

(Wasserstein) distributionally robust optim. in action

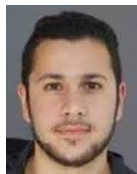
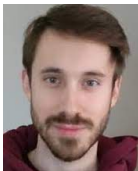
Jérôme MALICK

CNRS, Lab. Jean Kuntzmann & MIAI (Institut IA de Grenoble)



Data-driven optim. workshop – Ecole des Ponts – Oct. 2023

Based on joint work with
Waïss Azizian, Franck Iutzeler, Yassine Laguel, Zaid Harchaoui, Krishna Pillutla



Set-up: data-driven optimization under uncertainty

- Training data: ξ_1, \dots, ξ_N (in theory: sampled from $\mathbb{P}_{\text{train}}$ unknown)
e.g. in supervised learning: labeled data $\xi_i = (a_i, y_i)$ feature, label
- Train model: $f(x, \cdot)$ the loss function with x the parameter/decision $(\omega, \beta, \theta, \dots)$
e.g. least-square regression: $f(x, (a, y)) = (x^\top a - y)^2$
- Compute x via empirical risk minimization (a.k.a SAA)
(minimize the average loss on training data)

$$\min_x \frac{1}{N} \sum_{i=1}^N f(x, \xi_i)$$

- Prediction with x for different data ξ
 - Distributional shifts: $\mathbb{P}_{\text{train}} \neq \mathbb{P}_{\text{test}}$
 - Adversarial attacks
 - Generalization: computations with $\hat{\mathbb{P}}_N$ and guarantees on $\mathbb{P}_{\text{train}}$
 - Other situations, e.g. heterogeneous data
- Solution: take possible variations into account during training

Set-up: data-driven optimization under uncertainty

- Training data: ξ_1, \dots, ξ_N (in theory: sampled from $\mathbb{P}_{\text{train}}$ unknown)
e.g. in supervised learning: labeled data $\xi_i = (a_i, y_i)$ feature, label
- Train model: $f(x, \cdot)$ the loss function with x the parameter/decision $(\omega, \beta, \theta, \dots)$
e.g. least-square regression: $f(x, (a, y)) = (x^\top a - y)^2$
- Compute x via empirical risk minimization (a.k.a SAA)
(minimize the average loss on training data)

$$\min_x \frac{1}{N} \sum_{i=1}^N f(x, \xi_i) = \mathbb{E}_{\hat{\mathbb{P}}_N} [f(x, \xi)] \quad \text{with } \hat{\mathbb{P}}_N = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_i}$$

- Prediction with x for different data ξ
 - Distributional shifts: $\mathbb{P}_{\text{train}} \neq \mathbb{P}_{\text{test}}$
 - Adversarial attacks
 - Generalization: computations with $\hat{\mathbb{P}}_N$ and guarantees on $\mathbb{P}_{\text{train}}$
 - Other situations, e.g. heterogeneous data
- Solution: take possible variations into account during training

(Distributionally) robust optimization

Optimize expected loss for the worst probability in a set of perturbations

rather than $\min_x \mathbb{E}_{\hat{\mathbb{P}}_N}[f(x, \xi)]$ solve instead $\min_x \max_{Q \in \mathcal{U}} \mathbb{E}_Q[f(x, \xi)]$

with \mathcal{U} a neighborhood of $\hat{\mathbb{P}}_N$ (called ambiguity set)

(Distributionally) robust optimization

Optimize expected loss for the worst probability in a set of perturbations

rather than $\min_x \mathbb{E}_{\hat{\mathbb{P}}_N}[f(x, \xi)]$ solve instead $\min_x \max_{\mathbb{Q} \in \mathcal{U}} \mathbb{E}_{\mathbb{Q}}[f(x, \xi)]$

with \mathcal{U} a neighborhood of $\hat{\mathbb{P}}_N$ (called ambiguity set)

