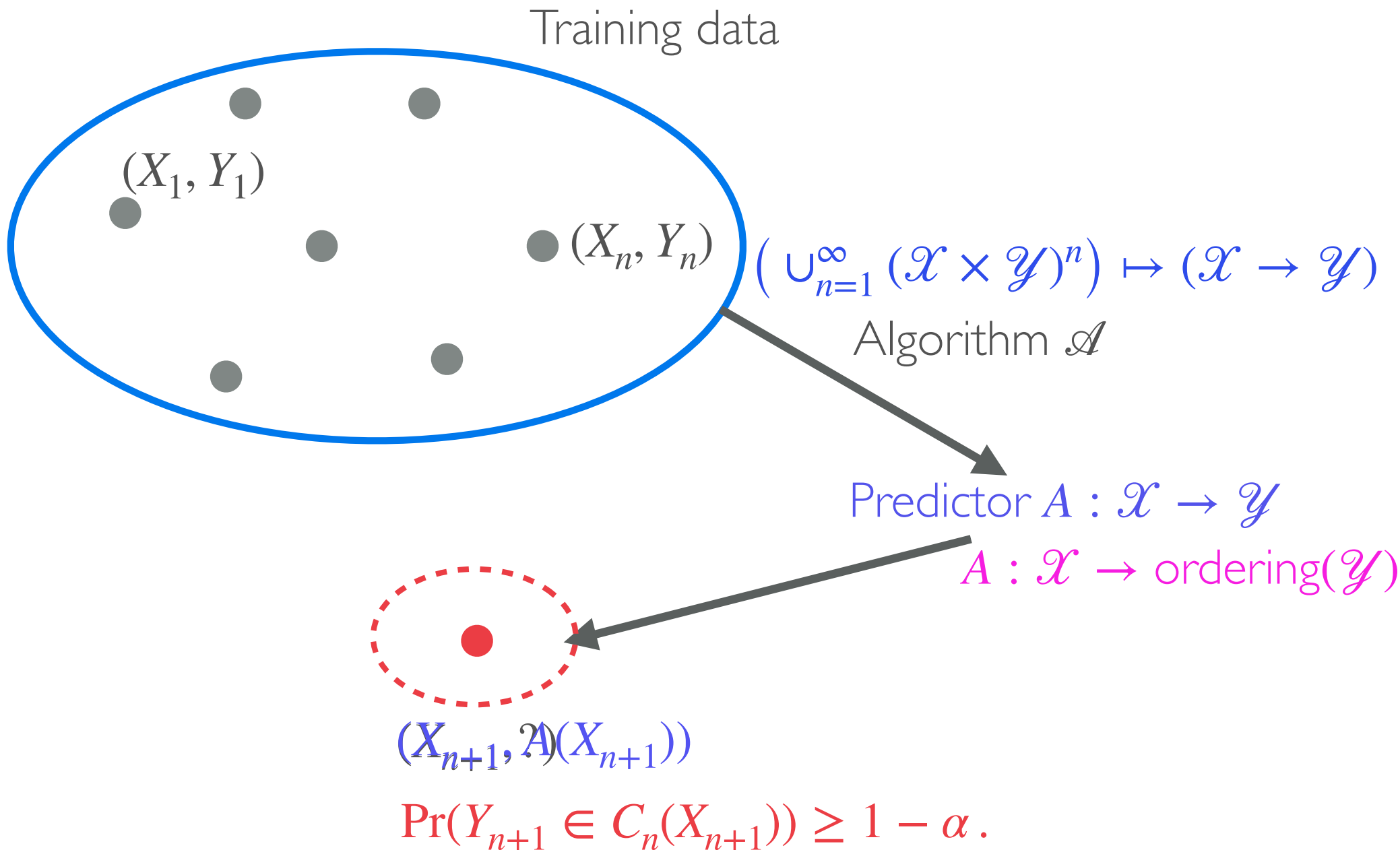# Uncertainty quantification
