

Guaranteed prediction sets via concentration inequalities1. Introduction

Split conformal prediction:

$(X_1, Y_1), \dots, (X_n, Y_n), (X, Y)$ exchangeable

$s(x, y)$ "non-conformity score" (e.g., for regression, $s(x, y) = |y - f(x)|$ with predictor f)

$s_i := s(X_i, Y_i)$, $i = 1, \dots, n$

$\mapsto s_{(1)} \leq s_{(2)} \leq \dots \leq s_{(n)}$ "order statistics"

Prediction set (for a fixed $\alpha \in [\frac{1}{n+1}, 1)$):

$$C(x) := \left\{ y \in \mathcal{Y} : s(x, y) \leq s_{(\lceil (1-\alpha)(n+1) \rceil)} \right\}$$

Marginal coverage guarantee: $\mathbb{P}(Y \in C(X)) \geq 1 - \alpha$

\hookrightarrow w.r.t to $\underbrace{(X_1, Y_1), \dots, (X_n, Y_n), (X, Y)}_{\mathcal{D}_n}$

Question: can we derive calibration-conditional guarantees of the form

$$\mathbb{P} \left[\mathbb{P}(Y \in C(X) | \mathcal{D}_n) \geq 1 - \alpha \right] \geq 1 - \delta \quad ?$$

\uparrow
for most calibration data $\mathcal{D}_n = (X_i, Y_i)_{1 \leq i \leq n}$
the coverage on test points is $\geq 1 - \alpha$

Somewhat stronger than

$$\mathbb{E} \left[\mathbb{P}(Y \in C(X) | \mathcal{D}_n) \right] = \mathbb{P}(Y \in C(X)) \geq 1 - \alpha.$$

2. A "tolerance region" viewpoint on split conformal prediction

2.1. Reminder: tolerance regions

Wilks (1941), Wald (1943), Scheffé and Tukey (1945), Tukey (1947), etc
Suppose Z_1, \dots, Z_n are i.i.d. copies of a real-valued r.v. Z

F : cdf of Z

ν : law of Z

Problem: construct a "tolerance region" $[q_1(Z_1, \dots, Z_n), q_2(Z_1, \dots, Z_n)] = [T_1, T_2]$
such that $\mathbb{P}(\nu([T_1, T_2]) \geq 1 - \alpha) \geq 1 - \delta$.

$[T_1, T_2]$ covers (at least) a $1 - \alpha$ fraction of the population mass. \uparrow for most observed n -samples

One-sided alternative: $\mathbb{P}(\nu((-\infty; T_1]) \geq 1 - \alpha) \geq 1 - \delta$.
 $\uparrow T_1 \geq \text{quantile}_{1-\alpha}(\nu)$

Solution: T_1 and T_2 of the form $T_i = Z_{(k_i)}$ (order statistics)

Classical (Wilks '41, Scheffé and Tukey '45) "Wilks' method"

• If F is continuous, $\nu([Z_{(r)}, Z_{(m+1-r)}]) \sim \text{Beta}(m+1-2r, 2r)$

$\nu((-\infty; Z_{(k)}]) \sim \text{Beta}(k, m+1-k)$

• If F is not continuous, the random masses $\nu(\cdot)$ above stochastically dominate the Beta distributions above. e.g.,
dominated by $\stackrel{\text{sto}}{\geq}$ below

$\mathbb{P}(\nu((-\infty; Z_{(k)}]) \geq t) \geq \mathbb{P}(B \geq t)$ with $B \sim \text{Beta}(k, m+1-k)$

In particular, if $\mathbb{P}(B \geq 1 - \alpha) \geq 1 - \delta$, then this also holds for $\nu((-\infty; Z_{(k)}])$

Proof ingredient: simulate the Z_i 's by inverse transform sampling,
and reduce the problem to the case $\mathcal{D} = \mathcal{U}([0,1])$, up to an
inequality if F is not entreeses.

