# Federated Conformal Prediction General

Min, Xia

May 6, 2024

## 1 Conformal Prediction General[1]

**Definition 1.1** (Exchangeability). *[3] For any r.v. $x_1, \cdots, x_k$, we say they are exchangeable if for any permutation $\sigma : [k] \to [k]$(bijection), $(x_1, \cdots, x_k) \overset{d.}{=} (x_{\sigma(1)}, \cdots, x_{\sigma(k)})$.*

**Definition 1.2** (Weighted Exchangeability). *[4] For any r.v. $x_1, \cdots, x_k$, we say they are weighted exchangeable if their joint density canbe factorized as*

$$f(x_1, \cdots, x_k) = \prod_{i=1}^{k} w_i(x_i) \cdot g(x_1, \cdots, x_k),$$

*where g is exchangeable, i.e., $g(x_1, \cdots, x_k) = g(x_{\sigma(1)}, \cdots, x_{\sigma(k)})$.*

For conformal prediction two classes of targets are studied.

**Definition 1.3** (Marginal Coverage). *$(X, Y) \in \mathbb{R}^p \times \mathbb{R} \sim P_{XY}$ which is unknown. Given training set $Tr = \{(X_i, Y_i)\}_{i=1}^{n}$, and test on $(X_{n+1}, Y_{n+1})$, both i.i.d.*

*$C_\alpha$ satisfies distribution-free marginal coverage at level $1 - \alpha$ if*

$$P(Y_{n+1} \in C_\alpha(X_{n+1})) \geq 1 - \alpha, \ \forall P_{XY}$$

*The probability is with respect to $\{(X_i, Y_i)\}_{i=1}^{n+1}$.*

**Definition 1.4** (Conditional Coverage). *$(X, Y) \in \mathbb{R}^p \times \mathbb{R} \sim P_{XY}$ which is unknown. Given training set $Tr = \{(X_i, Y_i)\}_{i=1}^{n}$, and test on $(X_{n+1}, Y_{n+1})$, both i.i.d.*

*$C_\alpha$ satisfies distribution-free marginal coverage at level $1 - \alpha$ if*

$$P\left(Y_{n+1} \in C_\alpha(X_{n+1}) \middle| X_{n+1} = x\right) \geq 1 - \alpha, \ \forall P_{XY}$$

*The probability is with respect to $\{(X_i, Y_i)\}_{i=1}^{n}$ and $Y_{n+1}$.*

# 2   Standard Split Conformal Prediction

- First divide training set $D$ into two sets: $D_1$ for proper training set and $D_2$ for calibration set. And let $n_i = |D_i|$, fit point predictor $\hat{f}_1$ on $D_1$.

- Calculate residuals on $D_2$: $R_i = \left| Y_i - \hat{f}_1(X_i) \right|$, $i \in D_2$.

- Find quantile on calibration residuals: $\hat{q}_2 = \lceil (1-\alpha)(n_2+1) \rceil$ smallest of $R_i$, $i \in D_2$.

- Construct a conformal set: $C_\alpha(x) = \left[ \hat{f}_1(x) - \hat{q}_2, \hat{f}_1(x) + \hat{q}_2 \right]$.

Let $R_{n+1} = \left| Y_{n+1} - \hat{f}_1(X_{n+1}) \right|$. Let rank statistic $R_{(j)}$ be the $j$-th smallest in $R_i$, $i \in D_2$, and $k_\alpha = \lceil (1-\alpha)(n_2+1) \rceil$. As

$$\{ Y_{n+1} \in C_\alpha(X_{n+1}) \} = \{ R_{n+1} \le \hat{q}_2 \} = \left\{ R_{n+1} \le R_{(k_\alpha)} \right\},$$

and $R_i$, $i \in D_2$, $R_{n+1}$ are exchangeable, we have

$$\mathbb{P} \left( Y_{n+1} \in C_\alpha(X_{n+1}) \Big| D_1 \right) \in \left[ 1 - \alpha, 1 - \alpha + \frac{1}{n_2 + 1} \right).$$

Assume a more general score function $V(x, y) = V((x, y); \hat{f}_1)$, define $R_i = V(X_i, Y_i)$ and change the conformal set to

$$C_\alpha(x) = \left\{ y : V(x, y) \le R_{(k_\alpha)} \right\}.$$

**Remark 2.1.** *Further condition on calibration set, which means conditioning on entire training set $D$ and assume $R = V(x, y)$ has distribution $F$. As*

$$\{ Y_{n+1} \in C_\alpha(X_{n+1}) \} = \left\{ R_{n+1} \le R_{(k_\alpha)} \right\},$$