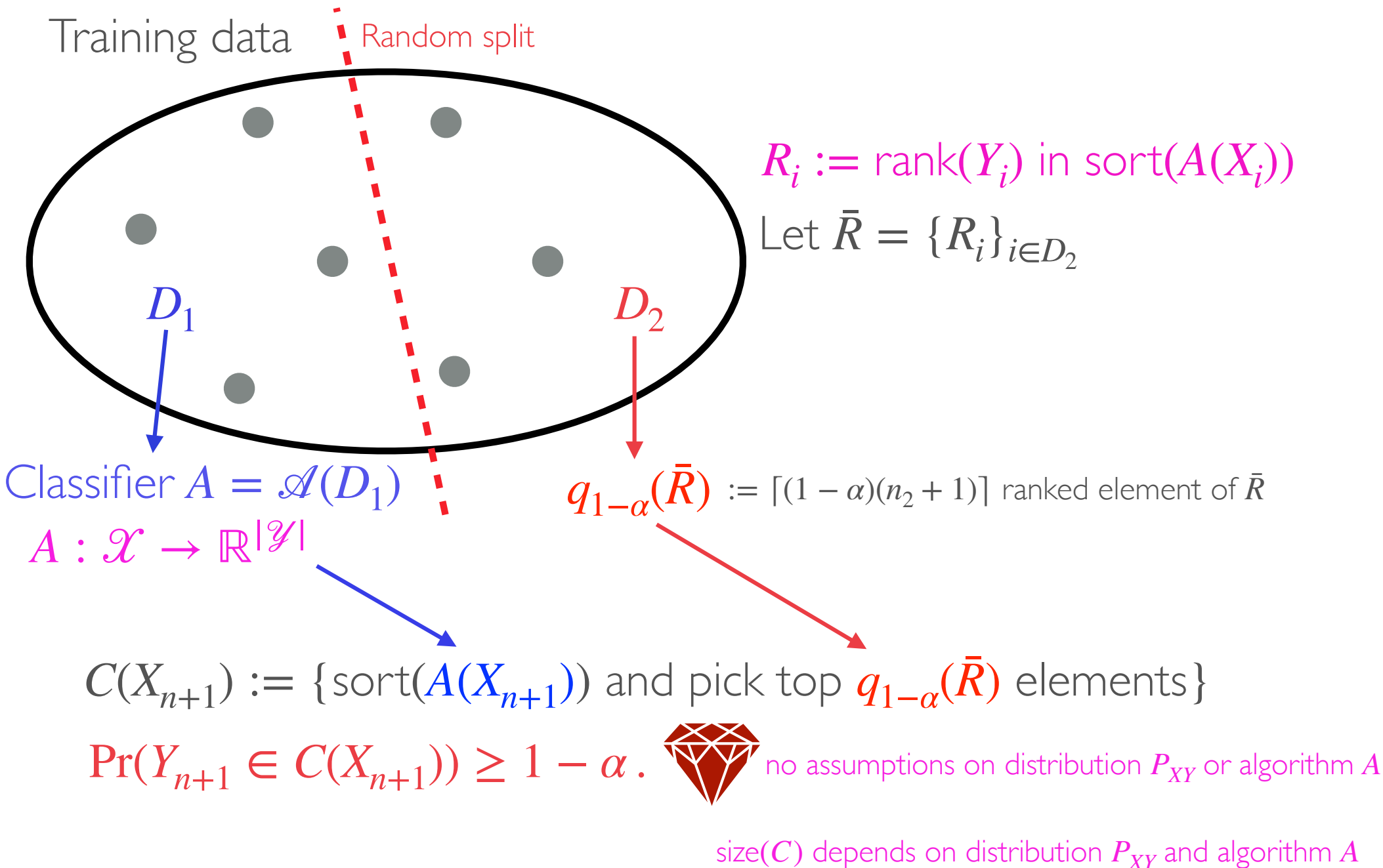# for black-box classifiers
# without distributional assumptions

