

Robust Yet Efficient Conformal Prediction Sets

Xiaoyu (Joy) Liang

Beihang University

Meeting with: **Prof. Aleksandar Bojchevski**

University of Cologne

Problem

Conformal Prediction: Returns prediction sets with a distribution-free guarantee to cover the true label.

$$\Pr[\text{airplane} \in \mathcal{C}(\text{airplane})] \geq 1 - \alpha$$

Attacks: Break the guarantee by perturbing the test inputs (**evasion**) or **poisoning** the calibration data.

$$\Pr[\text{airplane} \in \mathcal{C}(\text{airplane} + 0.001 \text{ noise})] \not\geq 1 - \alpha$$

How can we extend robustness certificates to conformal prediction sets?

Robust CP: Recover the same guarantee for the worst-case perturbed input and calibration.

$$\Pr[\text{airplane} \in \bar{\mathcal{C}}(\text{airplane} + 0.001 \text{ noise})] \geq 1 - \alpha$$

Intuition

Conformal Prediction

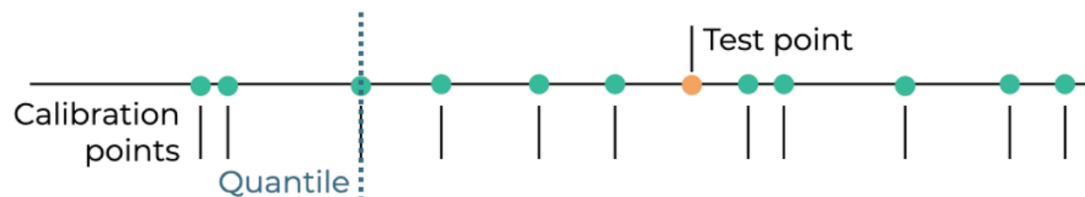
- Conformity score function $s: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, e.g. TPS, APS
- Calibration set $\mathcal{D}_{\text{cal}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$
- Coverage probability $1 - \alpha$, user-specified

Calibration. α -quantile $q = \text{Quant}(\alpha; \{s(\mathbf{x}_i, y_i)\}_{i=1}^n)$.

Prediction. Define the set $\mathcal{C}(x_{n+1}) = \{y: s(x_{n+1}, y) \geq q\}$.

Guarantee. If test and calibration are exchangeable then we have:

$$\Pr[y_{n+1} \in \mathcal{C}(x_{n+1})] \geq 1 - \alpha.$$



Randomized Smoothing Conformal Prediction

- Bound score: $\bar{s}(\xi(\tilde{x}), y) \geq s(x, y), \tilde{x} \in \mathcal{B}(x)$

$$\bar{s}(x, y) = \Phi^{-1}(\mathbb{E}_{\delta \sim \mathcal{N}(0, \sigma^2 I)}[s(x + \delta, y)])$$

1. Loose bound 2. Continuous

3. Poison attack

Calibration. $q = \text{Quant}(\alpha; \{\bar{s}(\mathbf{x}_i, y_i)\}_{i=1}^n) - M_{\mathcal{B}}$.

Prediction. $\bar{\mathcal{C}}(x_{n+1}) = \{y: \bar{s}(\xi(x_{n+1}), y) \geq q\}$.

Guarantee.

$$\Pr[y_{n+1} \in \bar{\mathcal{C}}(\xi(x_{n+1}))] \geq 1 - \alpha.$$

4. Finite sample error

1. CDF Aware smooth scores.
2. Sparse smoothing.
3. Lower bounds on calibration for poison.
4. Finite sample correction

Robust CP

How can we build conservative prediction set against evasion and poison?

