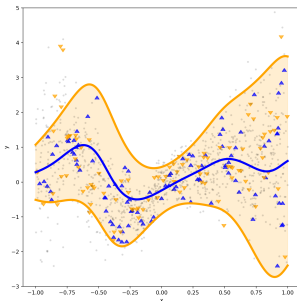
- i) Faithful estimation of the mean function
- ii) 100% coverage on the training sample - **convex** constraint (later adjusted with split CP)
- iii) Minimization of the interval mean width
- iv) Control of the regularity of the bands
 - lasso-type norm $\|\mathcal{A}\|_{\star}$
 - ridge-type norm $\|\mathcal{A}\|_F$

Representer theorem

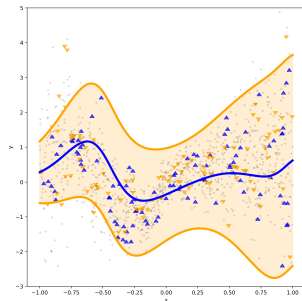
- We proved a representer theorem for this infinite dimensional problem
- It becomes a Semi-Definite Program (SDP) problem, solvable using off-the-shelves solvers



(a) $\theta^f = 0.1$



(b) $\theta^f = 0.5$



(c) $\theta^f = 0.9$

Note: θ^f is the vector of lengthscales for k^f , the kernel corresponding to \mathcal{H}^f

- The SDP problem is not scalable past 200 samples
- We proved that it admits a dual representation if $\lambda_2 > 0$, which is solvable using **accelerated gradient descent**

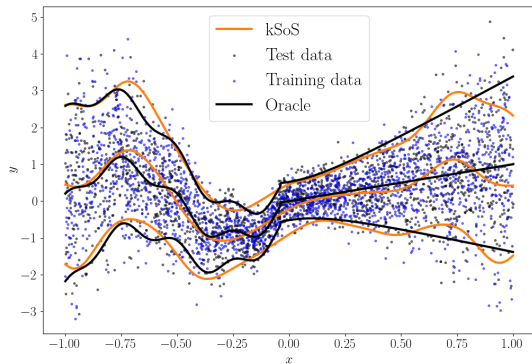


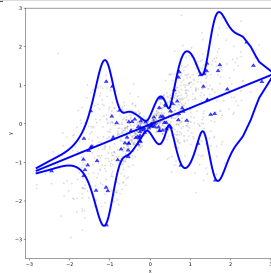
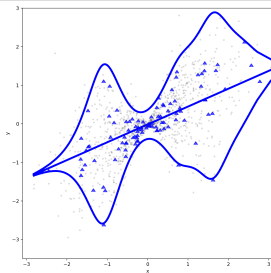
Figure 5: Dual solver with 2000 samples

- Hyperparameters a, λ_1, λ_2 do not have a huge impact on the prediction bands
- We fix θ^m and s using a preliminary Gaussian Process model
- We focus on the two most important hyperparameters
 - b : mean width
 - θ^f : lengthscales associated to f , control complexity of f

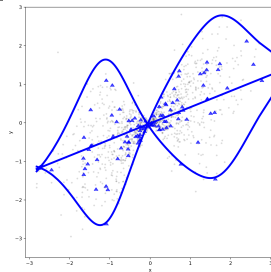
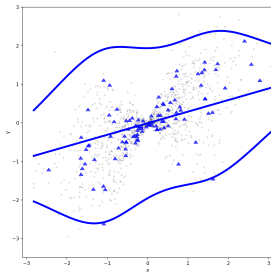
$b = 1$