Aaditya Ramdas

Dept. of Statistics and Data Science
Machine Learning Dept.
Carnegie Mellon University

# Prediction vs "Predictive Inference"

Training data



$(X_1, Y_1)$

$(X_n, Y_n)$

$\left( \cup_{n=1}^{\infty} (\mathscr{X} \times \mathscr{Y})^n \right) \mapsto (\mathscr{X} \to \mathscr{Y})$

Algorithm $\mathscr{A}$

Predictor $A : \mathscr{X} \to \mathscr{Y}$

$A : \mathscr{X} \to \text{ordering}(\mathscr{Y})$

$(X_{n+1}, A(X_{n+1}))$

$\Pr(Y_{n+1} \in C_n(X_{n+1})) \geq 1 - \alpha.$

# Split Conformal Prediction for classification

Training data     Random split



$R_i := \text{rank}(Y_i) \text{ in sort}(A(X_i))$

Let $\bar{R} = \{R_i\}_{i \in D_2}$

$D_1$        $D_2$

Classifier $A = \mathcal{A}(D_1)$

$q_{1-\alpha}(\bar{R}) := \lceil (1 - \alpha)(n_2 + 1) \rceil$ ranked element of $\bar{R}$

$A : \mathcal{X} \to \mathbb{R}^{|\mathcal{Y}|}$

$C(X_{n+1}) := \{\text{sort}(A(X_{n+1})) \text{ and pick top } q_{1-\alpha}(\bar{R}) \text{ elements}\}$

$\Pr(Y_{n+1} \in C(X_{n+1})) \geq 1 - \alpha .$     no assumptions on distribution $P_{XY}$ or algorithm $A$