modeling vs. computational tractability

- $\mathcal{U} = \{\hat{\mathbb{P}}_N\}$: $\min_x \frac{1}{N} \sum_{i=1}^N f(x, \xi_i)$ standard ERM
- $\mathcal{U} = \{\mathbb{Q} : \text{supp}(\mathbb{Q}) \subset U\}$: $\min_x \max_{\xi \in U} f(x, \xi)$ standard robust optimization
- \mathcal{U} defined by moments e.g. [Delage, Ye, '10]
- $\mathcal{U} = \{\mathbb{Q} : d(\hat{\mathbb{P}}_N, \mathbb{Q}) \leq \rho\}$ for various distances or divergences
E.g. KL-div., χ^2 -div., max-mean-discrepancy... e.g. [Namkoong, Duchi '17]
- $\mathcal{U} = \{\mathbb{Q} : W(\hat{\mathbb{P}}_N, \mathbb{Q}) \leq \rho\}$ Wasserstein distance [Kuhn *et al.* '18]

Wasserstein balls as ambiguity sets

Def: Wasserstein distance (given a cost function c)

$$W(\mathbb{P}, \mathbb{Q}) = \min_{\pi} \{ \mathbb{E}_{\pi}[c(\xi, \xi')] : \pi \text{ with marginals } [\pi]_1 = \mathbb{P} \text{ and } [\pi]_2 = \mathbb{Q} \}$$

Wasserstein balls as ambiguity sets

Def: Wasserstein distance (given a cost function c)

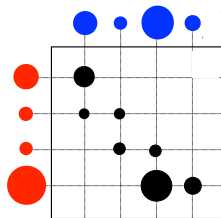
$$W(\mathbb{P}, \mathbb{Q}) = \min_{\pi} \{ \mathbb{E}_{\pi}[c(\xi, \xi')] : \pi \text{ with marginals } [\pi]_1 = \mathbb{P} \text{ and } [\pi]_2 = \mathbb{Q} \}$$

Demystification: in the discrete case

e.g. $\mathbb{P} = (p_1, \dots, p_N)$ and $\mathbb{Q} = (q_1, \dots, q_N)$ in the simplex

$$\left\{ \begin{array}{l} \min_{\pi} \sum_{i,j=1}^N c_{i,j} \pi_{i,j} \\ \sum_{j=1}^N \pi_{i,j} = p_i \quad i = 1, \dots, N \\ \sum_{i=1}^N \pi_{i,j} = q_j \quad j = 1, \dots, N \\ \pi_{i,j} \geq 0 \quad i, j = 1, \dots, N \end{array} \right.$$

linear assignment !



Wasserstein balls as ambiguity sets

Def: Wasserstein distance (given a cost function c)

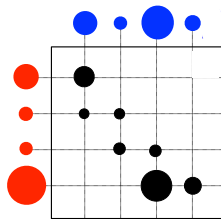
$$W(\mathbb{P}, \mathbb{Q}) = \min_{\pi} \{ \mathbb{E}_{\pi}[c(\xi, \xi')] : \pi \text{ with marginals } [\pi]_1 = \mathbb{P} \text{ and } [\pi]_2 = \mathbb{Q} \}$$

Demystification: in the discrete case

e.g. $\mathbb{P} = (p_1, \dots, p_N)$ and $\mathbb{Q} = (q_1, \dots, q_N)$ in the simplex

$$\begin{cases} \min_{\pi} \sum_{i,j=1}^N c_{i,j} \pi_{i,j} \\ \sum_{j=1}^N \pi_{i,j} = p_i \quad i = 1, \dots, N \\ \sum_{i=1}^N \pi_{i,j} = q_j \quad j = 1, \dots, N \\ \pi_{i,j} \geq 0 \quad i, j = 1, \dots, N \end{cases}$$

linear assignment !



