

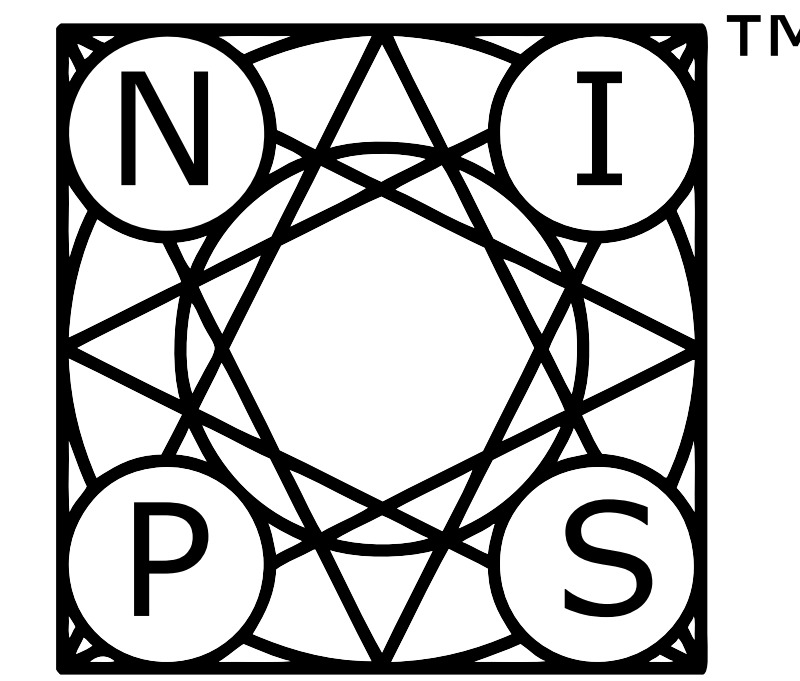


ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Distributionally Robust Logistic Regression

Soroosh Shafieezadeh-Abadeh, Peyman Mohajerin Esfahani, and Daniel Kuhn

École Polytechnique Fédérale de Lausanne (EPFL)



Abstract

Motivation: Classification problems face the following challenges

- Data-generating distribution \mathbf{P} is unknown
- Overfitting when # of training samples N is small
- Ad hoc regularization techniques lack theoretical justification

Our solution:

- Use **Distributionally Robust Optimization (DRO)**
- Use the **Wasserstein distance** to construct ambiguity sets

Logistic Regression (LR)

Assumption:

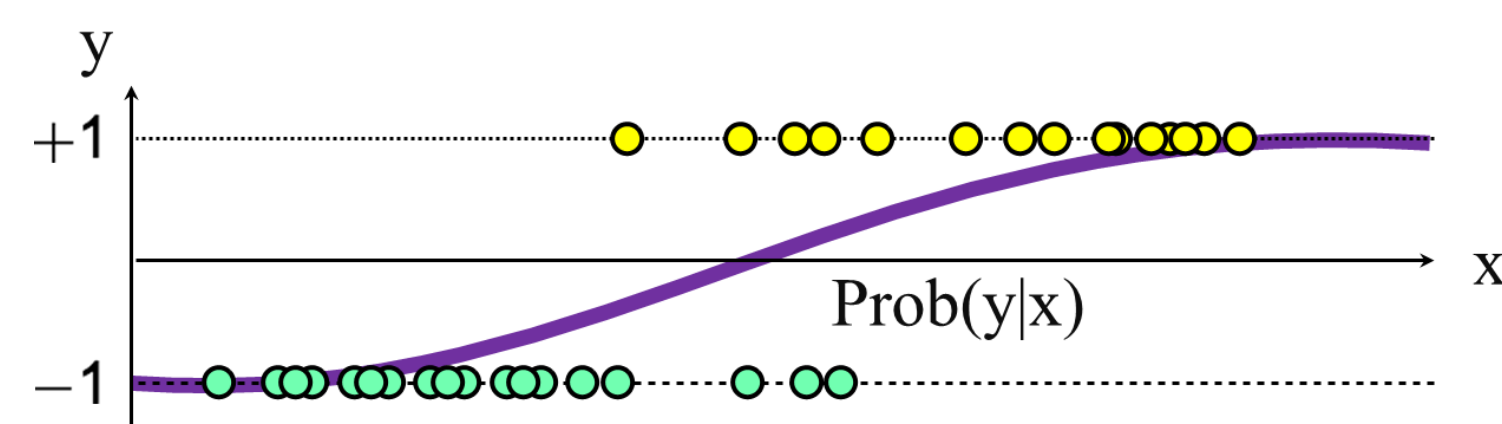
$$\text{Prob}(y|x) = [1 + \exp(-y\langle\beta, x\rangle)]^{-1}$$