size$(C)$ depends on distribution $P_{XY}$ and algorithm $A$

# Better residual for probabilistic classifiers

$$A : \mathscr{X} \to \Delta^{|\mathscr{Y}|}$$

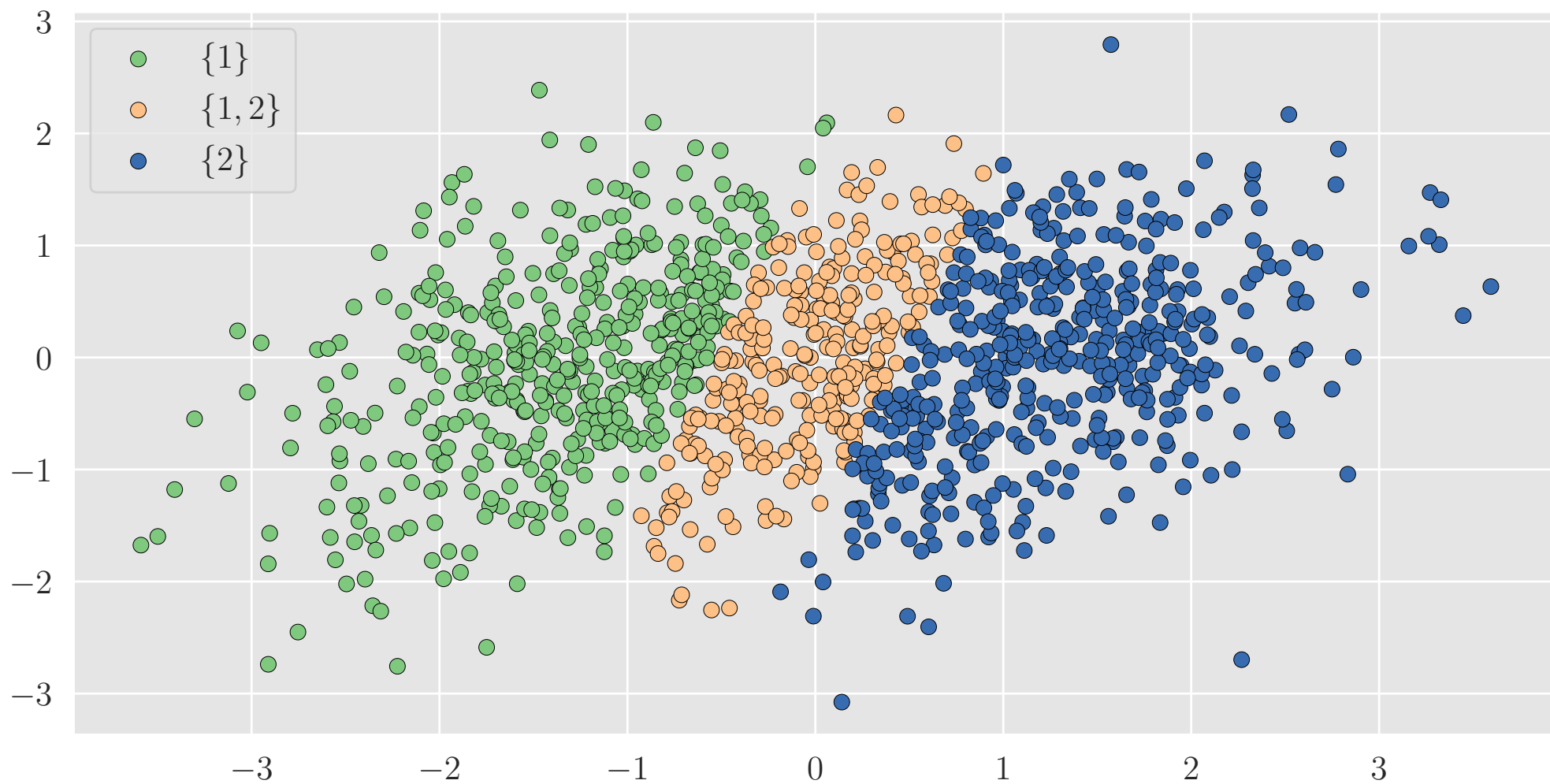$$R_i := \sum_{\ell \neq Y_i : [A(X_i)]_\ell > [A(X_i)]_{Y_i}} [A(X_i)]_\ell$$