*Assume the distribution function of $R_{(j)}$ is $F_{(j)}$, and we have*

$$\mathbb{P} \left( \mathbb{P} \left( Y_{n+1} \in C_\alpha(X_{n+1}) \Big| D \right) \le t \right) = \mathbb{P} \left( \mathbb{P} \left( R_{n+1} \le R_{(k_\alpha)} \Big| D \right) \le t \right)$$

*condition on $D$ randomness comes from $R_{n+1}$,* $= \mathbb{P} \left( F(R_{(k_\alpha)}) \le t \right)$

$$= \mathbb{P} \left( R_{(k_\alpha)} \le F^{-1}(t) \right)$$

$$= F_{(k_\alpha)}(F^{-1}(t)) \tag{1}$$

*rank statistic has density $F'_{(j)}(x) = jC^j_{n_2}x^{j-1}(1-x)^{n-j}f(x)$, thus take derivative on formula*
*(1), and $\mathbb{P}\left(Y_{n+1} \in C_\alpha(X_{n+1})\Big|D\right)$ has density*

$$k_\alpha C^{k_\alpha}_{n_2}t^{k_\alpha-1}(1-t)^{n-k_\alpha}.$$

# 3   Standard Full Conformal Prediction

Full CP has similar steps as split CP. It uses all data points for training.

- Fix any $x$ and trial data $y$ to construct training set $\{(X_1, Y_1), \cdots, (X_n, Y_n), (x, y)\}$.

- Train point predictor $\hat{f}$ on training set and define residuals $R_i = \left|Y_i - \hat{f}(X_i)\right|$, $i \in$ $[n]$, $R_{n+1} = \left|y - \hat{f}(x)\right|$.

- Define $j$-th rank statistic of $R_i$, $i \in [n]$ as $R_{(j)}$, $k_\alpha = \lceil(1-\alpha)(n_2+1)\rceil$, and conformal set

$$C_\alpha(x) = \left\{y : R_{n+1} \leq R_{(k_\alpha)}\right\}.$$

As $\{Y_{n+1} \in C_\alpha(X_{n+1})\} = \left\{R_{n+1} \leq R_{(k_\alpha)}\right\}$, and the exchangebility of data

$$\mathbb{P}\left(Y_{n+1} \in C_\alpha(X_{n+1})\right) \in \left[1 - \alpha, 1 - \alpha + \frac{1}{n+1}\right).$$

# 4   Standard CP under covariate shift

Follow the procedure of split CP and heterogeneity between training and test data[4]. Assume

$$Z_i = (X_i, Y_i) \sim P = P_X \times P(Y|X), i = 1, \cdots, n,$$
$$Z_{n+1} = (X_{n+1}, Y_{n+1}) \sim P' = P'_X \times P_{Y|X}.$$

- Fix any trial data $y$ to construct training set $\{(X_1, Y_1), \cdots, (X_n, Y_n), (X_{n+1}, y)\}$. Train point predictor $\hat{f}$ on new training set.

- Calculate nonconformity scores $R_i = V(X_i, Y_i)$, $i \in \{1, \cdots, n\}$, $R_{n+1} = V(X_{n+1}, y)$ based on $\hat{f}$.

- Calculate importance weights $p_i$ based on likelihood ratio $w$:

$$w(x) = \frac{dP'_X(x)}{dP_X(x)},$$

$$p_i = \frac{w(X_i)}{\sum\limits_{j=1}^{n+1} w(X_j)}, \ i = 1, \cdots, n+1.$$

- Calculate $1 - \alpha$ quantile of distribution $\sum\limits_{i=1}^{n} p_i \delta_{R_i} + p_{n+1} \delta_\infty$ as $q_\alpha$. Define conformal set $C_\alpha(x) = \{y : R_{n+1} \leq q_\alpha\}$.

All independent variables are weighted exchangeable. Let $E_Z$ be $\{Z_1, \cdots, Z_{n+1}\} = \{z_1, \cdots, z_{n+1}\}$. Assume joint density is $f(z_1, \cdots, z_{n+1}) = \prod\limits_{i=1}^{n+1} dP(z_i) \cdot w(x_{n+1})$. Condition on $E_Z$, calculate $R_i$ based on $\hat{f}$ and $z_i$, for all permutation $\sigma$

$$\mathbb{P}\left(R_{n+1} = r_i \Big| E_Z\right) = \mathbb{P}\left(Z_{n+1} = z_i \Big| E_Z\right) = \frac{\sum\limits_{\sigma(n+1)=i} f(z_{\sigma(1)}, \cdots, z_{\sigma(n+1)})}{\sum\limits_{\sigma} f(z_{\sigma(1)}, \cdots, z_{\sigma(n+1)})} = p_i,$$

which leads to $R_{n+1} \Big| E_Z \sim \sum\limits_{i=1}^{n+1} p_i \delta_{r_i}$. Let $Q(1 - \alpha, F)$ be the quantile function,

