

## ABSTRACT

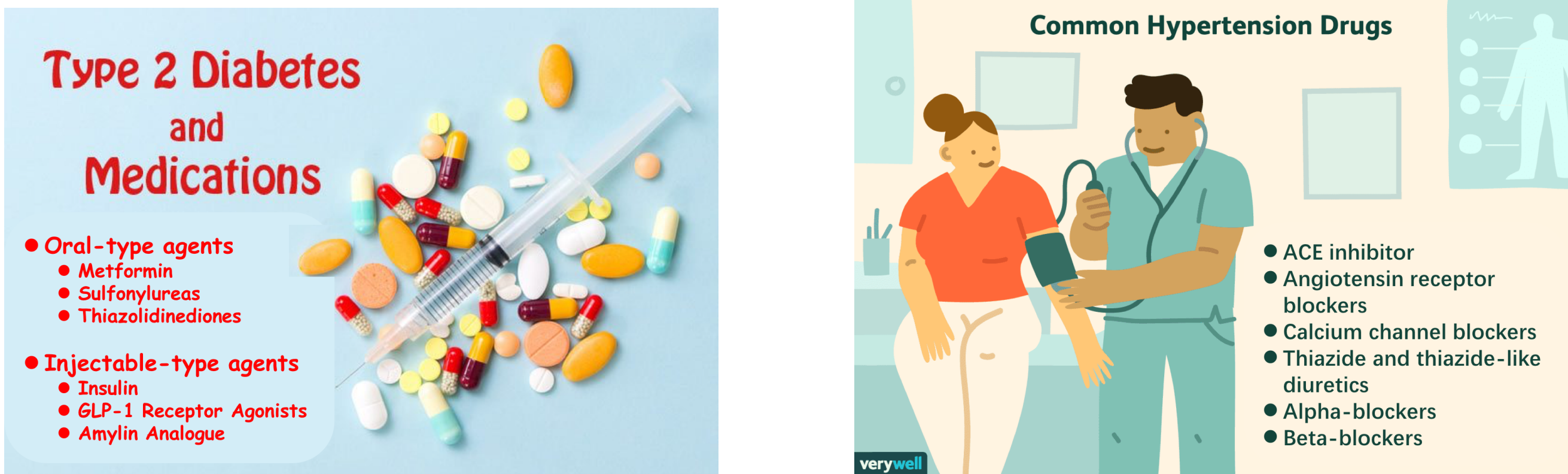
We develop a **prediction-based prescriptive model** for optimal decision making that

- predicts** the outcome under each action using a robust nonlinear model;
- prescribes** actions based on their predicted outcomes.

The *predictive* model combines *Distributionally Robust Linear Regression* (DRLR) with the *K-Nearest Neighbors* (K-NN) regression, which produces predictions that are robust to data perturbations and captures the nonlinearity embedded in the data. The *prescriptive* model selects each action with a probability inversely proportional to its exponentiated predicted outcome. We show theoretical guarantees on the out-of-sample performance of the predictive model, and prove the optimality of the randomized prescriptive policy in terms of the expected true future outcome. We demonstrate the proposed methodology on a diabetes and a hypertension dataset, showing that our prescribed treatment leads to a larger reduction in HbA<sub>1c</sub> and systolic blood pressure compared to a series of alternatives.

## PROBLEM DESCRIPTION

- Problem:** Given a set of drugs  $[M] \triangleq \{1, \dots, M\}$ , choose the one that yields the lowest future HbA<sub>1c</sub>/systolic blood pressure  $y$ , with the aid of patient data  $\mathbf{x}$  that is predictive of  $y$ .
- Idea:** *Predict* the outcome  $y_m$  for each drug  $m \in [M]$  using a **robust nonlinear** framework, and *prescribe* the actions based on their predictions.
- Applications:** Prescribe optimal treatments for patients with **diabetes** or **hypertension**.



*\*Publicly available Internet images.*

## ROBUST NONLINEAR PRESCRIPTION

- Assumption: For any  $m \in [M]$ ,  $y_m = \mathbf{x}_m' \beta_m^* + h_m(\mathbf{x}_m) + \varepsilon_m$ .
- Method:
  - For each  $m \in [M]$ , derive a robust estimate of  $\beta_m^*$ , denoted by  $\hat{\beta}_m$ , using *Wasserstein Distributionally Robust Optimization* (DRO)[1].
  - Given a new sample  $\mathbf{x}$ , find its  $K_m$  nearest neighbors, whose responses are denoted by  $y_{m(i)}$ ,  $i = [K_m]$ , in each action group  $m$  using the metric:

$$\|\mathbf{x} - \mathbf{x}_{m(i)}\|_{\hat{\mathbf{W}}_m} = \sqrt{(\mathbf{x} - \mathbf{x}_{m(i)})' \hat{\mathbf{W}}_m (\mathbf{x} - \mathbf{x}_{m(i)})},$$

where  $\hat{\mathbf{W}}_m = \text{diag}((\hat{\beta}_{m1})^2, \dots, (\hat{\beta}_{mp})^2)$ .