2.2. Application to split conformal predictors (see Wass 2012)

In all the sequel, we assume that $\underbrace{(X_1, Y_1), \dots, (X_m, Y_m)}_{\mathcal{D}_m}$ are independent copies of (X, Y) ,
and that $\alpha \in \left[\frac{\ln(1/\delta)}{m}, 1 \right)$

Recall that $C(x) = \{y \in \mathcal{Y} : s(x, y) \leq s(x)\}$

so that $\mathbb{P}(Y \in C(x) | \mathcal{D}_m) = \mathbb{P}(s(X, Y) \leq s(x) | \mathcal{D}_m)$
 $= \mathcal{D}((-\infty, s(x)))$ with $\mathcal{D} := \text{law of } s(X, Y)$
 $\stackrel{\text{to}}{\approx} \text{Beta}(k, m+1-k)$ $s(x) = k\text{-th order statistic of } s(X_i, Y_i), i=1, \dots, m$

So get $\mathbb{P}\left[\mathbb{P}(Y \in C(x) | \mathcal{D}_m) \geq 1-\alpha\right] \geq 1-\delta$, we thus pick $k \in \{1, \dots, m\}$ such that

of. Appendix A for details	$\left\{ \begin{array}{l} \text{quantile}_{1-\alpha}(\text{Beta}(k, m+1-k)) \geq 1-\alpha \\ \text{iff of Beta/Bin distributions} \\ \Downarrow \\ \Leftrightarrow \text{Beta}_{k, m+1-k}(1-\alpha) \leq \delta \\ \Leftrightarrow \text{Bin}_{m, 1-\alpha}(k-1) \geq 1-\delta \end{array} \right.$	\uparrow feasible iff $\alpha \geq 1-\delta^{1/m}$ (e.g., if $\alpha \geq \frac{\ln(1/\delta)}{m}$)	$\left\{ \begin{array}{l} \text{Lemma (e.g., Wass 2012)} \\ \forall m \geq 1, \forall k \in \{1, \dots, m\}, \forall p \in [0, 1], \\ \text{Bin}_{m, p}(k-1) = 1 - \text{Beta}_{k, m+1-k}(p) \end{array} \right.$

We thus pick $k := 1 + \text{quantile}_{1-\delta}(\text{Bin}(m, 1-\alpha))$ and set

$$C(x) := \{y : s(x, y) \leq s(x)\}.$$

Then: $\mathbb{P}\left[\mathbb{P}(Y \in C(x) | \mathcal{D}_m) \geq 1-\alpha\right] \geq 1-\delta.$

Note that the choice of $k = \lceil (n+1)(1-\alpha) \rceil$ recommended by the split CP method only yields: (when $s(x, y)$ has a continuous cdf).

$$P(Y \in C(x) | \mathcal{D}_n) \sim \text{Beta}(k, n+1-k), \text{ whose expectation equals } \frac{k}{n+1} \approx 1-\alpha$$

Therefore, the coverage $P(Y \in C(x) | \mathcal{D}_n)$ is "often" below $1-\alpha$.

(but not much below, since, by Bernstein's inequality (see Appendix A for details),

$$P \left[P(Y \in C(x) | \mathcal{D}_n) \geq 1-\alpha - \sqrt{\frac{2\alpha \ln(1/\delta)}{n} - \frac{2\epsilon \ln(1/\delta)}{n}} \right] \geq 1-\delta$$

$\ll \alpha$ if $n \gg \frac{\ln(1/\delta)}{\alpha}$

3. Beyond binary losses

Two key references:

- RCE algorithm \rightarrow
- Bates et al. (2021), Distribution-free, risk-controlling prediction sets.
 - Park et al. (2021), PAC confidence sets for DNNs via calibrated prediction.

3.1. Setting: set-valued prediction with nested sets

Let

- loss function: $L(y, S) \geq 0$ be nonincreasing in $S \subset \mathcal{Y}$, for any $y \in \mathcal{Y}$
- $(C_\lambda(\cdot))_{\lambda \in \Lambda}$ be a family of set-valued predictors, with $\Lambda \subset \mathbb{R}$ closed and satisfying the nesting property: $\lambda_1 \leq \lambda_2 \Rightarrow \forall x \in \mathcal{X}, C_{\lambda_1}(x) \subset C_{\lambda_2}(x)$
- assume also: $\lim_{\lambda \rightarrow \sup(\Lambda)} L(y, C_\lambda(x)) = 0$ for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$.