$b = 100$

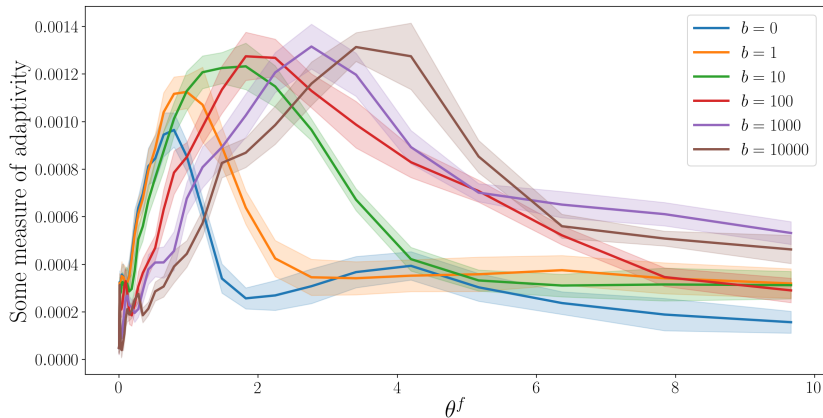
$\theta^f = 0.5$



$\theta^f = 1.74$



Before
calibration

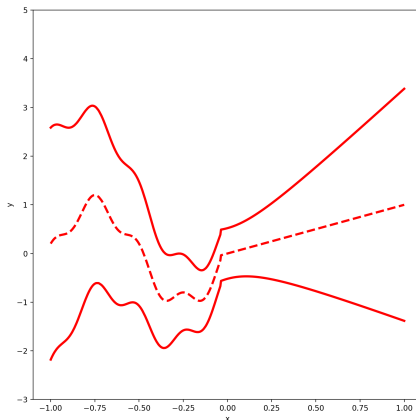


For each b , we can find an optimal value for θ^f that maximizes adaptivity

- Which adaptivity measure can we use to choose θ^f ?
- We propose the Hilbert-Schmidt Independence Criterion (HSIC), an independence measure between random variables
- What is the link between independence of random variables and adaptivity?

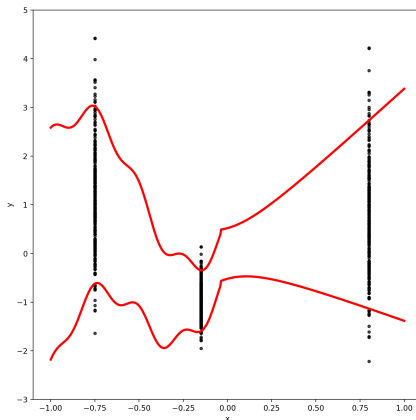
- Perfectly adaptive bands guarantee local coverage

$$\mathbb{P}(Y_{N+1} \in \hat{C}(X_{N+1}) \mid X_{N+1} = x) \geq 1 - \alpha$$



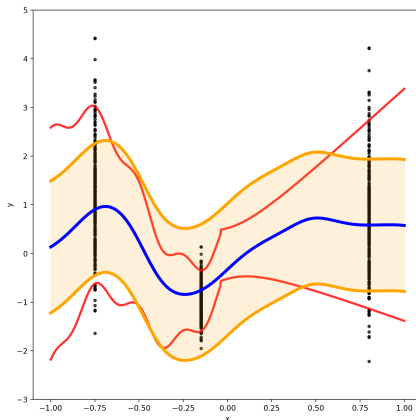
- Perfectly adaptive bands guarantee local coverage

$$\mathbb{P}(Y_{N+1} \in \hat{C}(X_{N+1}) \mid X_{N+1} = x) \geq 1 - \alpha$$



- Perfectly adaptive bands guarantee local coverage

$$\mathbb{P}(Y_{N+1} \in \hat{C}(X_{N+1}) \mid X_{N+1} = x) \geq 1 - \alpha$$



- Without hypothesis on the data, satisfying this local coverage leads to infinitely wide prediction bands [Vovk 2012; Barber et al. 2021]
- We can relax the local coverage by considering X in a small neighbourhood ω_X , such that $\forall x \in \mathcal{X}, \mathbb{P}(x \in \omega_X) \geq \delta$:

$$p_{\mathcal{D}_N} := \mathbb{P}(Y_{N+1} \in \hat{C}(X_{N+1}) | X_{N+1} \in \omega_X) \geq 1 - \alpha$$

- Deutschmann et al. 2024 proved a lower bound for $p_{\mathcal{D}_N}$, which involves $\text{MI}(X, S_{\theta_f}(X, Y))$, but MI is not robust numerically

- Using information theory results and recent inequalities result between the TV distance and the MMD, we proved a new bound

$$p_{\mathcal{D}_N} \geq 1 - \alpha - \frac{1}{\delta} \sqrt{1 - \frac{\alpha_1}{1 - \alpha_2 \text{HSIC}(r_{\mathcal{D}_n}(X_{N+1}, Y_{N+1}), \hat{f}_{\theta^f}(X_{N+1}))}}$$