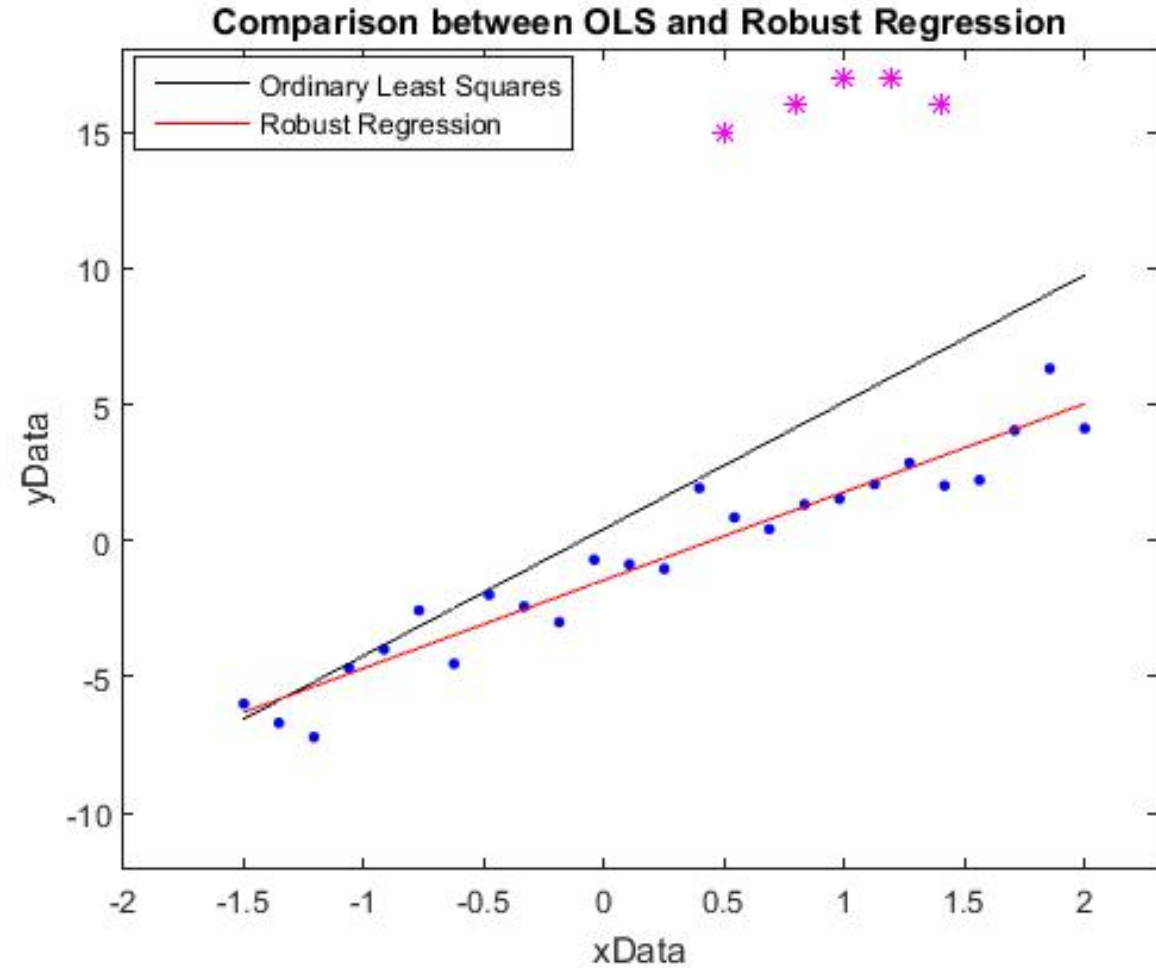
– Prediction:

$$\hat{y}_m(\mathbf{x}) = \frac{1}{K_m} \sum_{i=1}^{K_m} y_{m(i)}.$$

– Prescription: select action  $m$  with probability

$$e^{-\xi y_m(\mathbf{x})} / \sum_{j=1}^M e^{-\xi \hat{y}_j(\mathbf{x})}.$$

## ROBUST REGRESSION

- Goal:
    - Estimate the regression line that is not skewed by outliers.
- 
- The samples  $(\mathbf{x}_i, y_i)$ ,  $i = 1, \dots, N$ , may be contaminated with outliers.
  - Method: Inducing robustness by hedging against a set of uncertain parameters.