Risk: Let (X, Y) be a random pair in $\mathcal{X} \times \mathcal{Y}$. We set: for all $\lambda \in [0, 1]$,

$$R(\lambda) := \mathbb{E}[L(Y, C_\lambda(X))] = \int L(y, C_\lambda(x)) dP_{X, Y}(x, y) < +\infty$$
assumption

Exc 1: confound prediction!

$$C_\lambda(x) = \{y \in \mathcal{Y} : s(x, y) \leq \lambda\} \quad \leftarrow \text{NB: any } (C_\lambda(x))_{\lambda \in \mathbb{R}} \text{ is of this form if } \lim_{\lambda \downarrow \lambda_0} C_\lambda(x) = C_{\lambda_0}(x)$$

$$L(y, S) = \mathbb{1}_{y \notin S}$$

$$R(\lambda) = \mathbb{E}[L(Y, C_\lambda(X))] = \mathbb{P}(Y \notin C_\lambda(X))$$

Exc 2: multiclass classification with class-weighting loss

$$\mathcal{Y} = \{1, \dots, K\}$$

For $S \subset \mathcal{Y}$, $L(y, S) = l_y \cdot \mathbb{1}_{y \notin S}$ (some mistakes are more serious than others)

$C_\lambda(x) \stackrel{\text{e.g.}}{=} \{ \text{all classes that a trained DNN predicts as likely (softmax} \geq 1-\lambda) \}$

Exc 3: multilabel classification

$x \in \mathcal{X}$, $y \subset \{1, \dots, K\}$ (multiple classes for object x)

$$\text{For } S \subset \{1, \dots, K\}, \quad L(y, S) = \frac{|y \setminus S|}{|y|} = 1 - \frac{|y \cap S|}{|y|}$$

"recall"

(NB: in this example, S and y live in the same space = $\mathcal{P}(\{1, \dots, K\})$)

Other examples: image segmentation (close to Ex.3)
 hierarchical classification
 protein structure prediction
 etc... as long as $(C_X(i))_{i \in \mathcal{R}}$ are nested

• $L(y, S)$ is nonincreasing in S

3.2. The RCPS algorithm (RCPS = Risk Controlling Prediction Sets)

Let $\alpha > 0$ and $\delta \in (0, 1)$ be two predefined "risk levels".

Let (X_i, Y_i) , $1 \leq i \leq m$, be m copies of (X, Y) . We set $\mathcal{D}_m := (X_i, Y_i)_{i=1}^m$

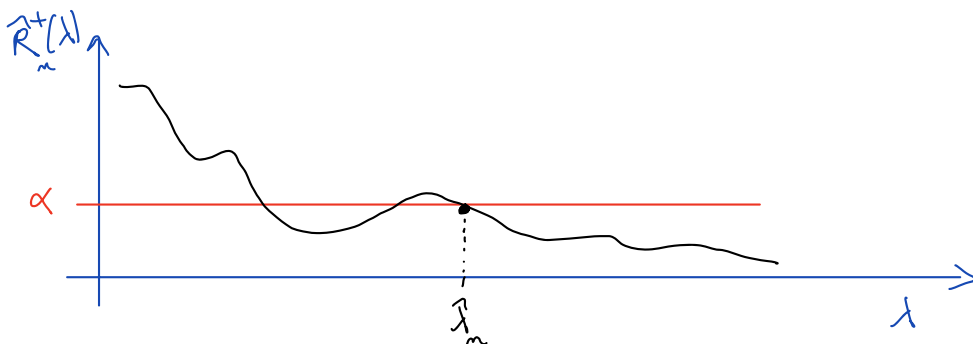
Assume $\lambda \mapsto \hat{R}^+(\lambda)$ is a random function (built from \mathcal{D}_m) such that:

$$\forall \lambda \in \Lambda, \mathbb{P}(R(\lambda) \leq \hat{R}^+(\lambda)) \geq 1 - \delta. \quad (1)$$

↑ pointwise high-probability bound

Then, choose $\hat{\lambda}_m$ from \mathcal{D}_m such that (assuming existence)

$$\forall \lambda' \geq \hat{\lambda}_m, \hat{R}_m^+(\lambda') \leq \alpha \quad (2)$$



(in some cases, $\hat{R}_m^+(\lambda)$ is nonincreasing)

Examples: we assume $L(y, S)$ bounded in, say, $[0, 1]$.