Wasserstein-DRO objective for given \mathbb{P} and ρ

$$\begin{cases} \max_{\mathbb{Q}} \mathbb{E}_{\mathbb{Q}}[f(\xi)] \\ W(\mathbb{P}, \mathbb{Q}) \leq \rho \end{cases}$$

Wasserstein balls as ambiguity sets

Def: Wasserstein distance (given a cost function c)

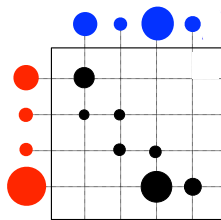
$$W(\mathbb{P}, \mathbb{Q}) = \min_{\pi} \{ \mathbb{E}_{\pi}[c(\xi, \xi')] : \pi \text{ with marginals } [\pi]_1 = \mathbb{P} \text{ and } [\pi]_2 = \mathbb{Q} \}$$

Demystification: in the discrete case

e.g. $\mathbb{P} = (p_1, \dots, p_N)$ and $\mathbb{Q} = (q_1, \dots, q_N)$ in the simplex

$$\begin{cases} \min_{\pi} \sum_{i,j=1}^N c_{i,j} \pi_{i,j} \\ \sum_{j=1}^N \pi_{i,j} = p_i \quad i = 1, \dots, N \\ \sum_{i=1}^N \pi_{i,j} = q_j \quad j = 1, \dots, N \\ \pi_{i,j} \geq 0 \quad i, j = 1, \dots, N \end{cases}$$

linear assignment !



Wasserstein-DRO objective for given \mathbb{P} and ρ

$$\begin{cases} \max_{\mathbb{Q}} \mathbb{E}_{\mathbb{Q}}[f(\xi)] \\ W(\mathbb{P}, \mathbb{Q}) \leq \rho \end{cases} \Leftrightarrow \begin{cases} \max_{\mathbb{Q}, \pi} \mathbb{E}_{\mathbb{Q}}[f(\xi)] \\ [\pi]_1 = \mathbb{P}, [\pi]_2 = \mathbb{Q} \\ \min_{\pi} \mathbb{E}_{\pi}[c(\xi, \xi')] \leq \rho \end{cases}$$

Wasserstein balls as ambiguity sets

Def: Wasserstein distance (given a cost function c)

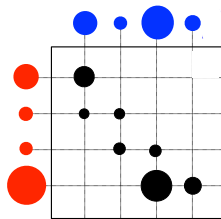
$$W(\mathbb{P}, \mathbb{Q}) = \min_{\pi} \{ \mathbb{E}_{\pi}[c(\xi, \xi')] : \pi \text{ with marginals } [\pi]_1 = \mathbb{P} \text{ and } [\pi]_2 = \mathbb{Q} \}$$

Demystification: in the discrete case

e.g. $\mathbb{P} = (p_1, \dots, p_N)$ and $\mathbb{Q} = (q_1, \dots, q_N)$ in the simplex

$$\begin{cases} \min_{\pi} \sum_{i,j=1}^N c_{i,j} \pi_{i,j} \\ \sum_{j=1}^N \pi_{i,j} = p_i \quad i = 1, \dots, N \\ \sum_{i=1}^N \pi_{i,j} = q_j \quad j = 1, \dots, N \\ \pi_{i,j} \geq 0 \quad i, j = 1, \dots, N \end{cases}$$

linear assignment !



Wasserstein-DRO objective for given \mathbb{P} and ρ

$$\begin{cases} \max_{\mathbb{Q}} \mathbb{E}_{\mathbb{Q}}[f(\xi)] \\ W(\mathbb{P}, \mathbb{Q}) \leq \rho \end{cases} \Leftrightarrow \begin{cases} \max_{\mathbb{Q}, \pi} \mathbb{E}_{\mathbb{Q}}[f(\xi)] \\ [\pi]_1 = \mathbb{P}, [\pi]_2 = \mathbb{Q} \\ \min_{\pi} \mathbb{E}_{\pi}[c(\xi, \xi')] \leq \rho \end{cases} \Leftrightarrow \begin{cases} \max_{\pi} \mathbb{E}_{[\pi]_2}[f(\xi)] \\ [\pi]_1 = \mathbb{P} \\ \mathbb{E}_{\pi}[c(\xi, \xi')] \leq \rho \end{cases}$$

Wasserstein DRO (WDRO)