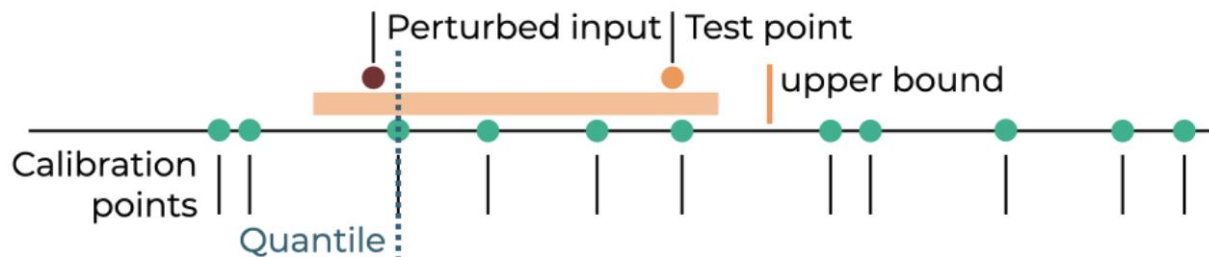
Robustness to Evasion Attacks

Goal: $\Pr[y_{n+1} \in \mathcal{C}(\tilde{x}_{n+1})] \geq 1 - \alpha$

Method: Upper bound of score $\bar{s}(\xi(\tilde{x}), y)$

Conservative prediction set

$$\bar{\mathcal{C}}(x_{n+1}) = \{y: \bar{s}(\xi(\tilde{x}_{n+1}), y) \geq q\}.$$



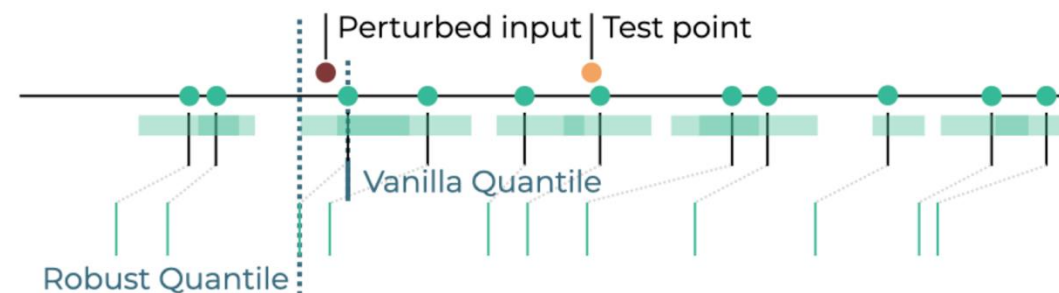
Robustness to Poison Attacks

Method: Lower bound of calibration

$$\left. \begin{array}{l} \text{Feature } \underline{q}_\alpha = \min_{z_i \in \mathcal{X}} \text{Quant}(\alpha; \{s(\mathbf{z}_i, y)\}_{i=1}^n) \\ \text{Label } \underline{q}_\alpha = \min_{z_i \in \mathcal{Y}} \text{Quant}(\alpha; \{s(\mathbf{x}_i, \mathbf{z}_i)\}_{i=1}^n) \end{array} \right\} \text{MILP}$$

Conservative prediction set

$$\bar{\mathcal{C}}(x_{n+1}) = \{y: s(x_{n+1}, y) \geq \underline{q}_\alpha\}.$$



Randomized Smoothing Bounds

Smooth scores

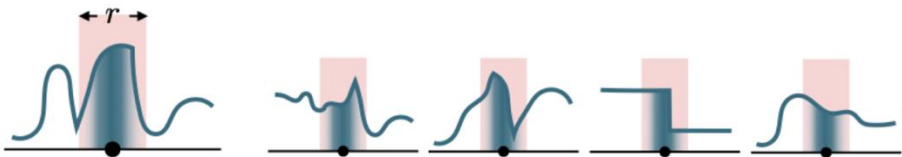
Gaussian/sparse smoothing score: $\hat{s}(x, y) := \mathbb{E}(s(\xi(x), y))$

Bound score for worst case input: $\max_{\tilde{x} \in \mathcal{B}(x)} \mathbb{E}(s(\xi(\tilde{x}), y))$

Relax by search worst case score function:

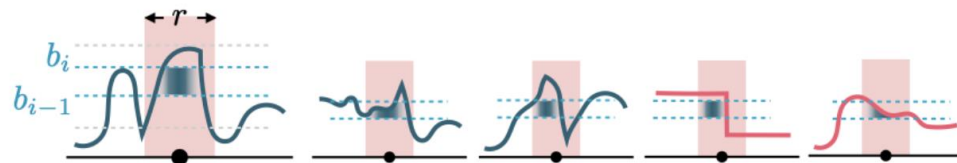
$$\max_{\tilde{x} \in \mathcal{B}(x)} \mathbb{E}(s(\xi(\tilde{x}), y)) \leq \max_{\tilde{x} \in \mathcal{B}(x), h \in \mathcal{H}} \mathbb{E}(h(\xi(\tilde{x}), y))$$

Bound h with mean: $\mathbb{E}(h(\xi(\tilde{x}), y)) = \mathbb{E}(s(\xi(x), y))$



Bound h with **CDF**: $\forall b_i: \Pr[h(\xi(\tilde{x}), y) \leq b_i] = p_i$

$$a = b_1 < b_2 \cdots b_m = b, p_i = \Pr[s(\xi(x), y) \leq b_i]$$



How can we build smaller conservative sets?

Bounds for smooth scores

Gaussian smoothing

$$\bar{s}_{\text{mean}}(\mathbf{x}, y) = \Phi(\Phi^{-1}(p) + r)$$

$$\bar{s}_{\text{cdf}}(\mathbf{x}, y) = b_m - \sum_{j=2}^{m-1} \Phi(\Phi^{-1}(p_j) - r)(b_{j+1} - b_j)$$

Sparse smoothing

$\bar{s}_{\text{mean}}(\mathbf{x}, y)$: Greedily distribute the p mass to each constant likelihood region.

$\bar{s}_{\text{cdf}}(\mathbf{x}, y)$: Distribute the p_j masses in each region and each bin $[b_j, b_j + 1]$.

$$\max_{\mathbf{H}} b_m - \mathbf{H} \tilde{\mathbf{t}} \mathbf{d} \quad \text{s.t. } \mathbf{H} \mathbf{t} = \mathbf{p}, 0 \leq \mathbf{H} \leq 1$$

Robust Yet Efficient CP

How can we build efficient and correct conservative sets?