Maximum likelihood estimation:

$$\text{LR: } \min_{\beta} \frac{1}{N} \sum_{i=1}^N l_{\beta}(\hat{x}_i, \hat{y}_i)$$

Logloss function:

$$l_{\beta}(x, y) = \log(1 + \exp(-y\langle\beta, x\rangle))$$



Statistical Learning

The ideal classifier would be a solution of

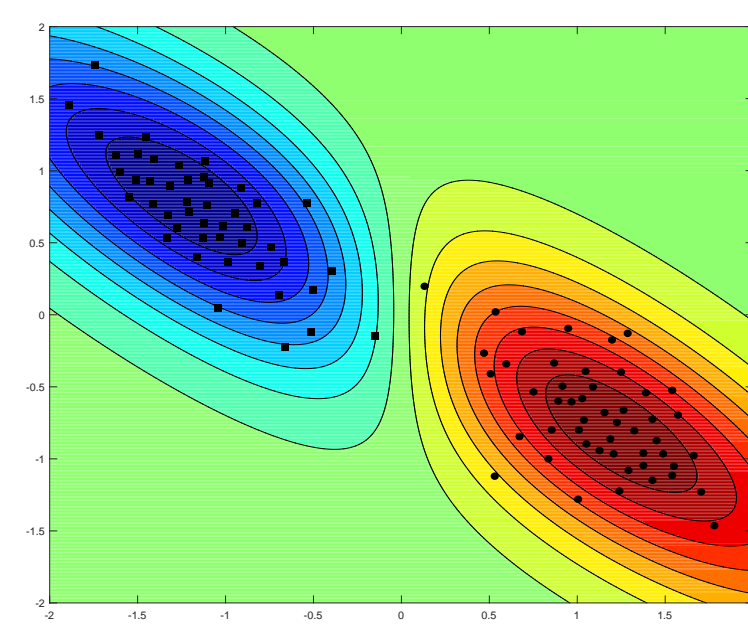
$$\min_{\beta} \mathbb{E}^{\mathbf{P}}[l_{\beta}(x, y)]$$

Challenges:

- \mathbf{P} is unknown
- Integration is hard even if \mathbf{P} is uniform on a box

Remedy: replace \mathbf{P} with the empirical distribution $\hat{\mathbf{P}}_N$

$$\text{LR} \Leftrightarrow \min_{\beta} \mathbb{E}^{\hat{\mathbf{P}}_N}[l_{\beta}(x, y)]$$



Regularized LR (RLR)