WDRO objective for given \mathbb{P} , ρ – and choice $c(\xi, \xi') = \|\xi - \xi'\|^2$

$$\text{(Primal)} \quad \begin{cases} \max_{\pi} \mathbb{E}_{[\pi]_2}[f(\xi)] \\ [\pi]_1 = \mathbb{P} \\ \mathbb{E}_{\pi}[\|\xi - \xi'\|^2] \leq \rho \end{cases}$$

$$\text{(Dual)} \quad \min_{\lambda \geq 0} \lambda \rho + \mathbb{E}_{\mathbb{P}}[\max_{\xi'} \{f(\xi') - \lambda \|\xi - \xi'\|^2\}]$$

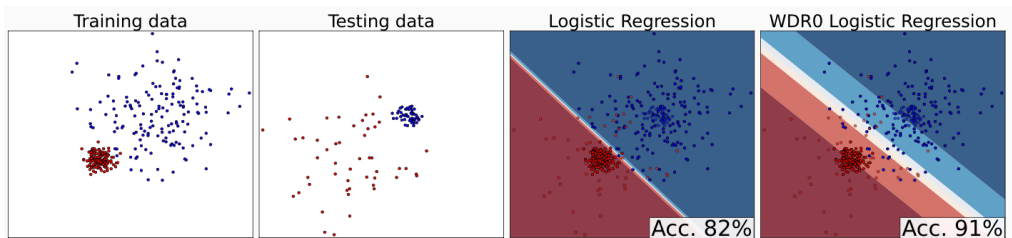
Wasserstein DRO (WDRO)

WDRO objective for given \mathbb{P}, ρ – and choice $c(\xi, \xi') = \|\xi - \xi'\|^2$

$$\text{(Primal)} \quad \begin{cases} \max_{\pi} \mathbb{E}_{[\pi]_2} [f(\xi)] \\ [\pi]_1 = \mathbb{P} \\ \mathbb{E}_{\pi} [\|\xi - \xi'\|^2] \leq \rho \end{cases}$$

$$\text{(Dual)} \quad \min_{\lambda \geq 0} \lambda \rho + \mathbb{E}_{\mathbb{P}} [\max_{\xi'} \{f(\xi') - \lambda \|\xi - \xi'\|^2\}]$$

Ex: classification. Compute separator θ by min. logistic loss $f_{\theta}(\xi) = \log(1 + \exp(-y a^{\top} \theta))$
 $\xi_i = (a_i, y_i) \in \mathbb{R}^2 \times \{-1, +1\}$ sampled from two Gaussian distributions
with variances $\sigma = 1$ and $\sigma = 5$ – reversed in testing !



WDRO : success !(?)

WDRO is very attractive

- Statistical/practical properties
e.g. [Blanchet *et al.* '18] and [Blanchet and Shapiro '23]
- Computable in many cases
e.g. [Kuhn *et al.* '18], [Zhao Guan '18]...
- Natural in many applications
back to [Scarf 1958] ! + (...) + recent trend in learning, e.g. [Kuhn *et al.* '20]
- Interprets up to first-order as a penalization by $\|\nabla_{\xi} f(x, \xi)\|$ e.g. [Gao *et al.* '18]

WDRO : success !(?)

WDRO is very attractive – in theory **and in practice ?**

- Statistical/practical properties – **warning : dimensionality ! (spotlight #1)**
e.g. [Blanchet *et al.* '18] and [Blanchet and Shapiro '23]
- Computable in many cases – **but not always ! (spotlight #2)**
e.g. [Kuhn *et al.* '18], [Zhao Guan '18]...
- Natural in many applications – **but not always ! (spotlight #3)**
back to [Scarf 1958] ! + (...) + recent trend in learning, e.g. [Kuhn *et al.* '20]
- Interprets up to first-order as a penalization by $\|\nabla_{\xi} f(x, \xi)\|$ e.g. [Gao *et al.* '18]

spotlight #1 : generalization guarantees



Azizian Waiss, Franck lutzeler, and Jérôme Malick

Exact generalization guarantees for (regularized) WDRO models

Just accepted in [NeurIPS, 2023](#)