- Using Hoeffding's inequality and setting $\hat{R}_n(\lambda) = \frac{1}{n} \sum_{i=1}^n L(X_i, C_\lambda(X_i))$,

$$\mathbb{P}\left(R(\lambda) \leq \underbrace{\hat{R}_n(\lambda) + \sqrt{\frac{\ln(2/\delta)}{2n}}}_{\hat{R}_n^+(\lambda) \text{ satisfies (1)}}\right) \geq 1 - \delta$$

- If instead we use an "empirical Bernstein bound" (Maurer and Pontil '2009), we obtain, with $\hat{\sigma}^2(\lambda) = \frac{1}{n-1} \sum_{i=1}^n (L(X_i, C_\lambda(X_i)) - \hat{R}_n(\lambda))^2$

$$\mathbb{P}\left(R(\lambda) \leq \underbrace{\hat{R}_n(\lambda) + \hat{\sigma}(\lambda) \sqrt{\frac{2 \ln(2/\delta)}{n}} + \frac{7 \ln(2/\delta)}{3(n-1)}}_{\hat{R}_n^+(\lambda) \text{ satisfies (1)}}\right) \geq 1 - \delta$$

- Many other non-asymptotic concentration inequalities! (beyond bounded loss)
- Let us just mention a tight $\hat{R}_n^+(\lambda)$ when $L(y, S) = \mathbb{1}_{y \notin S} \in \{0, 1\}$:

$$\hat{R}_n^+(\lambda) = \sup \left\{ p \in [0, 1] : \underset{\substack{\text{cdf of} \\ \text{Bin}(n, p)}}{\text{Bin}_{n, p}}(n \hat{R}_n(\lambda)) \geq \delta \right\} \quad (3)$$

See, e.g., Langford (2005, Lemma 3.3) for a proof, or Appendix A below.

NB: In that case, RCPs reduce to Wilks' method! [see Section 3.3 below]

Theorem (Bates et al. 2021) Assume $(X_i, Y_i)_{1 \leq i \leq n}$ are i.i.d. copies of $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, and that $\lambda \mapsto \hat{R}_n^+(\lambda)$ satisfies (1) with $\alpha > 0$, $\delta \in (0, 1)$. Then, any $\hat{\Lambda}_n$ such that (2) holds satisfies

$$\mathbb{P}(R(\hat{\Lambda}_n) \leq \alpha) \geq 1 - \delta \quad [\mathbb{P} \text{ is w.r.t. random draw of } (X_i, Y_i)_{1 \leq i \leq n}]$$

If R is right-continuous, the conclusion also holds with $\hat{\Lambda}_n := \inf \{ \lambda \in \Lambda : \forall \lambda' \geq \lambda, \hat{R}_n^+(\lambda') \leq \alpha \}$

Proof: We can assume wlog that $\{\lambda \in \Lambda : R(\lambda) > \alpha\} \neq \emptyset$.

Define

$$\lambda^* := \sup \{\lambda \in \Lambda : R(\lambda) > \alpha\} < +\infty \text{ since } \lim_{\lambda \rightarrow \sup(\Lambda)} R(\lambda) = 0$$

We distinguish two cases:

Case 1: $R(\lambda^*) > \alpha$

By (1), we have, with probability $\geq 1 - \delta$,

$$\alpha < R(\lambda^*) \leq \widehat{R}_n^+(\lambda^*)$$

so that, by (2): $\widehat{\lambda}_n > \lambda^*$, which implies $R(\widehat{\lambda}_n) \leq \alpha$.
by def of λ^*

Case 2: $R(\lambda^*) \leq \alpha$

Let $\lambda \in \Lambda$ such that $\lambda < \lambda^*$. In particular: $R(\lambda) > \alpha$.