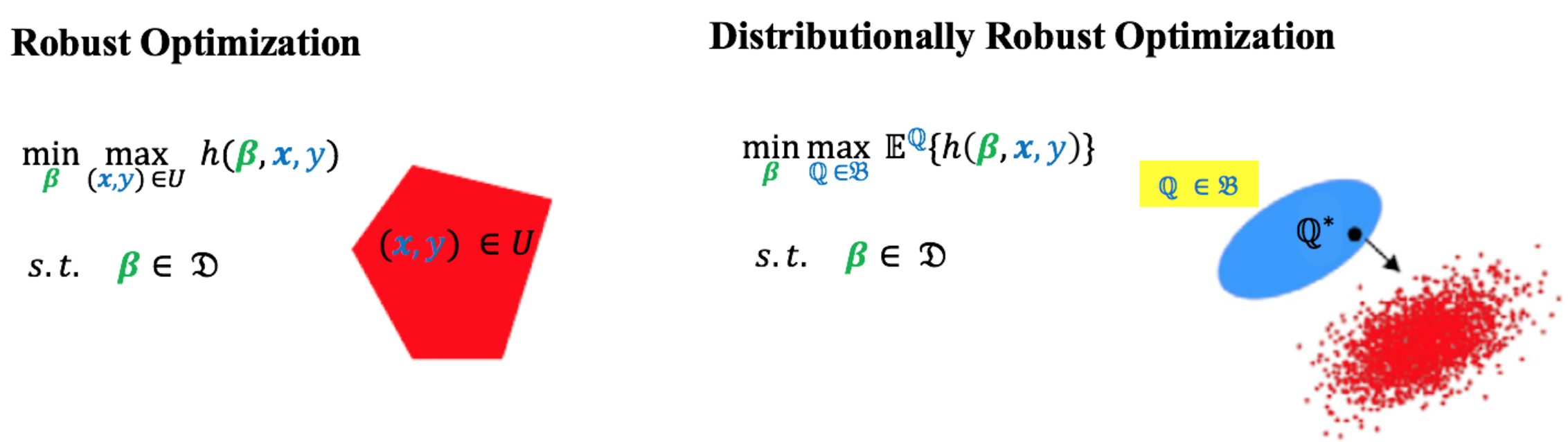


Fig. Comparisons of different optimization schemes.

## WASSERSTEIN DISTRIBUTIONALLY ROBUST OPTIMIZATION

- The Wasserstein DRO problem:

$$\inf_{\beta \in \mathcal{P}} \sup_{Q \in \mathcal{B}} E^Q[\|y - \mathbf{x}'\beta\|].$$

- Notation:
  - $\beta$ : the regression coefficient to be estimated;  $\mathbb{Q}$ : the probability distribution of  $(\mathbf{x}, y)$ .
  - $\mathcal{B}$ : the Wasserstein ball of distributions centered at the empirical distribution  $\hat{\mathbb{P}}_N$ :  $\mathcal{B} = \{\mathbb{Q} \in \mathcal{M}(\mathcal{X}) : W_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \varepsilon\}$ , where the Wasserstein distance is defined through,

$$W_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \triangleq \min_{\Pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X})} \left\{ \int_{\mathcal{X} \times \mathcal{X}} \|(\mathbf{x}_1, y_1) - (\mathbf{x}_2, y_2)\| \Pi(d(\mathbf{x}_1, y_1), d(\mathbf{x}_2, y_2)) \right\},$$