Existing generalization guarantees

- Suppose $\xi_1, \dots, \xi_N \sim \mathbb{P}_{\text{train}}$ (where $\xi \in \mathbb{R}^d$)
- Computations with $\hat{\mathbb{P}}_N = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_i} \dots$ guarantees with $\mathbb{P}_{\text{train}}$?
- We manipulate the WDRO risk $R_\rho(f) = \max_{W(\hat{\mathbb{P}}_N, \mathbb{Q}) \leq \rho} \mathbb{E}_{\mathbb{Q}}[f(\xi)]$
- Obviously, if ρ, N large enough such that $W(\mathbb{P}_{\text{train}}, \hat{\mathbb{P}}_N) \leq \rho$, then

$$\underbrace{R_\rho(f)}_{\text{can compute \& optimize}} \geq \underbrace{\mathbb{E}_{\mathbb{P}_{\text{train}}}[f(\xi)]}_{\text{cannot access}}$$

Existing generalization guarantees

- Suppose $\xi_1, \dots, \xi_N \sim \mathbb{P}_{\text{train}}$ (where $\xi \in \mathbb{R}^d$)
- Computations with $\hat{\mathbb{P}}_N = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_i}$... guarantees with $\mathbb{P}_{\text{train}}$?
- We manipulate the WDRO risk $R_\rho(f) = \max_{W(\hat{\mathbb{P}}_N, \mathbb{Q}) \leq \rho} \mathbb{E}_{\mathbb{Q}}[f(\xi)]$
- Obviously, if ρ, N large enough such that $W(\mathbb{P}_{\text{train}}, \hat{\mathbb{P}}_N) \leq \rho$, then

$$\underbrace{R_\rho(f)}_{\text{can compute \& optimize}} \geq \underbrace{\mathbb{E}_{\mathbb{P}_{\text{train}}}[f(\xi)]}_{\text{cannot access}}$$

- But it requires $\rho \propto 1/\sqrt[2]{N}$ [Fournier and Guillin '15]
- Not optimal: $\rho \propto 1/\sqrt{N}$ suffices
 - asymptotically [Blanchet *et al* '22]
 - in particular cases [Shafieez-Adehabadeh *et al* '19]
 - or with error terms [Gao '22]

Extended exact generalization guarantees

By a direct approach (work direct to get a concentration result on the (dual) objective)

Theorem ([Azizian, Iutzeler, M. '23])

Assumptions : compactness on ξ + compactness on f + quad. growth of f near its minimizers

For $\delta \in (0, 1)$, if $\rho \geq O\left(\sqrt{\frac{\log 1/\delta}{N}}\right)$

Generalization guarantee: w.p. $1 - \delta$, $R_\rho(f) \geq \mathbb{E}_{\mathbb{P}_{\text{train}}}[f(\xi)]$

Distribution shifts: w.p. $1 - \delta$,

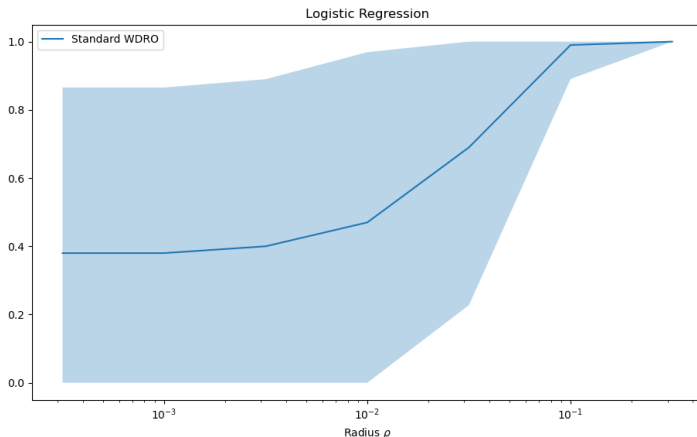
$$W(\mathbb{P}, \mathbb{Q})^2 \leq \rho \left(\rho - O\left(\sqrt{\frac{\log 1/\delta}{N}}\right) \right) \quad \text{it holds} \quad R_\rho(f) \geq \mathbb{E}_{\mathbb{Q}}[f(\xi)]$$