- HSIC is much more robust than MI

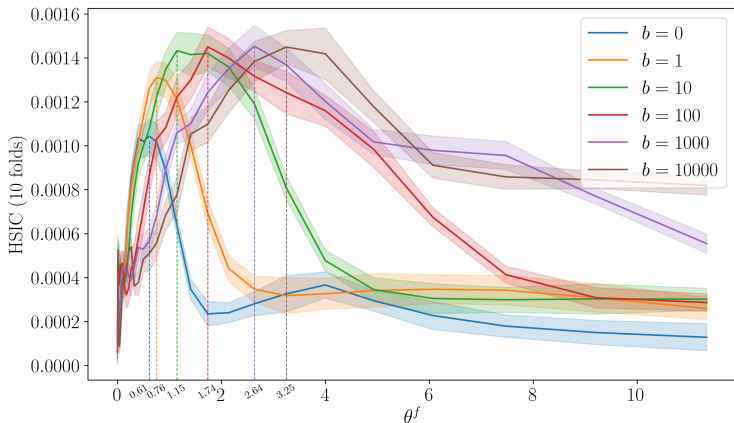


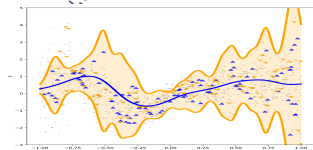
Figure 6: $\text{HSIC}(r_{\mathcal{D}_n}(X, Y), \hat{f}_{\theta^f}(X))$

Maximizing this HSIC, i.e. the dependence between the residuals and the interval widths, allows to target better local coverage

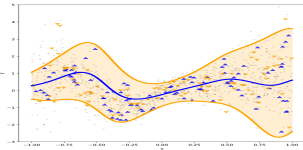


Table of Contents

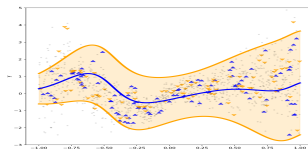
1. Introduction
2. Learning a score function
3. Experiments
4. Conclusion
5. References and appendix