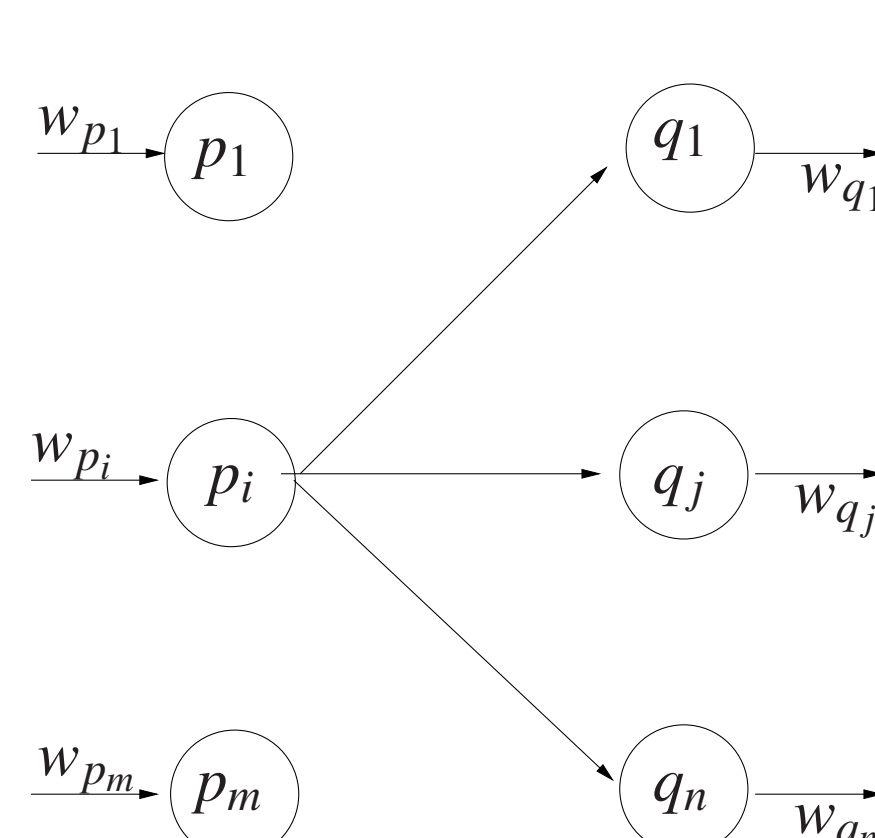
with  $\Pi$  the joint distribution of  $(\mathbf{x}_1, y_1)$  and  $(\mathbf{x}_2, y_2)$ , with marginals  $\mathbb{Q}$  and  $\hat{\mathbb{P}}_N$ .

## WHY THE WASSERSTEIN METRIC?

- Other options: Kullback-Leibler distance,  $f$ -divergences.
- Wasserstein metric incorporates a notion of **cost**:

$$W_1(\mathbb{Q}_1, \mathbb{Q}_2) = \inf_{\Pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X})} \left\{ \int_{\mathcal{X} \times \mathcal{X}} s(\mathbf{z}_1, \mathbf{z}_2) \Pi(d\mathbf{z}_1, d\mathbf{z}_2) \right\}.$$

- Allow support out of the observed samples.
- Wasserstein a.k.a. optimal mass transport, earth mover’s distance. Discrete case: **transportation problem**



$$W_1(\mathbb{P}, \mathbb{Q}) = \min_{\pi} \sum_{i=1}^m \sum_{j=1}^n \pi(i, j) s(i, j)$$

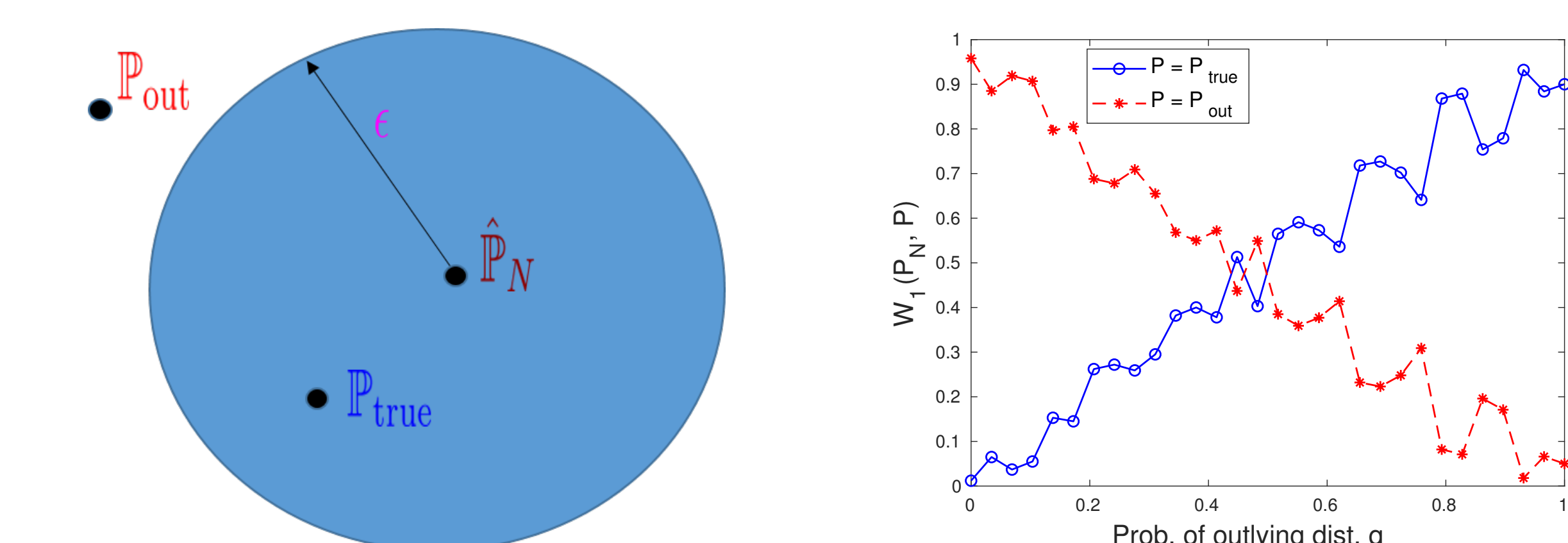
$$\text{s.t.} \quad \sum_{i=1}^m \pi(i, j) = w_{q_j}, \quad \forall j,$$