Assumptions valid in many cases: linear/logistic regression, kernel models, smooth neural networks, family of invertible mappings (e.g. normalizing flows)

Illustration

On logistic regression:

- for each ρ , sample 200 training datasets
- solve the WDRO problem on each of them [Blanchet *et al* '22]
- plot the proba of $R_\rho(f) - \mathbb{E}_{\mathbb{P}_{\text{train}}}[f] \geq 0$ (average, standard deviation)
- the training robust loss is indeed an upper-bound on the true loss



spotlight #2 – towards efficient computational toolkit: smoothed/regularized WDRO



Azizian Waiss, Franck Iutzeler, and Jérôme Malick

Regularization for Wasserstein distributionally robust optimization

ESAIM:COCV (Control Optim. Calculus of Variations), 2023

WDRO objective to be minimized

Dual WDRO is nonsmooth (which complicates resolution [Kuhn *et al.* '18])

$$R_\rho(f) = \min_{\lambda \geq 0} \lambda \rho + \mathbb{E}_{\mathbb{P}}[\max_{\xi'} \{f(\xi') - \lambda \|\xi - \xi'\|^2\}]$$

WDRO objective to be minimized

Dual WDRO is nonsmooth (which complicates resolution [Kuhn et al. '18])

$$R_\rho(f) = \min_{\lambda \geq 0} \lambda \rho + \mathbb{E}_{\mathbb{P}}[\max_{\xi'} \{f(\xi') - \lambda \|\xi - \xi'\|^2\}]$$

What about smoothing ? Smoothed counterpart

$$R_\rho^\varepsilon(f) = \min_{\lambda \geq 0} \lambda \rho + \varepsilon \mathbb{E}_{\mathbb{P}} \log \left(\mathbb{E}_{\xi' \sim \pi_0(\cdot|\xi)} e^{\frac{f(\xi') - \lambda \|\xi - \xi'\|^2}{\varepsilon}} \right)$$

(Nice interpretation as entropy-regularized WDRO)

Theorem (approximation bounds for WDRO [Azizian, Iutzeler, M. '21])

Under mild assumptions (non-degeneracy, Lipschitz), if the support of \mathbb{P} is contained in a compact convex set $\Xi \subset \mathbb{R}^d$, then

$$0 \leq R_\rho(f) - R_\rho^\varepsilon(f) \leq \left(C \varepsilon \log \frac{1}{\varepsilon} \right) d$$

Great ! but no computational results to show yet...

Smoothed ~~W~~DRO in action

Superquantile S_θ [Rockfellar *et al* '00] (a.k.a. Conditional Value-at-Risk)

Risk measure with dual formulation

$$R_\theta(x) = \max_{q \in \Delta_n} \left\{ \sum_{i=1}^n q_i \ell(y_i, \varphi(x, a_i)) : 0 \leq q_i \leq \frac{1}{n(1-\theta)} \right\}$$

DRO with (smoothed) superquantile in Pytorch

<https://github.com/krishnap25/sqwash>

```
import torch.nn.functional as F
from sqwash import reduce_superquantile

for x, y in dataloader:
    y_hat = model(x)
    batch_losses = F.cross_entropy(y_hat, y, reduction='none') # must set `reduction='none'`
    loss = reduce_superquantile(batch_losses, superquantile_tail_fraction=0.5) # Additional line
    loss.backward() # Proceed as usual from here
    ...
```

spotlight #3 : ~~W~~DRO for federated learning



Krishna Pillutla, Yassine Laguel, Jérôme Malick, Zaid Harchaoui

Federated Learning with Superquantile Aggregation for Heterogeneous Data

[Machine Learning Journal, 2023](#)

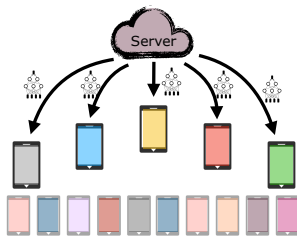
Federated learning in a nutshell

- Standard learning : get all the data and learn your model on it
- Efficient... but is privacy invasive (hospitals, european laws...)
- Idea : move the model not the data !