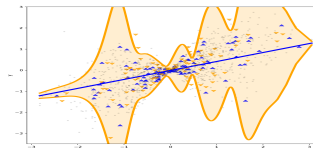
(a) $\theta^f = 0.1$



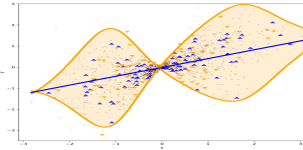
(b) $\theta_{\text{HSIC}}^f = 0.5$



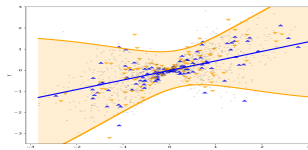
(c) $\theta^f = 0.9$



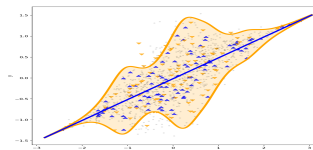
(d) $\theta^f = 0.5$



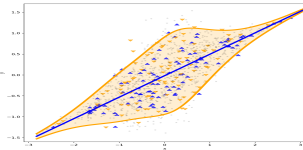
(e) $\theta_{\text{HSIC}}^f = 1.74$



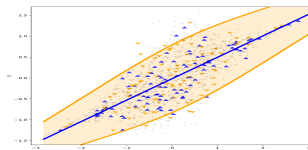
(f) $\theta^f = 10$



(g) $\theta^f = 0.4$

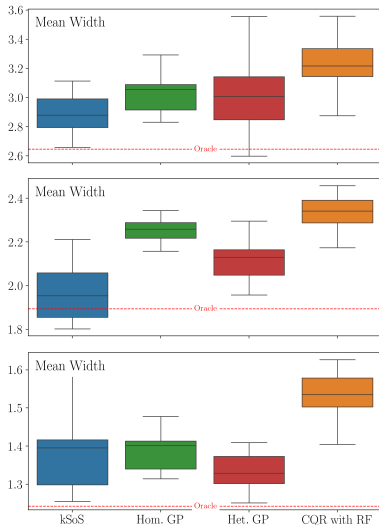


(h) $\theta_{\text{HSIC}}^f = 0.9$

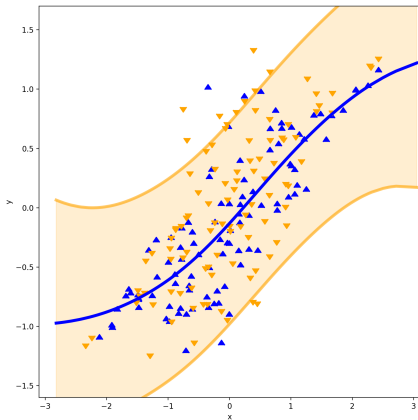


(i) $\theta^f = 3$

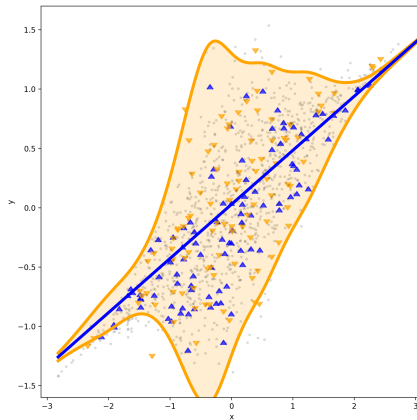
Mean width metric



- A common measure for adaptive prediction bands in the literature is **mean width**, which should be minimized
- kSoS leads to better or as good mean width as competitors
- **However, mean width does not always tell the full story**

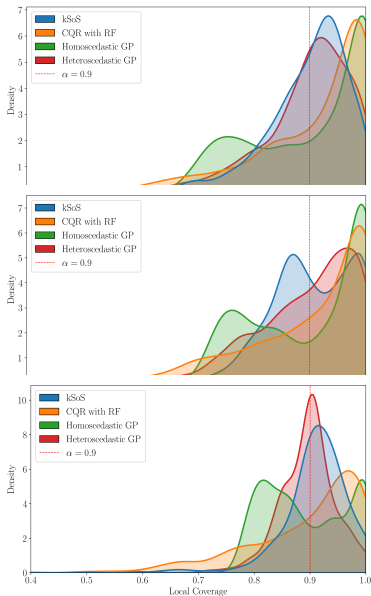


(a) Homoscedastic GP
MW = 1.712



(b) kSoS with Opt. HSIC
MW = 1.759

Local coverage metric



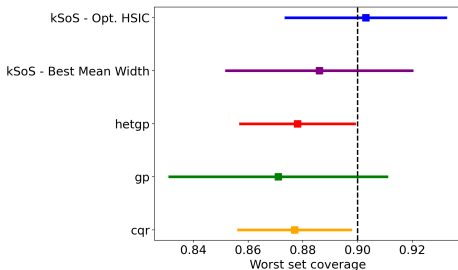
- The best measure of adaptivity is **local coverage**
- The target for local coverage is a Dirac at $1 - \alpha = 0.9$
- kSoS leads to better concentrated local coverage in general

Real world datasets - Comparison between mean widths

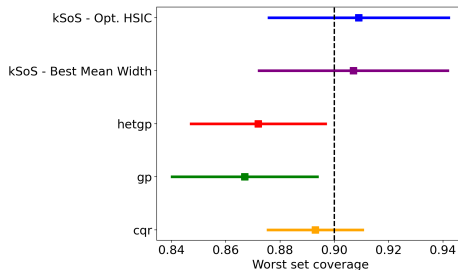
Dataset	CQR	Het GP	Hom GP	kSoS	kSoS
				Best mean width	Opt. HSIC
Concrete	0.586 ± 0.032	0.508 ± 0.052	0.543 ± 0.044	0.556 ± 0.044	0.568 ± 0.06
Bike	1.114 ± 0.062	1.000 ± 0.079	0.809 ± 0.024	0.804 ± 0.032	0.803 ± 0.032
Bio	1.879 ± 0.046	2.21 ± 0.100	2.194 ± 0.119	2.03 ± 0.07	—
Diabetes	188.62 ± 9.33	191.24 ± 11.95	190.58 ± 11.19	185.83 ± 14.47	187.6 ± 16.18
MPG	9.89 ± 0.82	9.70 ± 1.06	9.71 ± 0.73	9.15 ± 0.8	9.36 ± 0.82
Housing	1.816 ± 0.045	1.585 ± 0.099	1.453 ± 0.099	1.468 ± 0.094	1.586 ± 0.104

- Mean width for six real-world datasets, kSoS with HSIC-optimized θ^f achieves best mean width on almost every datasets against competitors
- **Again, mean width does not tell the full story**