= total prob. of more likely, wrong labels
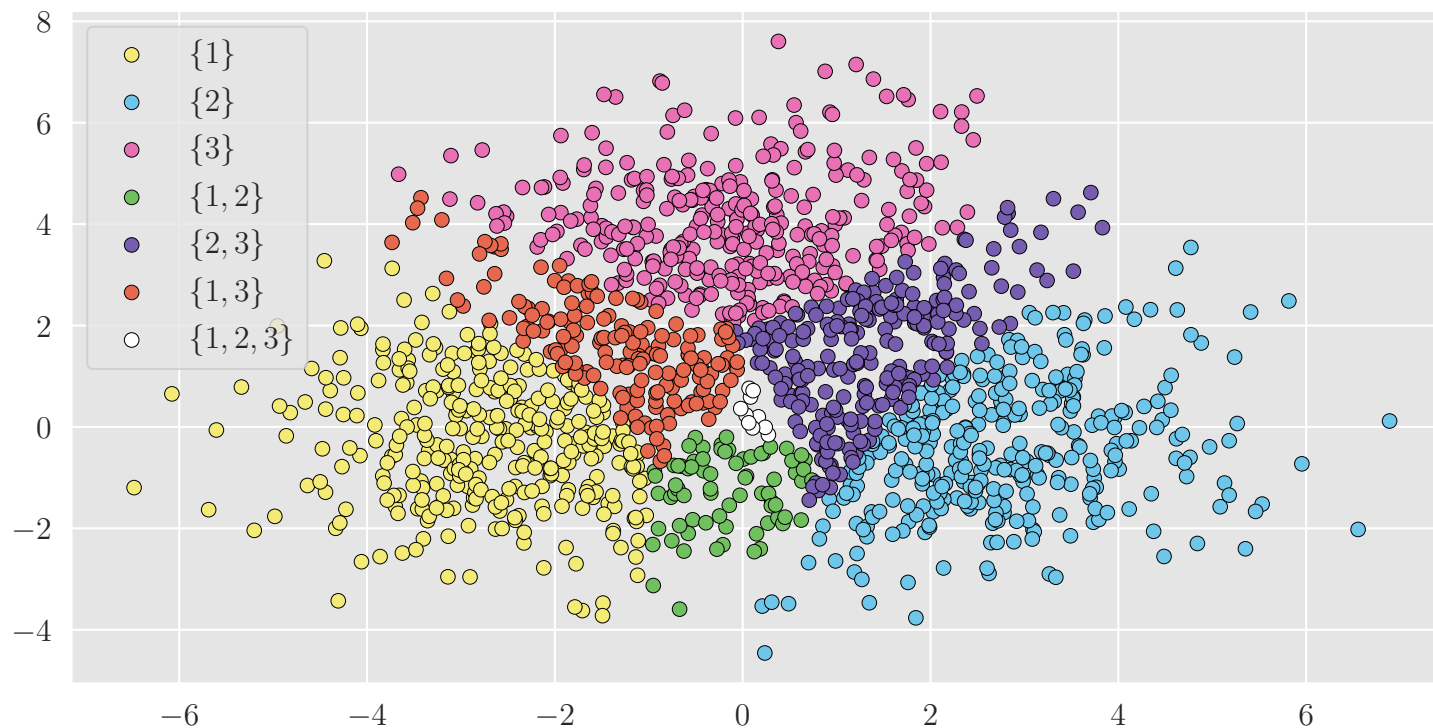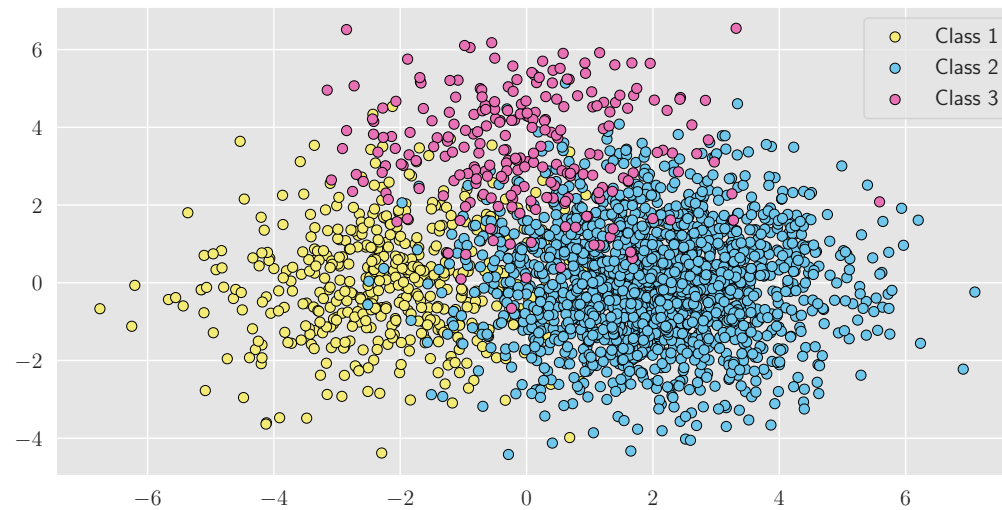
Let $\bar{R} = \{R_i\}_{i \in D_2}$

$C(X_{n+1}) := \{$least number of labels whose total mass $\geq q_{1-\alpha}(\bar{R})\}$

$= \{$first $k$ labels in sort$(A(X_{n+1}))$ with cumulative prob. $\geq q_{1-\alpha}(\bar{R})\}$

(Can be "smoothed" out with randomization)

(Mixture of two Gaussians in two dimensions)

(Mixture of three Gaussians in two dimensions)

# Part 2: Calibrated probabilities

# Calibration in the binary setting

A function $f : \mathcal{X} \to [0,1]$ returns calibrated probabilities if

$$\mathbb{E}[Y_{n+1} \mid f(X_{n+1})] = f(X_{n+1})$$

Eg: Suppose we predict $f(X_{n+1}) \approx 0.3$ for 100 points,
then $\approx 30$ of those will have label one, and the rest label zero.