CAS: CDF-Aware Sets

Algorithm 1 CDF-Aware Sets (CAS, Evasion)

$q_\alpha = \text{Quant}(\alpha; \{\hat{s}(\mathbf{x}, y)\}_{(\mathbf{x}, y) \in \mathcal{D}_{\text{cal}}}) \triangleright$ Clean quantile
Compute $\bar{s}_{\text{cdf}}(\mathbf{x}, y)$, e.g. with Eq. 10 \triangleright Upper bound
Return $\bar{\mathcal{C}}_\alpha = \{y : \bar{s}_{\text{cdf}}(\mathbf{x}, y) \geq q_\alpha\} \triangleright$ Conservative set

Calibration-time variant: Compare the smooth test score $\hat{s}(\mathbf{x}, y)$ against conservative (lower) quantile.

$$\underline{q}_\alpha = \text{Quant}(\alpha; \{\underline{s}(\mathbf{x}_i, y_i)\})$$

Poisoning: Conservative threshold $\underline{q}_\alpha = \min_z \text{Quant}$

Combined: $\bar{\mathcal{C}}_\alpha = \{y : \bar{s}_{\text{cdf}}(x, y) \geq \underline{q}_\alpha\}$

Finite Sample Correction

Monte-Carlo sample for mean/CDF: $1 - \eta$ confidence

Calibrate $\alpha' = \alpha - \eta$

Corrected CDF score: $\underline{s}_{\text{cdf}+}(\mathbf{x}_i, y_i) \leq \underline{s}_{\text{cdf}}$ with $1 - \eta/(2|\mathcal{D}_{\text{cal}}|)$
DKW inequality

Corrected smooth score: $\hat{s}_+(\tilde{\mathbf{x}}_i, y_i) \geq \hat{s}$ with $1 - \eta/(2|y|)$
Bernstein bound

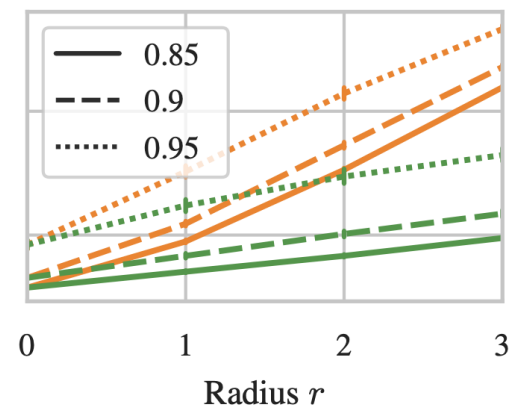
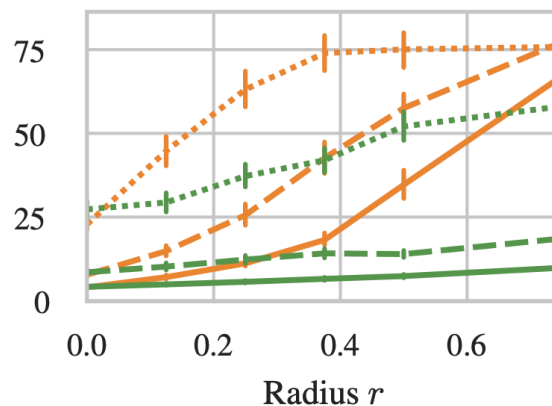
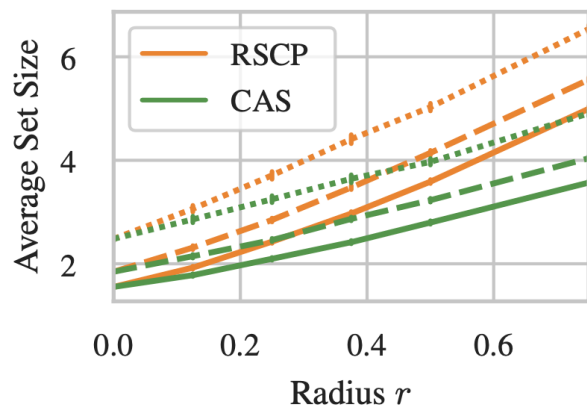
Conservative threshold: $\underline{q}_{\alpha+} = \text{Quant}(\alpha - \eta; \{\underline{s}_{\text{cdf}+}(\mathbf{x}_i, y_i)\})$

Conservative set: $\bar{\mathcal{C}}_{\alpha+}(x_{n+1}) = \{y : \hat{s}_+(x_{n+1}, y) \geq \underline{q}_{\alpha+}\}$