$$\sum_{j=1}^n \pi(i, j) = w_{p_i}, \quad \forall i,$$

$$\pi(i, j) \geq 0, \quad \forall i, j.$$

## ROBUSTNESS TO OUTLIERS

- Suppose we generate training data from a mixture of  $\mathbb{P}_{\text{true}}$  (w.p.  $1 - q$ ) and  $\mathbb{P}_{\text{out}}$  (w.p.  $q$ ).
- Then, for  $q < 0.5$ ,  $W_1(\mathbb{P}_{\text{true}}, \hat{\mathbb{P}}_N) < W_1(\mathbb{P}_{\text{out}}, \hat{\mathbb{P}}_N)$ . **Can exclude outliers!**



## ROBUSTNESS OF THE WASSERSTEIN SET

**Theorem 1** Suppose we are given two probability distributions  $\mathbb{P}_{\text{true}}$  and  $\mathbb{P}_{\text{out}}$ , and the mixture distribution  $\mathbb{P}_{\text{mix}}$  is a convex combination of the two:  $\mathbb{P}_{\text{mix}} = q\mathbb{P}_{\text{out}} + (1 - q)\mathbb{P}_{\text{true}}$ . Then,

$$\frac{W_1(\mathbb{P}_{\text{out}}, \mathbb{P}_{\text{mix}})}{W_1(\mathbb{P}_{\text{true}}, \mathbb{P}_{\text{mix}})} = \frac{1 - q}{q}.$$

- When  $q < 0.5$ , and  $W_1(\mathbb{P}_{\text{true}}, \mathbb{P}_{\text{mix}}) < W_1(\mathbb{P}_{\text{out}}, \mathbb{P}_{\text{mix}}) \implies$  the set  $\mathcal{B}$  will include the true distribution  $\mathbb{P}_{\text{true}}$  and exclude the outlying one  $\mathbb{P}_{\text{out}}$ .

## TRACTABLE RELAXATION OF WASSERSTEIN DRO

- The Wasserstein DRO problem could be relaxed to ( $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$  used in the Wasserstein metric)
$$\inf_{\beta \in \mathcal{P}} \varepsilon \|(-\beta, 1)\|_* + \frac{1}{N} \sum_{i=1}^N |y_i - \mathbf{x}_i' \beta|.$$
- Incorporates a class of models, e.g., **regularized LAD**, GLASSO with  $\ell_1$ -loss.
  - Connects **sparsity** with **robustness**.
  - New interpretation for the regularization coefficient  $\varepsilon$ .
  - The regularizer controls the **amount of ambiguity** in the data.

## ESTIMATION BIAS OF THE WASSERSTEIN DRO

**Theorem 2** Under mild conditions, when the sample size  $N_m \geq n_m$ , with probability at least  $\delta_m$ ,

$$\|\beta_m^* - \hat{\beta}_m\|_2 \leq \tau_m.$$

- The parameters  $n_m, \delta_m, \tau_m$  are related to the **Gaussian width** of the unit ball in  $\|\cdot\|_\infty$ , the **sub-Gaussian norm** of  $(\mathbf{x}_m, y_m)$ , the eigenvalues of the covariance matrix of  $(\mathbf{x}_m, y_m)$ , and the geometric structure of the true regression coefficient  $\beta_m^*$ .