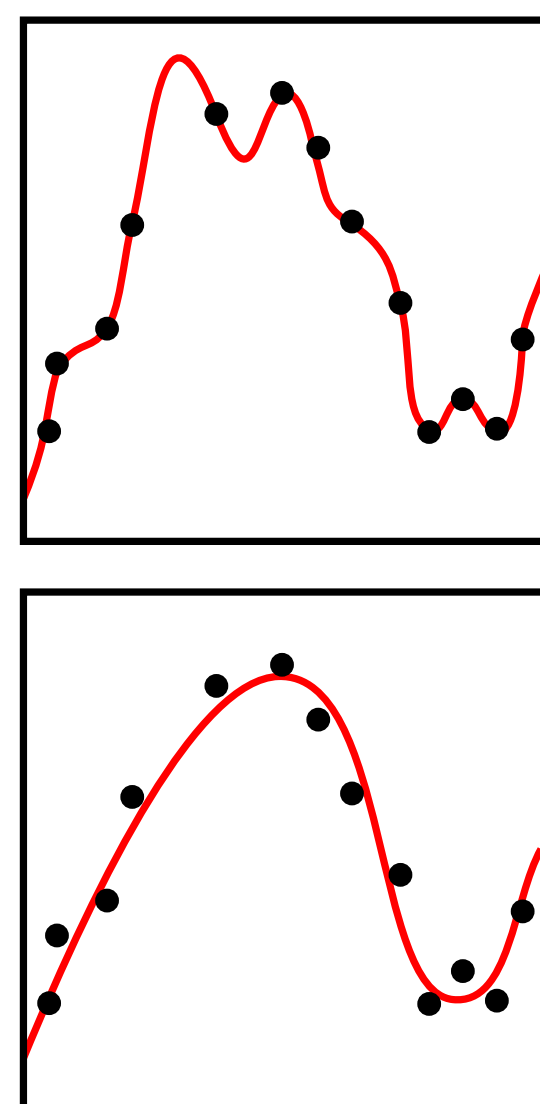
Challenge: LR suffers from overfitting

Remedy: *ad hoc* regularization techniques

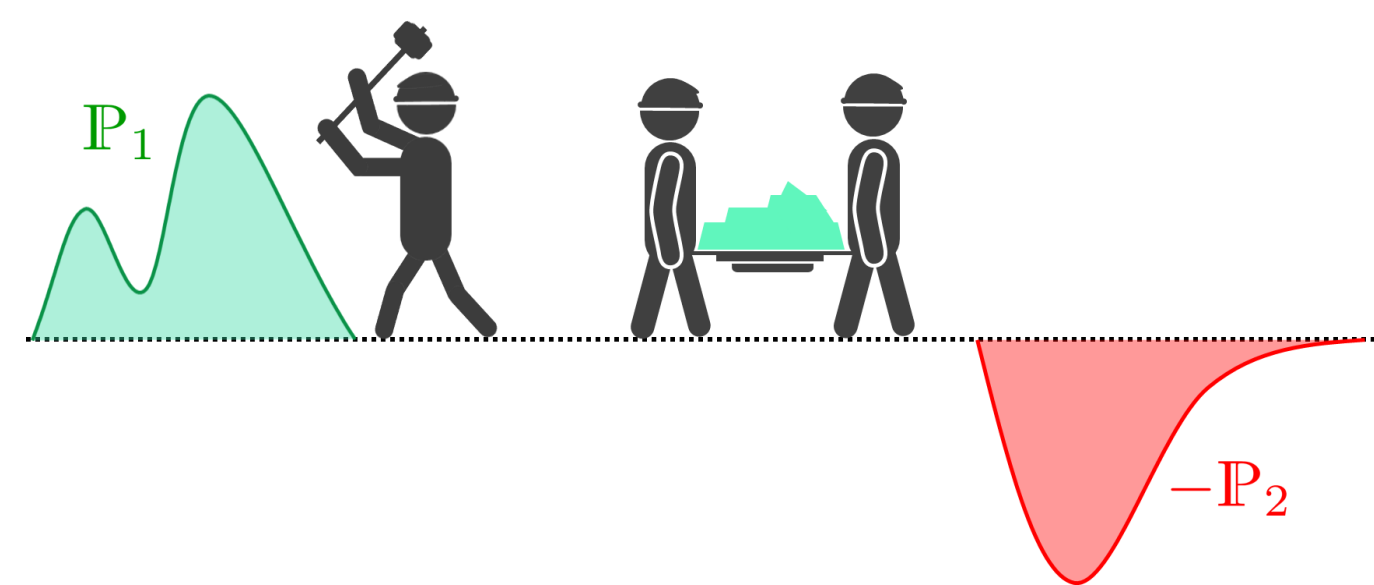
$$\text{RLR: } \min_{\beta} \frac{1}{N} \sum_{i=1}^N l_{\beta}(\hat{x}_i, \hat{y}_i) + c R(\beta)$$

Questions:

- Probabilistic interpretation for regularization?
- How to choose c and $R(\beta)$?