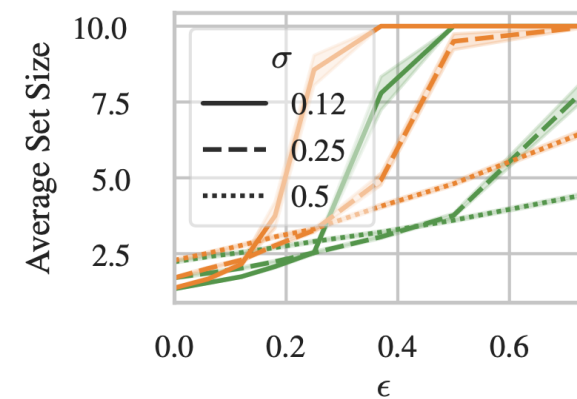
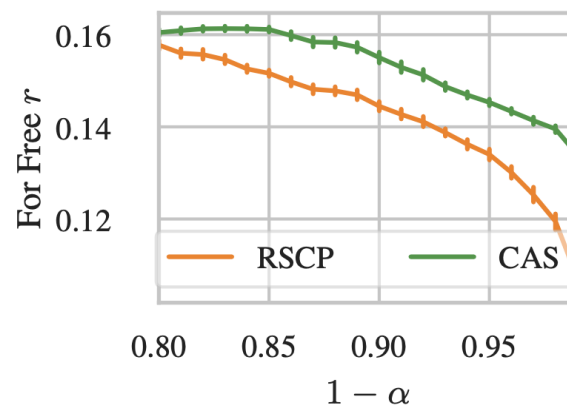
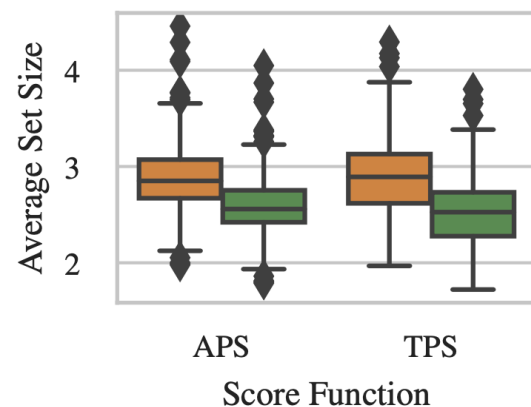
Experiments

CAS returns smaller prediction set with the same guarantee. CAS results in larger “for free” radius.