## MSE OF WASSERSTEIN DRO INFORMED K-NN

- The bias-variance decomposition implies ( $\eta_m$  is the standard deviation of the noise  $\varepsilon_m$ ):

$$\begin{aligned} \text{MSE}(\hat{y}_m(\mathbf{x}) | \mathbf{x}, \mathbf{x}_{m(i)}, i = [N_m]) &\triangleq \mathbb{E} \left[ (\hat{y}_m(\mathbf{x}) - y_m(\mathbf{x}))^2 | \mathbf{x}, \mathbf{x}_{m(i)}, i = [N_m] \right] \\ &= \left( \frac{1}{K_m} \sum_{i=1}^{K_m} ((\mathbf{x} - \mathbf{x}_{m(i)})' \beta_m^* + h_m(\mathbf{x}) - h_m(\mathbf{x}_{m(i)})) \right)^2 + \frac{\eta_m^2}{K_m} + \eta_m^2. \end{aligned}$$

- For MSE to be small:
  - $\|\beta_m^* - \hat{\beta}_m\|_2$  is small;
  - $\|\mathbf{x} - \mathbf{x}_{m(i)}\|_{\hat{\mathbf{W}}_m}$  is small for  $i = [K_m]$ ; and
  - $h_m(\mathbf{x}) - h_m(\mathbf{x}_{m(i)})$  is small for  $i = [K_m]$ .

## DISTANCE TO THE K NEAREST NEIGHBORS

**Theorem 3** Suppose we are given  $N_m$  i.i.d. samples  $(\mathbf{x}_{mi}, y_{mi})$ ,  $i \in [N_m]$ , drawn from some unknown probability distribution with finite fourth moment. Every  $\mathbf{x}_{mi}$  has independent, centered coordinates:

$$\mathbb{E}(\mathbf{x}_{mi}) = \mathbf{0}, \quad \text{cov}(\mathbf{x}_{mi}) = \text{diag}(\sigma_{m1}^2, \dots, \sigma_{mp}^2), \forall i \in [N_m].$$

For a fixed predictor  $\mathbf{x}$ , and any given positive definite diagonal matrix  $\mathbf{W} \in \mathbb{R}^{p \times p}$  with diagonal elements  $w_j$ ,  $j \in [p]$ , and  $|w_j| \leq \bar{B}^2$ , suppose:

$$|(x_{mij} - x_j)^2 - (\sigma_{mj}^2 + x_j^2)| \leq T_m, \text{ a.s., } \forall i \in [N_m], j \in [p],$$

where  $x_{mij}, x_j$  are the  $j$ -th components of  $\mathbf{x}_{mi}$  and  $\mathbf{x}$ , respectively. Under the condition that  $\bar{w}_m^2 > \bar{B}^2 \sum_{j=1}^p (\sigma_{mj}^2 + x_j^2)$ , with probability at least  $1 - I_{1-p_{m0}}(N_m - K_m + 1, K_m)$ ,

$$\|\mathbf{x} - \mathbf{x}_{m(i)}\|_{\mathbf{W}} \leq \bar{w}_m, i \in [K_m].$$

## PREDICTIVE PERFORMANCE

**Theorem 4** Given a fixed predictor  $\mathbf{x} = (x_1, \dots, x_p)$ , and some scalar  $\bar{w}_m$ , assuming

- $h_m(\cdot)$  is Lipschitz continuous with a Lipschitz constant  $L_m$ .
- $\bar{w}_m^2 > \bar{B}_m^2 \sum_{j=1}^p (\sigma_{mj}^2 + x_j^2)$ .
- $|(x_{mij} - x_j)^2 - (\sigma_{mj}^2 + x_j^2)| \leq T_m, \forall i, j$ .
- The coordinates of any feasible solution to Wasserstein DRO have absolute values greater than or equal to some positive number  $b_m$  (dense estimators).

When  $N_m \geq n_m$ , with probability at least  $\delta_m - I_{1-p_{m0}}(N_m - K_m + 1, K_m)$  w.r.t. the measure of samples,

$$\mathbb{E}[(\hat{y}_m(\mathbf{x}) - y_m(\mathbf{x}))^2 | \mathbf{x}, \mathbf{x}_{m(i)}, i = [N_m]] \leq \left( \frac{\bar{w}_m \bar{\varepsilon}_m}{b_m} + \sqrt{p} \bar{w}_m + \frac{L_m \bar{w}_m}{\bar{B}_m} \right)^2 + \frac{\eta_m^2}{K_m} + \eta_m^2.$$