$$\mathbb{P}\left(R_{n+1} \leq Q(1 - \alpha, \sum\limits_{i=1}^{n+1} p_i \delta_{r_i}) \Big| E_Z\right) \geq 1 - \alpha,$$

means

$$\mathbb{P}\left(R_{n+1} \leq Q(1 - \alpha, \sum\limits_{i=1}^{n} p_i \delta_{r_i} + p_{n+1} \delta_\infty) \Big| E_Z\right) \geq 1 - \alpha,$$

as condition on $E_Z$, $\sum\limits_{i=1}^{n} p_i \delta_{r_i} + p_{n+1} \delta_\infty = \sum\limits_{i=1}^{n} p_i \delta_{R_i} + p_{n+1} \delta_\infty$ (left $p$ is based on $z$ and right based on $Z$). The $p_i$ in following formula is different from previous one.

$$\mathbb{P}\left(R_{n+1} \leq Q(1 - \alpha, \sum\limits_{i=1}^{n} p_i \delta_{R_i} + p_{n+1} \delta_\infty) \Big| E_Z\right) \geq 1 - \alpha,$$

thus taking expectation on all $E_Z$,

$$\mathbb{P}\left(Y_{n+1} \in C_\alpha(X_{n+1})\right) = \mathbb{P}\left(R_{n+1} \leq q_\alpha\right) \geq 1 - \alpha$$

# 5   Federated Conformal Prediction Article1

Efficient Conformal Prediction under Data Heterogeneity[2]

Idea: The marginal coverage is measured over all training data and test points. However, if there is a high variability in the coverage probability as a function of the training data, the test coverage probability may be substantially below $1 - \alpha$ for a particular training set.

**Definition 5.1** (empirical miscoverage rate). $\alpha(Tr) = P(Y_{n+1} \notin C_\alpha(X_{n+1}) \big| Tr)$

In this article, assume $n$ agents each has calibration data $(X_k^i, Y_k^i) \sim P_X^i P_{Y|X}$, $k = 1, \cdots, n^i$, $i = 1, \cdots, n$, and calibration set $D_i = \{(X_k^i, Y_k^i)\}_{k=1}^{n_i}$, $i = 1, \cdots, n$. Let calibration distribution be $P^{cal} = \sum_{i=1}^{n} \pi_i P_X^i P_{Y|X}$, where $\pi_i = n_i / \left( \sum_{j=1}^{n} n_j \right)$, and the test distribution $P^{test} = P_X^{n+1} P_{Y|X}$. Let the general density ratio be $w(x,y) = \dfrac{dP_X^{n+1}(x)}{\sum_{i=1}^{n} \pi_i dP_X^i(x)}$.

- Utilize the GMM to compute parameters $\{\pi_y^i, \mu_y^i, \Sigma_y^i\}_{y \in \mathcal{Y}^i}$ on $D_i$. Note that $P_X^i$ is approximated by $|\mathcal{Y}^i|$ centers mixed GMM, $P_X^i = \sum_{y \in \mathcal{Y}^i} \pi_y^i N(\phi(x); \mu_y^i, \Sigma_y^i)$, $i = 1, \cdots, n+1$, where $\phi()$ be some latent map used while training $\hat{f}$. Further $w(x,y)$ can be calculated.

- Fix any trial data $y$, similar to covariate shift setting, a common idea should be calculate importance weight $p_k^i = w(X_k^i, Y_k^i)/W$, $k = 1, \cdots, n^i$, $i = 1, \cdots, n$, where $W = \sum_{i=1}^{n} \sum_{k=1}^{n^i} w(X_k^i, Y_k^i) + w(X_{n+1}, y)$, and $p_{n+1} = w(X_{n+1}, y)/W$. Similarly define residuals $R_k^i$, $R_{n+1}$.

- Conformal set: $C_\alpha(X_{n+1}) = \left\{ y : R_{n+1} \leq Q(1 - \alpha, \sum_{i=1}^{n} \sum_{k=1}^{n^i} p_k^i \delta_{R_k^i} + p_{n+1} \delta_\infty) \right\}$.

# References

[1] Anastasios N Angelopoulos, Stephen Bates, et al. Conformal prediction: A gentle introduction. *Foundations and Trends® in Machine Learning*, 16(4):494–591, 2023.

[2] Vincent Plassier, Nikita Kotelevskii, Aleksandr Rubashevskii, Fedor Noskov, Maksim Velikanov, Alexander Fishkov, Samuel Horvath, Martin Takac, Eric Moulines, and Maxim Panov. Efficient conformal prediction under data heterogeneity. In *International Conference on Artificial Intelligence and Statistics*, pages 4879–4887. PMLR, 2024.

[3] Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3), 2008.

[4] Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas. Conformal prediction under covariate shift. *Advances in neural information processing systems*, 32, 2019.