All radii
All dataset
All guarantee

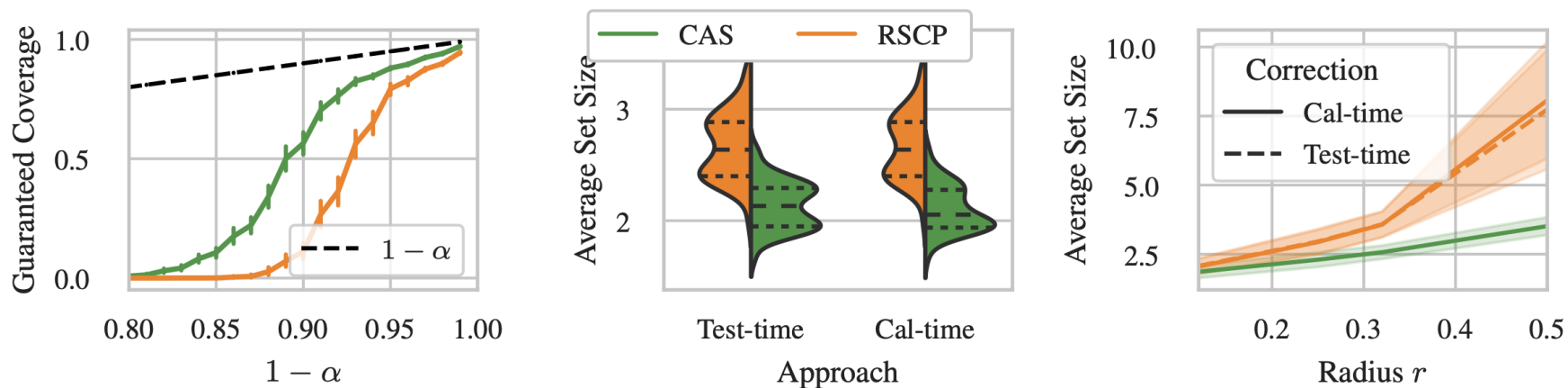


All score
All smoothing

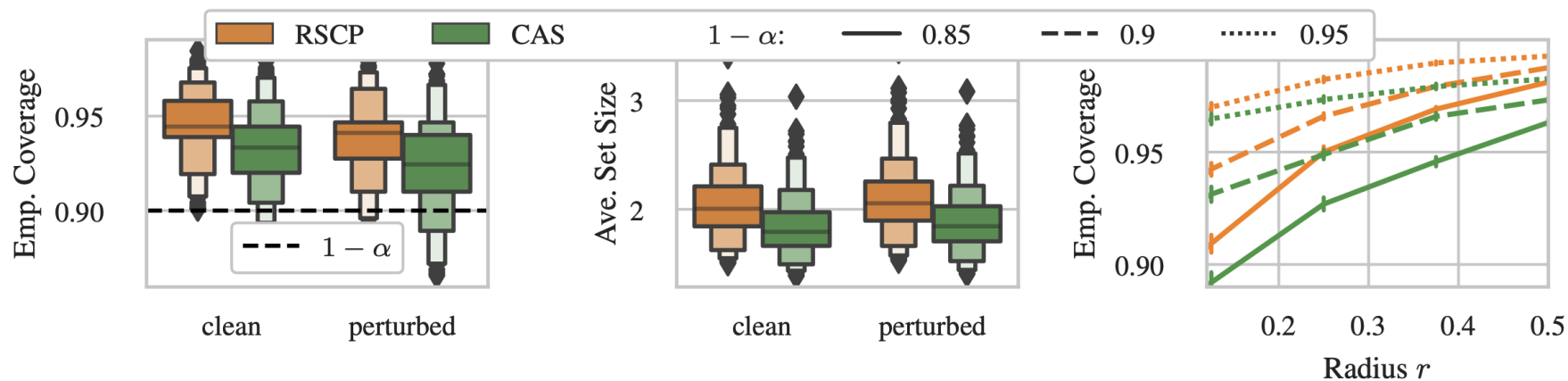


Experiments

CAS certifies a larger lower bound for vanilla CP. CAS shows smaller sets for the cal-time certificate and with correction.

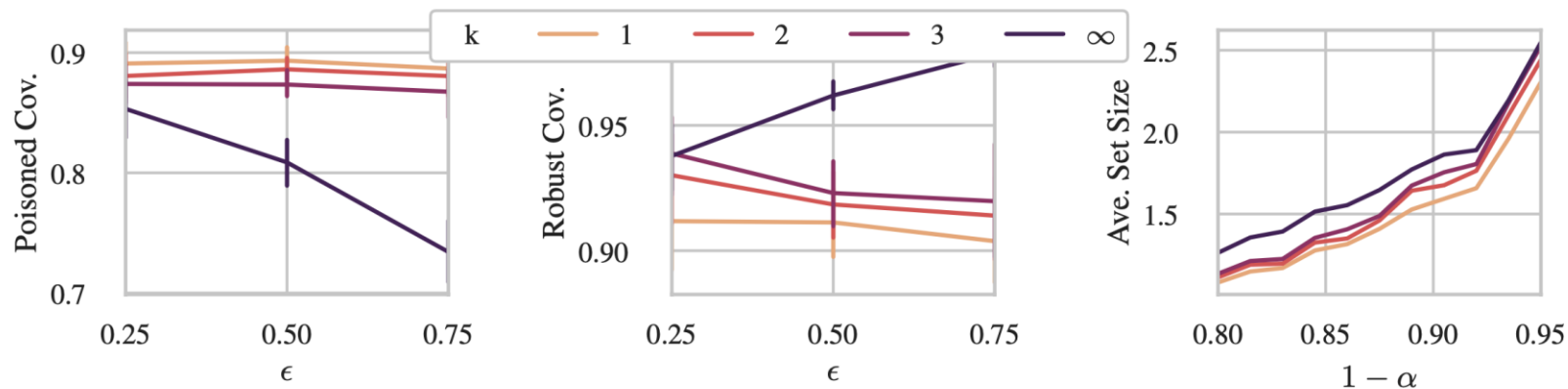


CAS is less conservative since it is closer to the nominal $1 - \alpha$, and has smaller sets.



Experiments

CAS makes CP robust to feature poisoning, coming at only a small cost



CAS makes CP robust to label poisoning, with a predictably larger set size. CAS is also faster.

Table 1. Label poisoning for CIFAR-10.

k	Cov. (Clean)	Vanilla Cov. (Pert)	Robust Cov. (Pert)	Set Size (Clean)
0	0.897	0.897	0.897	1.41
1	0.916	0.872	0.900	1.58
2	0.923	0.859	0.901	1.62

	Time (seconds)		No. Data
	Baseline	CAS (Ours)	
Calib.	0.15	0.79	204
Testing	2.93	0.15	100
Total	3.08	0.94	

Conclusion

Robust Yet Efficient Conformal Prediction Sets

- We provide certified robustness for conformal prediction both for evasion and poisoning attacks.

Upper bound of score & Lower bound of calibration

- We propose a CDF-aware bound on the conformity scores under adversarial perturbation.

$$Pr[h(\xi(\tilde{x}), y) \leq b_i] = Pr[s(\xi(x), y) \leq b_i]$$

- We generalize both results to discrete and binary (sparse) data.

Sparse smoothing

- We show how we can correct for finite-sample error.

Tighter DKW & Bernstein bound

Robust Yet Efficient Conformal Prediction Sets

Thank you for your attention. Looking forward to your questions!