As before, with probability $\geq 1 - \delta$,

$$\alpha < R(\lambda) \leq \widehat{R}_n^+(\lambda)$$

so that $\widehat{\lambda}_n > \lambda$ on this event.

Therefore, $P(\widehat{\lambda}_n \geq \lambda^*) = \lim_{\lambda \nearrow \lambda^*} P(\widehat{\lambda}_n > \lambda) \geq 1 - \delta$.

In particular, $P(R(\widehat{\lambda}_n) \leq \alpha) \geq P(R(\widehat{\lambda}_n) \leq R(\lambda^*)) \geq 1 - \delta$

This concludes the first part of the theorem.

The second part follows as a consequence. ■

3.3. For Bernoulli losses, RCPS with (3) empotes Wilks' tolerance region!

We assume here that $L(y, S) = \mathbb{1}_{y \notin S}$, so that $\hat{R}_m(\lambda) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{x_i \in C_\lambda(x_i)}$

$\hat{R}_m^+(\lambda) := \sup \{ p: \text{Bin}_{m,p}(m \hat{R}_m(\lambda)) \geq \delta \}$ satisfies (1).
 ↑ by, e.g., Peng (2005, Lem 3.5)

We also assume that $\forall x \in \mathcal{X}, \forall \lambda_0 \in \Lambda, \lim_{\lambda \downarrow \lambda_0} C_\lambda(x) = C_{\lambda_0}(x)$ (often true)

Note that: (the case $m \hat{R}_m(\lambda) = m$ is treated separately)

$$\hat{R}_m^+(\lambda) \leq \alpha \Leftrightarrow \text{Bin}_{m, \alpha + \varepsilon}(m \hat{R}_m(\lambda)) < \delta \quad \forall \varepsilon > 0 \text{ small enough}$$

$$\Leftrightarrow \text{Bin}_{m, 1 - \alpha - \varepsilon}(m - m \hat{R}_m(\lambda) - 1) > 1 - \delta \quad \forall \varepsilon > 0 \text{ small enough}$$

$$\Leftrightarrow \text{Bin}_{m, 1 - \alpha}(m - m \hat{R}_m(\lambda) - 1) \geq 1 - \delta \quad (p \mapsto \text{Bin}_{m,p}(k) \text{ is continuous and decreasing in } p, \text{ for } k \leq m-1)$$

$$\Leftrightarrow m - m \hat{R}_m(\lambda) - 1 \geq \underbrace{\min \{ k \in \{0, \dots, m\} : \text{Bin}_{m, 1 - \alpha}(k) \geq 1 - \delta \}}_{= \text{quantile}_{1-\delta}(\text{Bin}(m, 1 - \alpha))}$$

$$\Leftrightarrow \sum_{i=1}^m \mathbb{1}_{x_i \in C_\lambda(x_i)} \geq 1 + \text{quantile}_{1-\delta}(\text{Bin}(m, 1 - \alpha)) =: k_{m, \alpha, \delta}$$

$$\Leftrightarrow \lambda \geq \lambda_{(k_{m, \alpha, \delta})} \quad \text{where } \begin{cases} \lambda_i := \inf \{ \lambda \in \Lambda : x_i \in C_\lambda(x_i) \} \\ \lambda_{(1)} \leq \lambda_{(2)} \leq \dots \leq \lambda_{(m)} \end{cases} \quad \text{assuming existence}$$

$$\text{Thus: } \underbrace{\inf \{ \lambda : \hat{R}_m^+(\lambda) \leq \alpha \}}_{\text{RCPS}} = \underbrace{\lambda_{(k_{m, \alpha, \delta})}}_{\text{tolerance region seen in Section 1.2.}}$$

Appendix A: Some useful facts about Binomial and Beta distributions

We write:

- $t \in [0,1] \mapsto \text{Beta}_{a,b}(t)$ for the cumulative distribution function (cdf) of the $\text{Beta}(a,b)$ distribution ($a, b > 0$).
- $k \in \{0,1,\dots,m\} \mapsto \text{Bin}_{m,p}(k)$ for the cdf of the binomial $\text{Bin}(m,p)$ distribution ($m \geq 1, p \in [0,1]$).