Real world datasets - Comparison with worst-set coverage



(a) Housing dataset



(b) Concrete dataset

- Worst set coverage is a substitute for local coverage for real datasets³
- **kSoS achieves better or equal worst-set coverage than competitors with better mean width**

³Thurin et al. 2025



Table of Contents

1. Introduction
2. Learning a score function
3. Experiments
4. Conclusion
5. References and appendix

- Learning setting for a score function in the context of split CP
- Representer theorem to make the problem tractable
- Solvable in practice with the primal (small n) using SDP or the dual (big n) using AGD
- Brand new adaptivity measure based on HSIC, that allows to automatically choose hyperparameters of the model
- **Paper accepted at NeurIPS 2025**, preprint available on arXiv, final version in the proceedings



Q&R





Thank you for listening!




Your feedback will be highly appreciated!





Table of Contents

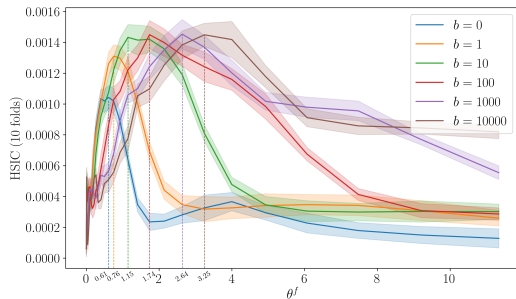
1. Introduction
2. Learning a score function
3. Experiments
4. Conclusion
5. References and appendix

-  Barber, Rina Foygel et al. (2021). “The limits of distribution-free conditional predictive inference”. In: Information and Inference: A Journal of the IMA. DOI: 10.1093/imaiai/iaaa017.
-  Deutschmann, Nicolas et al. (2024). “Adaptive Conformal Regression with Split-Jackknife+ Scores”. In: Transactions on Machine Learning Research. URL: <https://openreview.net/forum?id=1fbTGC3BUD>.
-  Jaber, Edgar et al. (2024). Conformal Approach To Gaussian Process Surrogate Evaluation With Coverage Guarantees. arXiv: 2401.07733 [stat.ML]. URL: <https://arxiv.org/abs/2401.07733>.
-  Johansson, U. et al. (2014). “Regression conformal prediction with random forests”. In: Machine Learning. URL: <https://api.semanticscholar.org/CorpusID:14015369>.

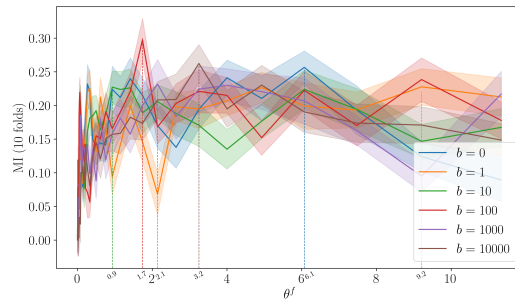
- 
 Lei, Jing et al. (2014). “Distribution-free Prediction Bands for Non-parametric Regression”. In:
Journal of the Royal Statistical Society Series B: Statistical Methodology. DOI:
 10.1111/rssb.12021.
- 
 Papadopoulos, Harris (2024). “Guaranteed Coverage Prediction Intervals With Gaussian Process Regression”. In:
IEEE Transactions on Pattern Analysis and Machine Intelligence. DOI:
 10.1109/tpami.2024.3418214.
- 
 Romano, Yaniv et al. (2019). “Conformalized Quantile Regression”. In:
Advances in Neural Information Processing Systems. Ed. by H. Wallach et al.
 Curran Associates, Inc. URL: https://proceedings.neurips.cc/paper_files/paper/2019/file/5103c3584b063c431bd1268e9b5e76fb-Paper.pdf.

-  Thurin, Gauthier et al. (2025). “Optimal transport-based conformal prediction”. In: Forty-second International Conference on Machine Learning. URL: <https://openreview.net/forum?id=kEAYffH3tn>.
-  Vovk, Vladimir (2012). “Conditional Validity of Inductive Conformal Predictors”. In: Proceedings of the Asian Conference on Machine Learning. Proceedings of Machine Learning Research. PMLR. URL: <https://proceedings.mlr.press/v25/vovk12.html>.

Bound on conditional coverage



(a) HSIC



(b) MI