## PREScriptive PERFORMANCE

**Theorem 5** Given any  $\mathbf{x} \in \mathbb{R}^p$ , denote its *predicted* and *true* future outcome under action  $m$  by  $\hat{y}_m(\mathbf{x})$  and  $y_m(\mathbf{x})$ , respectively. For any  $k \in [M]$ , the *expected true outcome* under the *randomized prescriptive policy* satisfies:

$$\begin{aligned} \sum_{m=1}^M \frac{e^{-\xi \hat{y}_m(\mathbf{x})}}{\sum_j e^{-\xi \hat{y}_j(\mathbf{x})}} y_m(\mathbf{x}) &\leq y_k(\mathbf{x}) + \left( \hat{y}_k(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \hat{y}_m(\mathbf{x}) \right) \\ &\quad + \xi \left( \frac{1}{M} \sum_{m=1}^M \hat{y}_m^2(\mathbf{x}) + \sum_{m=1}^M \frac{e^{-\xi \hat{y}_m(\mathbf{x})}}{\sum_j e^{-\xi \hat{y}_j(\mathbf{x})}} y_m^2(\mathbf{x}) \right) + \frac{\log M}{\xi}. \end{aligned}$$

## ACTIVATE THE RANDOMIZED STRATEGY

- In consideration of the **health care costs** and **treatment transients**, we do not want to switch patients’ treatments too frequently.
- Threshold  $T(\mathbf{x})$  to activate the randomized strategy:

$$m_t(\mathbf{x}) = \begin{cases} m, \text{ w.p. } \frac{e^{-\xi \hat{y}_m(\mathbf{x})}}{\sum_{j=1}^M e^{-\xi \hat{y}_j(\mathbf{x})}}, & \text{if } \sum_k \frac{e^{-\xi \hat{y}_k(\mathbf{x})}}{\sum_j e^{-\xi \hat{y}_j(\mathbf{x})}} \hat{y}_k(\mathbf{x}) \leq x_{co} - T(\mathbf{x}), \\ m_c(\mathbf{x}), & \text{otherwise.} \end{cases}$$

- Find the largest  $T(\mathbf{x})$  such that the probability of the expected improvement being less than  $T(\mathbf{x})$  is small.
**Theorem 6** Assume that the distribution of the predicted outcome  $\hat{y}_m(\mathbf{x})$  conditional on  $\mathbf{x}$ , is *sub-Gaussian*, and its *psi<sub>2</sub>-norm* is equal to  $\sqrt{2} C_m(\mathbf{x})$ , for any  $m \in [M]$  and any  $\mathbf{x}$ . Given a small  $0 < \bar{\varepsilon} < 1$ , in order to satisfy

$$\mathbb{P} \left( \sum_k \frac{e^{-\xi \hat{y}_k(\mathbf{x})}}{\sum_j e^{-\xi \hat{y}_j(\mathbf{x})}} \hat{y}_k(\mathbf{x}) > x_{co} - T(\mathbf{x}) \right) \leq \bar{\varepsilon},$$

it suffices to set a threshold

$$T(\mathbf{x}) = \max \left( 0, \min_m \left( x_{co} - \mu_{\hat{y}_m}(\mathbf{x}) - \sqrt{-2C_m^2(\mathbf{x}) \log(\bar{\varepsilon}/M)} \right) \right),$$

where  $\mu_{\hat{y}_m}(\mathbf{x}) = \mathbb{E}[\hat{y}_m(\mathbf{x}) | \mathbf{x}]$ .

- As  $\xi \rightarrow \infty$ , the randomized policy becomes deterministic.

$$m_t(\mathbf{x}) = \begin{cases} \arg \min_m \hat{y}_m(\mathbf{x}), & \text{if } \min_m \hat{y}_m(\mathbf{x}) \leq x_{co} - T(\mathbf{x}), \\ m_c(\mathbf{x}), & \text{otherwise.} \end{cases}$$

A slight modification to the threshold level  $T(\mathbf{x})$  is given below:

$$T(\mathbf{x}) = \max \left( 0, \min_m \left( x_{co} - \mu_{\hat{y}_m}(\mathbf{x}) - \sqrt{-2C_m^2(\mathbf{x}) \log \bar{\varepsilon}} \right) \right).$$

## ESTIMATE $\mu_{\hat{y}_m}(\mathbf{x})$ AND $C_m^2(\mathbf{x})$

**Algorithm 1** Estimating the conditional mean and standard deviation of the predicted outcome.

**Input:** a feature vector  $\mathbf{x}$ ;  $a_m$ : the number of subsamples used to compute  $\hat{\beta}_m$ ,  $a_m < N_m$ ;  $d_m$ : the number of repetitions.  
**for**  $i = 1, \dots, d_m$  **do**  
    Randomly pick  $a_m$  samples from group  $m$ , and use them to estimate a robust regression coefficient  $\hat{\beta}_{m_i}$  through solving Wasserstein DRO.  
    The future outcome for  $\mathbf{x}$  under action  $m$  is predicted as  $\hat{y}_{m_i}(\mathbf{x}) = \mathbf{x}' \hat{\beta}_{m_i}$ .  
**end for**  
**Output:** Estimate the conditional mean of  $\hat{y}_m(\mathbf{x})$  as:

$$\mu_{\hat{y}_m}(\mathbf{x}) = \frac{1}{d_m} \sum_{i=1}^{d_m} \hat{y}_{m_i}(\mathbf{x}),$$

and the conditional standard deviation as:

$$C_m(\mathbf{x}) = \sqrt{\frac{1}{d_m - 1} \sum_{i=1}^{d_m} \left( \hat{y}_{m_i}(\mathbf{x}) - \mu_{\hat{y}_m}(\mathbf{x}) \right)^2}.$$