Lemma (e.g., Task 12, Lemma 3 in earlier version)

$\forall m \geq 1, \forall k \in \{1,\dots,m\}, \forall p \in [0,1]$,

$$\text{Bin}_{m,p}(k-1) = 1 - \text{Beta}_{k,m+1-k}(p)$$

Proof: let $U_1, \dots, U_m \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0,1])$, and set $U_{(1)} \leq \dots \leq U_{(m)}$.

We know that $X := \sum_{i=1}^m \mathbb{1}_{\{U_i \leq p\}} \sim \text{Bin}(m,p)$ and $U_{(k)} \sim \text{Beta}(k, m+1-k)$

Therefore,

$$\begin{aligned} \text{Bin}_{m,p}(k-1) &= 1 - \mathbb{P}(X \geq k) \\ &= 1 - \mathbb{P}(U_{(k)} \leq p) \\ &= 1 - \text{Beta}_{k,m+1-k}(p). \end{aligned}$$

↑
cf Page 2 with $\mathcal{U} = \text{Unif}([0,1])$

Fact: for $m \geq 1$ and $\alpha, \delta \in (0,1)$, we have:

$$\exists k \in \{1,\dots,m\} \text{ s.t. } \text{quantile}_\alpha(\text{Beta}(k, m+1-k)) \geq 1 - \alpha \quad (\text{i})$$

$$\Leftrightarrow \text{quantile}_\alpha(\text{Beta}(m, 1)) \geq 1 - \alpha \quad (\text{ii})$$

$$\Leftrightarrow \alpha \geq 1 - \delta^{1/m} \quad (\text{iii})$$

$$\Leftrightarrow \alpha \geq \frac{\ln(1/\delta)}{m} \quad (\text{iv})$$

Proof: we have a natural coupling of all distributions $\text{Beta}(k, m+1-k)$, $k \in \{1,\dots,m\}$, given by the order statistics $U_{(1)} \leq \dots \leq U_{(m)}$ with $U_1, \dots, U_m \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0,1])$

(recall that $U_{(k)} \sim \text{Beta}(k, m+1-k)$ for each $k \in \{1, \dots, m\}$)

Therefore,

$$\text{Beta}(m, 1) = \text{law}(U_{(m)}) \stackrel{\text{st}}{\geq} \text{law}(U_{(k)}) = \text{Beta}(k, m+1-k)$$

so that

$$\text{quantile}_\delta(\text{Beta}(m, 1)) \geq \text{quantile}_\delta(\text{Beta}(k, m+1-k))$$

This proves (i) \Leftrightarrow (ii).

Furthermore,

$$(ii) \Leftrightarrow \text{Beta}_{m,1}(1-\alpha) \leq \delta$$

$$\Leftrightarrow \mathbb{P}(U_{(m)} \leq 1-\alpha) \leq \delta$$

$$\Leftrightarrow \mathbb{P}\left(\max_{1 \leq i \leq m} U_i \leq 1-\alpha\right) \leq \delta$$

$$\Leftrightarrow (1-\alpha)^m \leq \delta$$

$$\Leftrightarrow \alpha \geq 1 - \delta^{1/m}$$

which proves (ii) \Leftrightarrow (iii).