Wasserstein distance



$$W(\mathbf{P}_1, \mathbf{P}_2) := \inf_{\Pi} \mathbb{E}^{\Pi}[d(\xi_1, \xi_2)],$$

where \mathbf{P}_1 and \mathbf{P}_2 are the marginals of ξ_1 and ξ_2 under Π

Metric on feature-label-space:

$$d((x_1, y_1), (x_2, y_2)) = \|x_1 - x_2\| + \frac{\kappa}{2}|y_1 - y_2|$$

- κ : trust in labels
- $\kappa = \infty \Rightarrow$ exact labels

$W(\mathbf{P}_1, \mathbf{P}_2)$ = least cost of moving \mathbf{P}_1 to \mathbf{P}_2

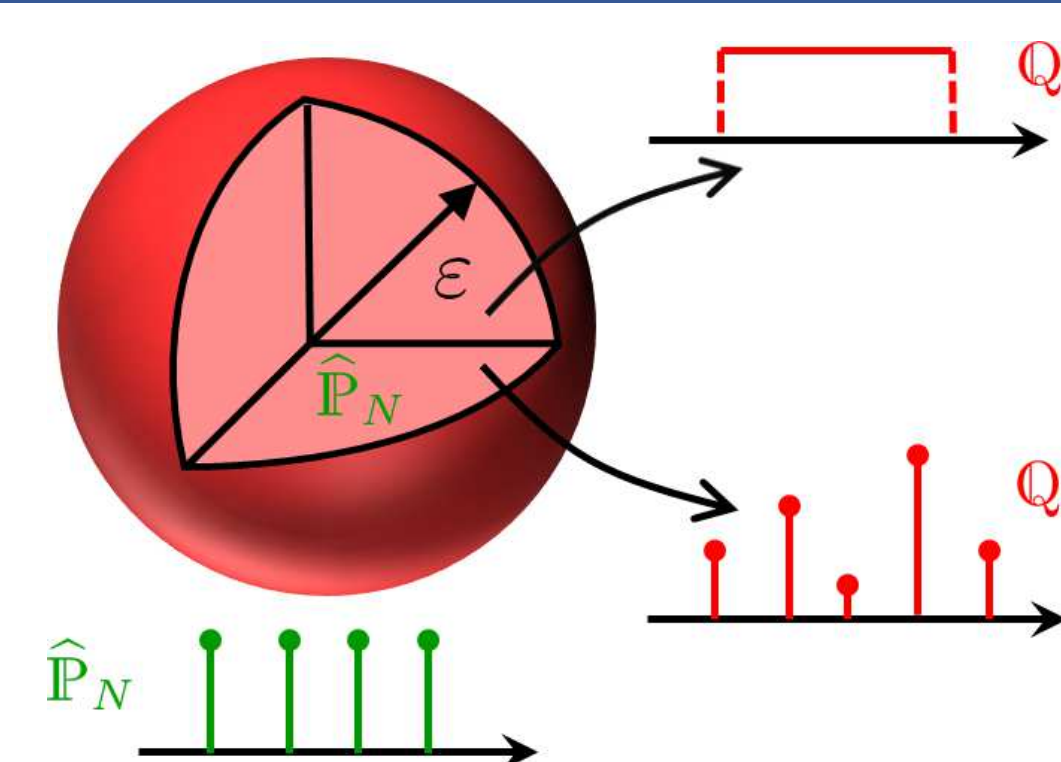
Distributionally Robust LR (DRLR)

Minimize the **worst-case** expected logloss

$$\text{DRLR: } \min_{\beta} \sup_{Q \in \mathcal{B}_{\varepsilon}(\hat{\mathbf{P}}_N)} \mathbb{E}^Q[l_{\beta}(x, y)]$$

The worst-case is taken over all distributions Q in the **Wasserstein ball**

$$\mathcal{B}_{\varepsilon}(\hat{\mathbf{P}}_N) := \{Q : W(Q, \hat{\mathbf{P}}_N) \leq \varepsilon\}$$



Tractability

Theorem 1: DRLR is equivalent to the finite convex program

$$\min_{\beta, \lambda, s_i} \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \quad \text{s.t.} \quad \begin{cases} l_{\beta}(\hat{x}_i, \hat{y}_i) \leq s_i & \forall i \leq N \\ l_{\beta}(\hat{x}_i, -\hat{y}_i) - \lambda \kappa \leq s_i & \forall i \leq N \\ \|\beta\|_* \leq \lambda \end{cases}$$