## PREScribe OPTIMAL TREATMENTS

- Goal: develop optimal prescriptions for patients with **type-2 diabetes** and hypertension using the Electronic Health Records (EHRs).
- Predictors:** demographics, diagnoses, lab tests, and past admission records.
- Response:** HbA<sub>1c</sub>, and systolic blood pressure.
- The reduction in HbA<sub>1c</sub>/systolic blood pressure, mean (std.):

	Diabetes		Hypertension	
	Deterministic	Randomized	Deterministic	Randomized
LASSO	-0.51 (0.16)	-0.51 (0.16)	-4.22 (0.20)	-4.22 (0.19)
CART	-0.45 (0.13)	-0.42 (0.14)	-4.48 (0.55)	-4.51 (0.49)
OLS+K-NN	-0.53 (0.13)	-0.53 (0.13)	-4.27 (0.32)	-4.29 (0.31)
DRO+K-NN	-0.56 (0.06)	-0.55 (0.08)	-6.58 (0.70)	-6.78 (0.73)
Current therapy	-0.22 (0.04)		-2.50 (0.16)	
Standard of care	-0.22 (0.03)		-2.37 (0.11)	

## REFINEMENT OF THE POLICY

- K-NN is sensitive to the number of neighbors  $K_m$ .
- Propose a patient-specific number of neighbors  $K'_m$ , where the neighbors that are relatively far away from the patient in query are discarded.
- Denote by  $d_m^i$  the distance between the patient in query and her  $i$ -th closest neighbor in group  $m$ , and define  $j_m^* = \arg \max_j \left( d_j^m - \sum_{i=1}^{j-1} \frac{d_i^m}{j-1} \right)$ .

$$K'_m = \begin{cases} j_m^* - 1, & \text{if } \frac{d_m^m - \sum_{i=1}^{j_m^*-1} \frac{d_i^m}{j_m^*-1}}{\sum_{i=1}^{j_m^*-1} \frac{d_i^m}{j_m^*-1}} > \bar{T}, \\ K_m, & \text{otherwise,} \end{cases}$$

where  $\bar{T}$  is some threshold that can be tuned using cross-validation.

- Results on the diabetes and hypertension datasets:

	Diabetes		Hypertension	
	Deterministic	Randomized	Deterministic	Randomized
LASSO	-0.54 (0.19)	-0.54 (0.20)	-4.34 (0.28)	-4.33 (0.28)
CART	-0.62 (0.32)	-0.57 (0.27)	-4.46 (0.46)	-4.49 (0.50)
OLS+K-NN	-0.65 (0.25)	-0.64 (0.25)	-4.30 (0.35)	-4.30 (0.32)
DRO+K-NN	-0.68 (0.20)	-0.67 (0.23)	-7.42 (0.46)	-7.58 (0.51)
Current therapy	-0.23 (0.05)		-2.56 (0.14)	
Standard of care	-0.22 (0.03)		-2.37 (0.11)	

## CONCLUSIONS

- All models outperform the current prescription and the standard of care.
- The **Wasserstein DRO+K-NN** model leads to the **largest reduction** in outcomes with a relatively stable performance.
- The best **DRO+K-NN** model leads to a **69%** reduction in future systolic blood pressure compared to the 2nd best model.
- Using a patient-specific  $K'_m$  in general leads to a more significant reduction in outcomes.
- The **randomized policy** achieves a similar (slightly better) performance than the deterministic one.

## REFERENCES

- Chen, R., and Paschalidis, I.C. (2018). A robust learning algorithm for regression models using distributionally robust optimization under the Wasserstein metric, **Journal of Machine Learning Research**, 19, 1-48.
- Chen, R., and Paschalidis, I.C. (2018). Learning optimal personalized treatment rules using robust regression informed K-NN, **NIPS Machine Learning for Health (ML4H) workshop**, Montreal, Canada.
- Bertsimas, D., Kallus, N., Weinstein, A.M., and Zhuo, Y.D. (2017). Personalized diabetes management using electronic medical records, **Diabetes Care**.