Fact: if $f$ is calibrated, then $f(X) = \mathbb{E}[Y \mid g(X)]$ for some $g$.
Reality: exact calibration is impossible with a finite data of size $n$.

We say that $f_n : \mathcal{X} \to [0,1]$ is distribution-free $(\epsilon_n, \alpha)$-calibrated if

$$\forall P_{XY}, \ \Pr(\,|\mathbb{E}[Y_{n+1} \mid f_n(X_{n+1})] - f_n(X_{n+1})| > \epsilon_n) \le \alpha,$$

and asymptotically calibrated if $\epsilon_n \to 0$.

Theorem: Asymptotic distribution-free calibration is impossible
if $\lim_{n \to \infty} \text{Range}(f_n)$ is uncountable.

# **Split** Binning
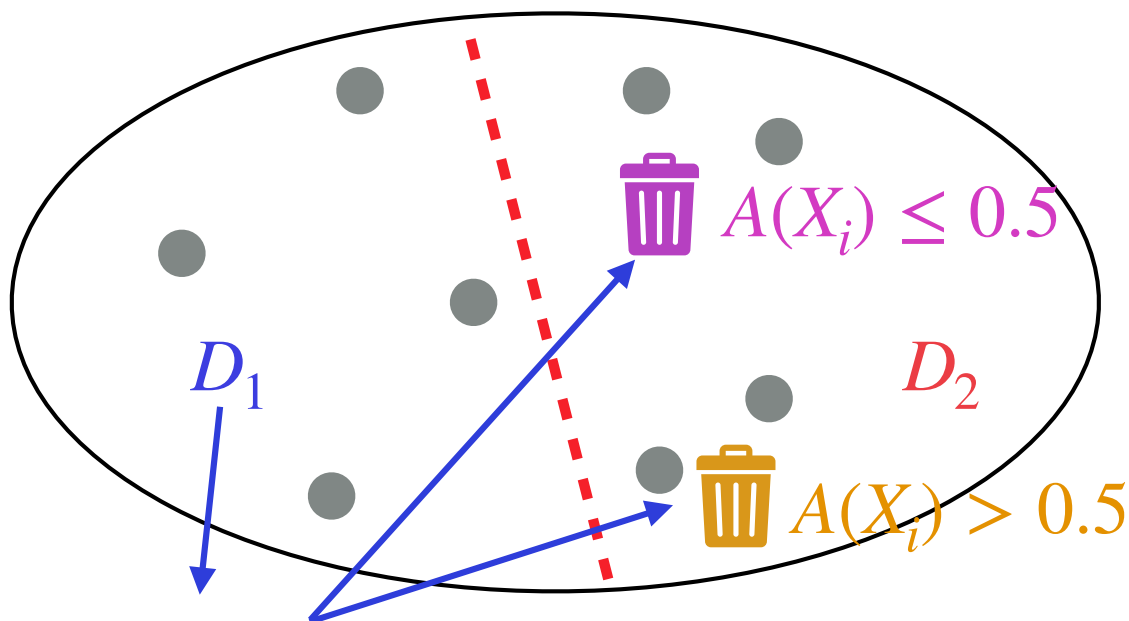
Training data

Random split



Improve to "uniform mass binning"