Special Case: DRLR reduces to RLR when $\kappa = \infty$

$$\inf_{\beta} \frac{1}{N} \sum_{i=1}^N \ell_{\beta}(\hat{x}_i, \hat{y}_i) + \varepsilon \|\beta\|_*$$

Out-of-Sample Performance

Theorem 2: If the tail of \mathbf{P} decays as $\exp(-\|2x\|^a)$ and the radius ε of the Wasserstein ball is set to

$$\varepsilon_N(\eta) = \left(\frac{\log(c_1 \eta^{-1})}{c_2 N} \right)^{\frac{1}{a}} \mathbb{1}_{\{N < \frac{\log(c_1 \eta^{-1})}{c_2 c_3}\}} + \left(\frac{\log(c_1 \eta^{-1})}{c_2 N} \right)^{\frac{1}{n}} \mathbb{1}_{\{N \geq \frac{\log(c_1 \eta^{-1})}{c_2 c_3}\}},$$

then we have

$$\mathbb{P}^N\{\mathbf{P} \in \mathcal{B}_{\varepsilon}(\hat{\mathbf{P}}_N)\} \geq 1 - \eta \quad \Rightarrow \quad \mathbb{P}^N\{\mathbb{E}^{\mathbf{P}}[l_{\beta}(x, y)] \leq \hat{J}\} \geq 1 - \eta \quad (*)$$

Note: $1 - \eta$ = confidence that \mathbf{P} is in the Wasserstein ball.

Risk Estimation

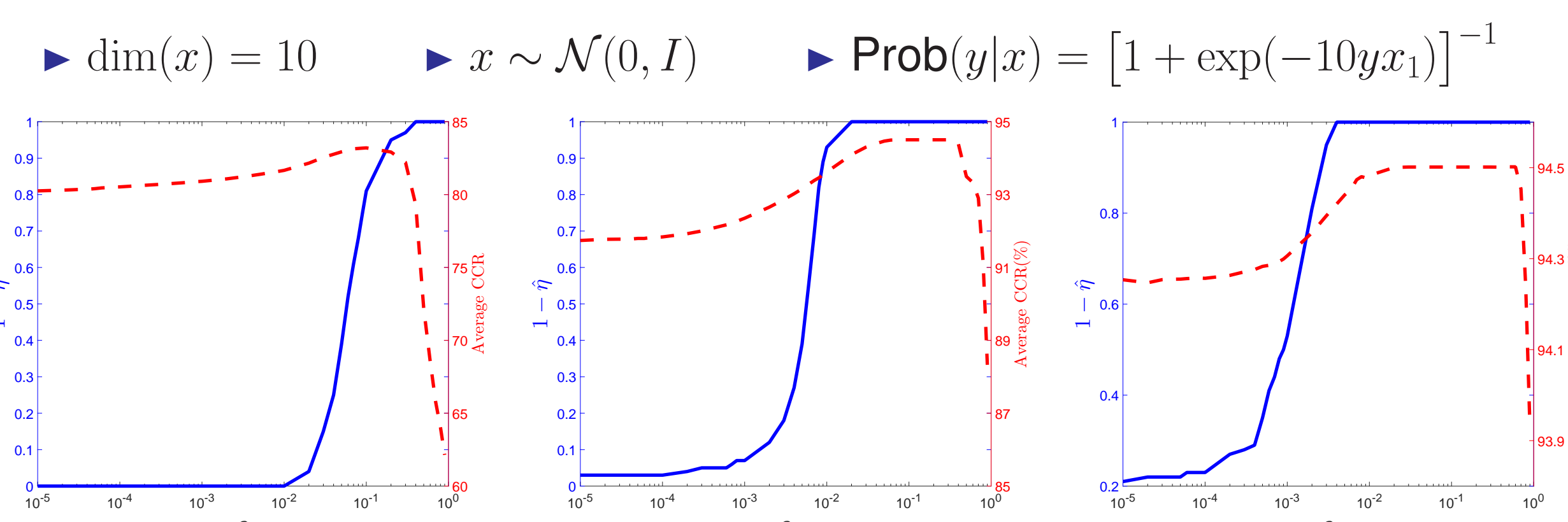
Theorem 3: The **worst/best-case** risk $\mathfrak{R}_{\max/\min} := \sup_{Q \in \mathcal{B}_{\varepsilon}(\hat{\mathbf{P}}_N)} \inf_{\beta} \mathbb{E}^Q[\mathbb{1}_{\{y(\beta, x) \leq 0\}}]$ is given by