Supplement

Upper-bound score function

Proposition 3.1. Define $\bar{s}(\mathbf{x}, y)$ as the upper bound for $\{s(\tilde{\mathbf{x}}, y) : \tilde{\mathbf{x}} \in \mathcal{B}(\mathbf{x})\}$. With q_α as the α -quantile of the true (clean) calibration scores, let $\bar{\mathcal{C}}_\alpha(\mathbf{x}) = \{y : \bar{s}(\mathbf{x}, y) \geq q_\alpha\}$. For all $\tilde{\mathbf{x}}_{n+1} \in \mathcal{B}(\mathbf{x}_{n+1})$, if $(\mathbf{x}_{n+1}, y_{n+1})$ is exchangeable with \mathcal{D}_{cal} then we have $\Pr [y_{n+1} \in \bar{\mathcal{C}}_\alpha(\tilde{\mathbf{x}}_{n+1})] \geq 1 - \alpha$.

Lower-bound quantile

Proposition 3.2. Let \underline{q}_α to be the solution to the optimization problem in Eq. 3. With the conservative prediction sets

$$\bar{\mathcal{C}}_\alpha(\mathbf{x}_{n+1}) = \{y_i : s(\mathbf{x}_{n+1}, y_i) \geq \underline{q}_\alpha\} \quad (4)$$

for any $(\mathbf{x}_{n+1}, y_{n+1})$ exchangeable with (clean) \mathcal{D}_{cal} we have $\Pr [y_{n+1} \in \bar{\mathcal{C}}_\alpha(\mathbf{x}_{n+1})] \geq 1 - \alpha$.

Score function mean bound

Baseline bound. A straightforward approach only incorporates the **expected smoothed score (mean)** for the given input \mathbf{x} . Let $p = \mathbb{E}[s(\xi(\mathbf{x}), y)]$ for simplicity. With $\tilde{\mathbf{x}} \in \mathcal{B}(\mathbf{x})$ the baseline upper-bound for $\hat{s}(\tilde{\mathbf{x}}, y) = \mathbb{E}[s(\xi(\tilde{\mathbf{x}}), y)]$ is determined by the following problem:

$$\begin{aligned} \bar{s}_{\text{mean}}(\mathbf{x}, y) &= \max_{h \in \mathcal{H}} \mathbb{E}[h(\xi(\tilde{\mathbf{x}}), y)] \\ \text{s.t. } &\mathbb{E}[h(\xi(\mathbf{x}), y)] = p \end{aligned} \quad (7)$$

Score function CDF bound

CDF-based bound. Let $a = b_1 < b_2 \leq \dots \leq b_{m-1} < b_m = b$ be **m real numbers that partition the output space.** Let $p_i = \Pr [s(\xi(\mathbf{x}), y) \leq b_i]$. We define the problem:

$$\begin{aligned} \bar{s}_{\text{cdf}}(\tilde{\mathbf{x}}, y) &= \max_{h \in \mathcal{H}} \mathbb{E}[h(\xi(\tilde{\mathbf{x}}), y)] \\ \text{s.t. } &\forall b_i : \Pr [h(\xi(\mathbf{x}), y) \leq b_i] = p_i \end{aligned} \quad (8)$$

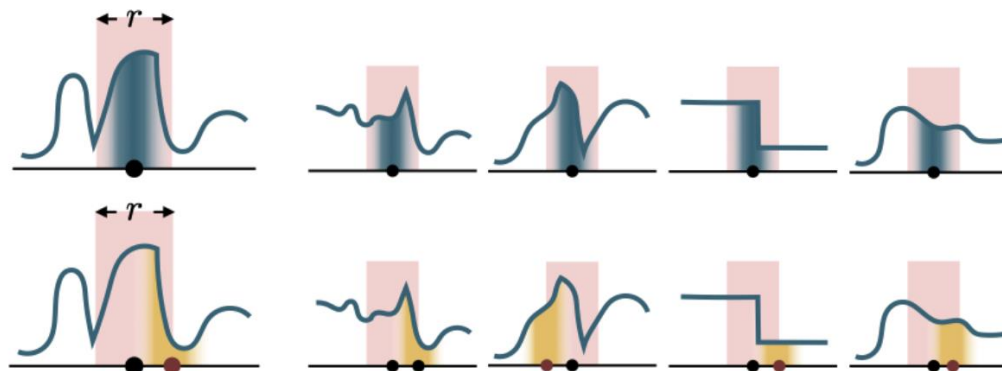
Supplement

CDF Aware Sets (CAS)

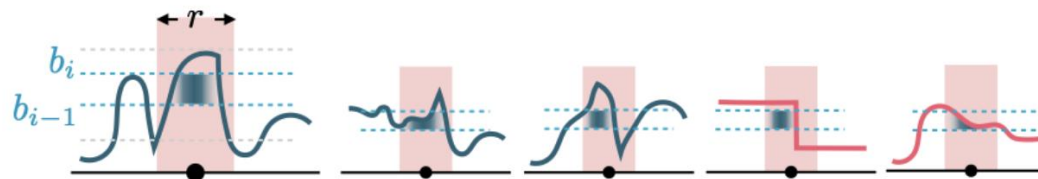
Bound Through Smoothing. We search over the space of all possible functions (due to black-box access).