$A(X_i) \leq 0.5$

$\rightarrow$ smaller half of $A(X_i)$

$D_1$

$D_2$

$A(X_i) > 0.5$

$\rightarrow$ larger half of $A(X_i)$

$A = \mathscr{A}(D_1) : \mathscr{X} \rightarrow [0,1]$

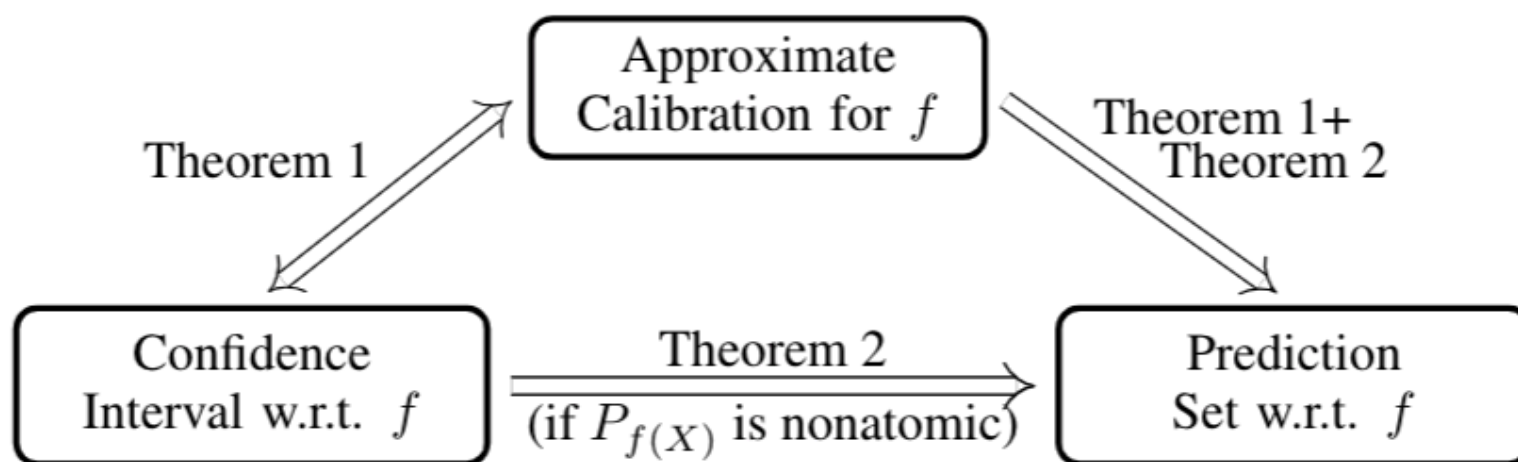$A(X_{n+1}) = $ 🗑  $\qquad f(X_{n+1}) = \dfrac{\sum_{i \in 🗑} Y_i}{|🗑|}$

(generalize to any number of bins)

$$\Pr\left( \,|\mathbb{E}[Y_{n+1} \,|\, f(X_{n+1}] - f(X_{n+1})| \leq c\hat{\sigma}\sqrt{\dfrac{\ln(1/\alpha)}{n}} \right) \geq 1 - \alpha\,.$$

no assumptions on distribution $P_{XY}$ or algorithm $A$

# Distribution-free binary classification: prediction sets, confidence intervals and calibration

Chirag Gupta[*,1], Aleksandr Podkopaev[*,1,2], Aaditya Ramdas[1,2]

**Theorem 3.** *Let $\alpha \in (0, 0.5)$ be a fixed threshold. If a sequence of scoring functions $\{f_n\}_{n \in \mathbb{N}}$ is asymptotically calibrated at level $\alpha$ for every distribution $P$ then*

$$\limsup_{n \to \infty} |\mathcal{X}^{(f_n)}| \leqslant \aleph_0.$$

## 4.3   Distribution-free calibration in the online setting

## 4.4   Calibration under covariate shift

# Sharpness?

One cannot guarantee sharpness without distributional assumptions.



$|\mathcal{Y}| = 2$

Number of bins, properties of $P_{XY}$ and quality of original classifier, all together determine sharpness, but not calibration.

Eg: consider the setting where $P(Y = 1|X) = 0.5$, i.e. $Y \perp X$. No classifier can be sharp, and not all classifiers are calibrated.