$$\mathfrak{R}_{\max/\min} = 1 - \begin{cases} \min_{\lambda, s_i, r_i, t_i} \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} \quad \begin{cases} 1 - r_i \hat{y}_i \langle \hat{\beta}, \hat{x}_i \rangle \leq s_i & \forall i \leq N \\ 1 \pm t_i \hat{y}_i \langle \hat{\beta}, \hat{x}_i \rangle - \lambda \kappa \leq s_i & \forall i \leq N \\ r_i \|\hat{\beta}\|_* \leq \lambda, \quad t_i \|\hat{\beta}\|_* \leq \lambda & \forall i \leq N \\ r_i, t_i, s_i \geq 0 & \forall i \leq N \end{cases} \end{cases}$$

If $\varepsilon \geq \varepsilon_N(\eta)$, then $\mathfrak{R}_{\min}(\hat{\beta}) \leq \mathfrak{R}(\hat{\beta}) \leq \mathfrak{R}_{\max}(\hat{\beta})$ with probability $1 - 2\eta$.

Results

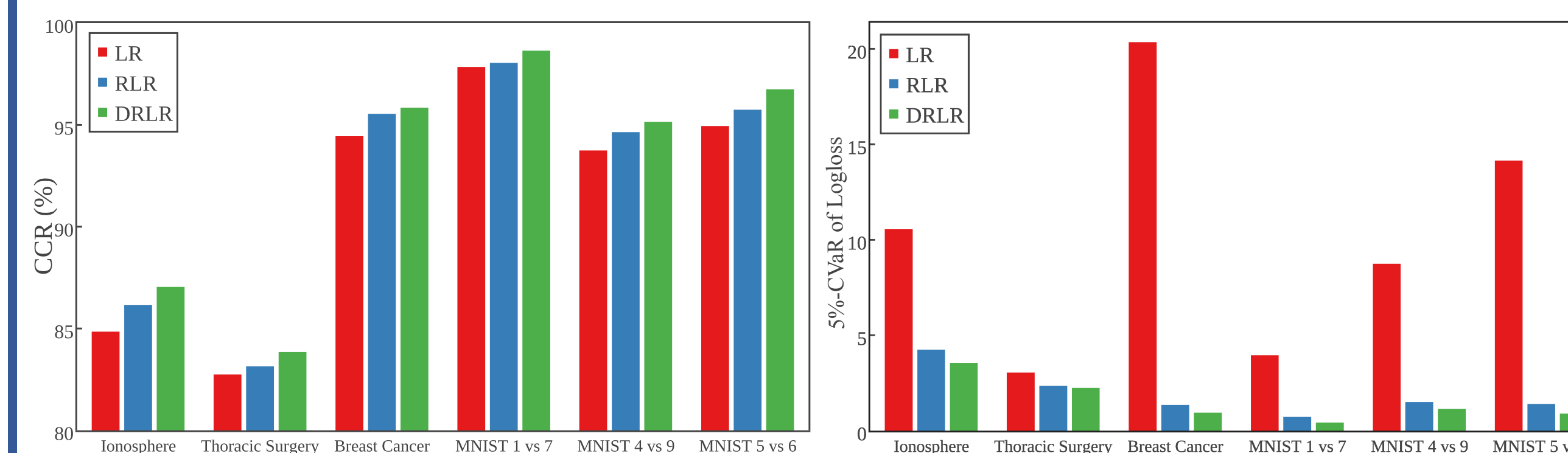
Synthetic Experiments (out-of-sample performance):



Observations:

- Empirical confidence that $\mathbf{P} \in \mathcal{B}_{\varepsilon}(\hat{\mathbf{P}}_N)$ saturates when the out-of-sample CCR is maximal
- The saturation point scales with N^{-1} consistent with $(*)$

Empirical Experiments: (MNIST & UCI datasets)



► Improvement of DRLR over RLR \approx improvement of RLR over LR

Ionosphere dataset: (UCI datasets)