$$\max_{\tilde{\mathbf{x}} \in \mathcal{B}(\mathbf{x})} \mathbb{E}[s(\tilde{\mathbf{x}} + \epsilon, y)] \leq \max_{\tilde{\mathbf{x}} \in \mathcal{B}(\mathbf{x}), h \in \mathcal{H}} \mathbb{E}[h(\tilde{\mathbf{x}} + \epsilon, y)]$$

Mean Constraint. We limit the search space to all functions with the same mean around the reference point. $\mathbb{E}[h(\mathbf{x} + \epsilon, y)] = \mathbb{E}[s(\mathbf{x} + \epsilon, y)]$



CDF Constraint. We limit the space to all functions with same CDF around the reference point. $\Pr[h(\mathbf{x} + \epsilon, y) \leq b_i] = \Pr[s(\mathbf{x} + \epsilon, y) \leq b_i]$



Supplement

Bound sparse smoothing

For the baseline bound we can rewrite Eq. 7 as a linear program by partitioning the input space \mathcal{X} into regions of constant likelihood ratio (Lee et al., 2019). Let $\mathcal{X} = \bigcup_i \mathcal{R}_i$ and $\mathcal{R}_i \cap \mathcal{R}_j = \emptyset$ be a partitioning into disjoint regions of constant likelihood ratio such that for ever $\mathbf{z} \in \mathcal{R}_i$ it holds $\frac{\Pr[\xi(\mathbf{x})=\mathbf{z}]}{\Pr[\xi(\tilde{\mathbf{x}})=\mathbf{z}]}$ = c_i for some constant c_i . Let $t_i = \Pr[\xi(\mathbf{x}) \in \mathcal{R}_i]$ and $\tilde{t}_i = \Pr[\xi(\tilde{\mathbf{x}}) \in \mathcal{R}_i]$ for for each region \mathcal{R}_i . Then Eq. 7 is equivalent to:

$$\max_{\mathbf{h}} \mathbf{h}^T \tilde{\mathbf{t}} \quad \text{s.t.} \quad \mathbf{h}^T \mathbf{t} = p, \quad 0 \leq \mathbf{h} \leq 1 \quad (17)$$

where $\mathbf{h} \in [0, 1]^I$ is the vector that we are optimizing over corresponding to the score function $h \in \mathcal{H}$, \mathbf{t} and $\tilde{\mathbf{t}}$ are the vectors with t_i and \tilde{t}_i as elements, and $I = r_a + r_d + 1$ is the number of regions. Note, that by replacing the constraint with $a \leq \mathbf{h} \leq b$ we can handle score functions that are bounded in $[a, b]$. The exact solution to this LP can be easily obtained with a simple algorithm. We visit each region in increasing order w.r.t. c_i where

$$c_i = \left[\frac{p_0}{1 - p_1} \right]^{i - r_d} \left[\frac{p_1}{1 - p_0} \right]^{i - r_a} \quad (18)$$

and assign $h_i = 1$ for all regions \mathcal{R}_i until the budget constraint is met, and $h_i = 0$ for the remaining regions, with the exception of the region in between where h_i is a value between 0 and 1 such that the equality constraint is exactly met. Since the likelihood ratios c_i are monotonic in i , the regions are automatically sorted so the solution to the LP can be obtain in linear $O(I)$ time. See Bojchevski et al. (2020) for more details and the pseudo-code.

For the CDF-based bound we can similarly rewrite Eq. 8 as the following linear program:

$$\max_{\mathbf{H}} b_m - \mathbf{H} \tilde{\mathbf{t}} \mathbf{d} \quad \text{s.t.} \quad \mathbf{H} \mathbf{t} = \mathbf{p}, \quad 0 \leq \mathbf{H} \leq 1 \quad (19)$$

where $\mathbf{H} \in [0, 1]^{(m-1) \times I}$ is the matrix that we are optimizing over with H_{ji} being the score that we assign to the j -th bin and the i -th region, \mathbf{d} is the vector of bin widths such that $d_j = b_j - b_{j-1}$, and \mathbf{p} is a vector where $p_i = \Pr[s(\xi(\mathbf{x}), y) \leq b_i]$. Intuitively, for each bin and each region the worst-case score function $h \in \mathcal{H}$ assigns the same score to all \mathbf{z} in that region since the likelihood ratio is constant. As before we have a simple algorithm to obtain the exact solution to this LP. Observe that Eq. 19 can be decomposed into $m - 1$ separate LPs similar to Eq. 17 which can be solved in parallel using the same algorithm as above. The reason is that there is no interaction between the different bins (different rows of \mathbf{H}) in neither the constraint nor the objective function. Therefore, the solution can be obtain in $O(m \times I)$ with serial computation and $O(I)$ with parallel computation.