Finally, note that

$$1 - \delta^{1/m} = 1 - \exp\left(-\frac{\ln(1/\delta)}{m}\right) \leq \frac{\ln(1/\delta)}{m}$$

so that (iv) \Rightarrow (iii). ■

Lemma (e.g., Lem 3.3. by Langford 2005): A tight high-probability bound on p

Let $m \geq 1$, $p^* \in [0, 1]$, and $Z \sim \text{Bin}(m, p^*)$

Then, the r.o.

$$\hat{p}_m^+ := \sup \left\{ p \in [0, 1] : \underbrace{\text{Bin}_{m,p}(Z)}_{\substack{\uparrow \\ \text{cdf of } \text{Bin}(m,p)}}} \geq \delta \right\}$$

satisfies:

$$\mathbb{P}(p^* \leq \hat{p}_m^+) \geq 1 - \delta.$$

Proof: since $Z \sim \text{Bin}_{m,p^*}$, we have $\underbrace{\text{Bin}_{m,p^*}(Z)}_{\substack{\text{sto} \\ \leftarrow \text{cdf of } \text{Bin}(m,p^*)}} \geq \text{Unif}([0,1])$

so that $\mathbb{P}(\text{Bin}_{m,p^*}(Z) \geq \delta) \geq 1 - \delta$

Noting that $\{\text{Bin}_{m,p^*}(Z) \geq \delta\} \subset \{p^* \leq \hat{p}_m^{1+\delta}\}$ concludes the proof. \square

Lemma (concentration of beta distribution) [See Vovk'2012 for a somewhat equivalent statement]

Let $m \geq 1$, $\alpha \in [\frac{1}{m+1}, 1)$, $k := \lceil (m+1)(1-\alpha) \rceil \in \{1, \dots, m\}$.

Then, for $X \sim \text{Beta}(k, m+1-k)$ and any $\delta \in (0, 1)$,

$$\mathbb{P}\left(X \geq 1 - \alpha - \sqrt{\frac{2\alpha \ln(1/\delta)}{m}} - \frac{2 \ln(1/\delta)}{m}\right) \geq 1 - \delta.$$

Proof (sketch): let $\varepsilon := \sqrt{\frac{2\alpha \ln(1/\delta)}{m}} + \frac{2 \ln(1/\delta)}{m}$, and assume that $\alpha + \varepsilon < 1$.

$$\begin{aligned} \mathbb{P}(X \geq 1 - \alpha - \varepsilon) &= 1 - \text{Beta}_{k, m+1-k}(1 - \alpha - \varepsilon) \\ &= \text{Bin}_{m, 1 - \alpha - \varepsilon}(k-1) \quad \text{by a lemma above} \\ &= \mathbb{P}(S \leq k-1) \quad \text{for } S \sim \text{Bin}(m, 1 - \alpha - \varepsilon) \\ &= 1 - \mathbb{P}(S \geq k) \end{aligned}$$

But,

$$\begin{aligned} \mathbb{P}(S \geq k) &= \mathbb{P}(m - S \leq m - k) \\ &= \text{Bin}_{m, \alpha + \varepsilon}(m - k) \\ &\leq \delta, \end{aligned}$$

where the last inequality follows from $\text{Bin}_{m, \alpha + \varepsilon}(m - k) \leq \text{Bin}_{m, p}(m - k)$ with $p := \frac{m - k}{m} + \sqrt{\frac{2 \frac{m - k}{m} \ln(1/\delta)}{m}} + \frac{2 \ln(1/\delta)}{m} < \alpha + \varepsilon$ (by the choice of $k = \lceil (m+1)(1-\alpha) \rceil$) and from $\text{Bin}_{m, p}(m - k) \leq \delta$ by Bernstein's inequality. \square

References [author (year), title]

- Wilks (1941), Determination of sample sizes for setting tolerance limits.
- Wald (1943), An extension of Wilks' method for setting tolerance limits.
- Scheffé and Zidek (1945), Non-parametric estimation I. Validation of order statistics.
- Zidek (1947), Non-parametric estimation II. Historically equivalent blocks and tolerance regions - the continuous case.
- Wash (2012), Conditional validity of inductive conformal predictors.
- Bates et al. (2021), Distribution-free, risk-controlling prediction sets.
- Park et al. (2021), PAC confidence sets for DNNs via calibrated prediction.
- Mauer and Portil (2009), Empirical Bernstein bounds and sample variance regularization.
- Langford (2005), Tutorial on practical prediction theory for classification.