Supplement

Calibration-time

Proposition 5.1. For $\tilde{\mathbf{x}}_{n+1} \in \mathcal{B}(\mathbf{x}_{n+1})$ and $(\mathbf{x}_{n+1}, y_{n+1})$ exchangeable with \mathcal{D}_{cal} , define

$$\underline{q}_\alpha = \text{Quant}(\alpha; \{\underline{s}(\mathbf{x}_i, y_i) : (\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{cal}}\}) \quad (11)$$

For prediction sets $\bar{\mathcal{C}}_\alpha(\tilde{\mathbf{x}}_{n+1}) = \{y : \hat{s}(\tilde{\mathbf{x}}_{n+1}, y) \geq \underline{q}_\alpha\}$ we have $\Pr[y_{n+1} \in \bar{\mathcal{C}}_\alpha(\tilde{\mathbf{x}}_{n+1})] \geq 1 - \alpha$. Moreover, the vanilla CP covers the true label with probability $\geq 1 - \beta$ for

$$\beta = \text{Quant}^{-1}(q_\alpha; \{\underline{s}(\mathbf{x}_i, y_i) : (\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{cal}}\}) \quad (12)$$

and $\text{Quant}^{-1}(t; A) = \min\{\tau' : \text{Quant}(\tau'; A) \geq t\}$.

Finite sample correction - Evasion

Proposition 6.1. Let $\underline{s}_{\text{cdf}+}(\mathbf{x}_i, y_i) \leq \underline{s}_{\text{cdf}}(\mathbf{x}_i, y_i)$ hold with $1 - \eta/(2|\mathcal{D}_{\text{cal}}|)$ probability for each $(\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{cal}}$, and $\hat{s}_+(\tilde{\mathbf{x}}_{n+1}, y_{n+1}) \geq \hat{s}(\tilde{\mathbf{x}}_{n+1}, y_{n+1})$ hold with $1 - \eta/(2|\mathcal{Y}|)$ probability. Define the conservative $\underline{q}_{\alpha+} = \text{Quant}(\alpha - \eta; \{\underline{s}_{\text{cdf}+}(\mathbf{x}_i, y_i) : (\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{cal}}\})$ and $\bar{\mathcal{C}}_{\alpha+}(\mathbf{x}_{n+1}) = \{y : \hat{s}_+(\mathbf{x}_{n+1}, y) \geq \underline{q}_{\alpha+}\}$. Then

$$\Pr[y_{n+1} \in \bar{\mathcal{C}}_{\alpha+}(\tilde{\mathbf{x}}_{n+1})] \geq 1 - \alpha \quad (13)$$

Finite sample correction - Poison

Proposition 6.2. Let $\underline{s}_{\text{cdf}+}(\tilde{\mathbf{x}}_i, y_i) \leq \underline{s}_{\text{cdf}}(\tilde{\mathbf{x}}_i, y_i)$ hold with $1 - \eta/(2|\mathcal{D}_{\text{cal}}|)$ probability for all $(\tilde{\mathbf{x}}_i, y_i) \in \mathcal{D}_{\text{cal}}$. Let $\underline{q}_{\alpha+}$ be the solution to Eq. 24 (Eq. 3 with CDF bounds) for $\alpha = \alpha' - \eta$ with $\underline{s}_i = \underline{s}_{\text{cdf}+}(\tilde{\mathbf{x}}_i, y_i) - \epsilon_{\text{hoef}}$. Then for each new test point \mathbf{x}_{n+1} exchangeable with \mathcal{D}_{cal} the prediction set defined as $\bar{\mathcal{C}}(\mathbf{x}_{n+1}) = \{y : s_{\text{mc}}(\mathbf{x}_{n+1}, y) \geq \underline{q}_{\alpha+}\}$ has $1 - \alpha'$ coverage.

Supplement

Finite sample correction bound

any adjustable $1 - \eta$ probability;

$$|\mathbb{E}[z] - \mathbb{E}_m[z]| \leq \sqrt{\frac{\log(\frac{2}{\eta})}{2m}}$$

This bound only accesses to the empirical mean and not the samples. Therefore, in cases where we want to account for the distance of empirical mean, and the true expectation for an unknown variable, we can use this bound. An example of this case is the test-time correction where the upper bound on the mean of the unseen point is computed while the empirical mean is not computable (since there are no samples).

Let σ_m^2 be the variance of the MC samples, then **empirical Bernstein inequality** produces variance-dependent confidence intervals as following:

$$|\mathbb{E}[z] - \mathbb{E}_m[z]| \leq \sqrt{2\sigma_m^2 \frac{\ln(\frac{4}{\eta})}{m}} + \frac{7 \ln(\frac{4}{\eta})}{3(m-1)}$$

Similar to the mean, the **empirical CDF** is also bounded between an upper and a lower CDF, via the Dvoretzky–Kiefer–Wolfowitz **(DKW) inequality** (Dvoretzky et al., 1956). Let $F(b_i) = \Pr[z \leq b_i]$ and $F_m(b_i) = \sum_{j=1}^m \mathbf{1}[z_j \leq b_i]$,

$$|F(b_i) - F_m(b_i)| \leq \sqrt{\frac{\log(\frac{2}{\eta})}{2m}}$$

The above inequality holds simultaneously for all b_i .

For Eq. 7 we use the Bernstein inequality as it has shown a better empirical result compared to Hoeffding's inequality. For Eq. 8 we use the DKW inequality to find confidence intervals for the empirical CDF.