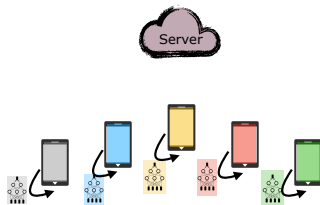
Federated learning in a nutshell

- Standard learning : get all the data and learn your model on it
- Efficient... but is privacy invasive (hospitals, european laws...)
- Idea : move the model not the data !
- Usual approach : FedAvg [McMahan *et al* 2017]
(based on old ideas, e.g. [Mangasarian 1995])

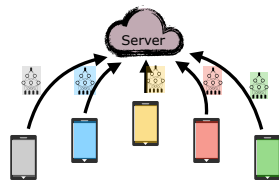
Step 1 of 3: Server broadcasts global model to sampled clients



Step 2 of 3: Clients perform some local SGD steps on their local data

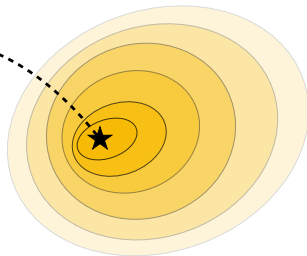
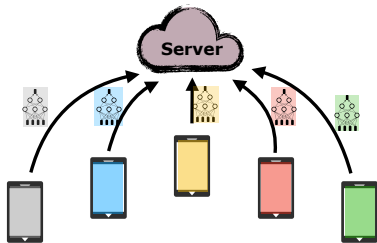


Step 3 of 3: Aggregate client updates securely



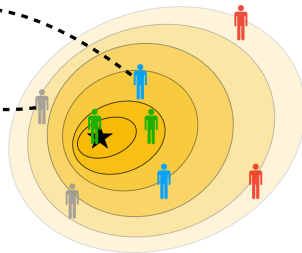
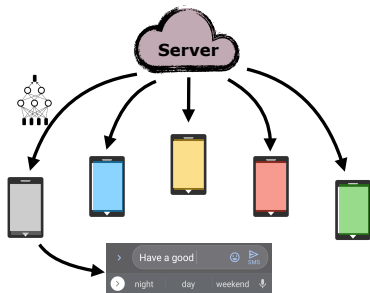
Issue of heterogeneous users

Global model is trained on *average distribution* across clients (ERM)



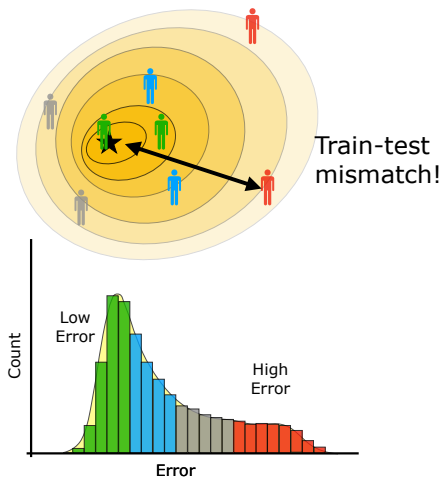
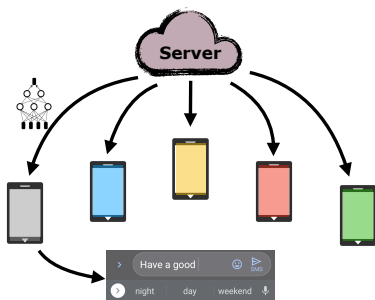
Issue of heterogeneous users

Global model is deployed on *individual* clients



Issue of heterogeneous users

Global model is deployed on *individual* clients

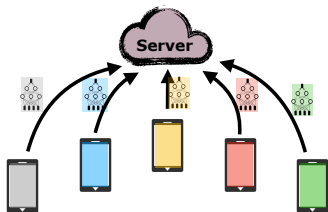


Our solution : superquantile minimization

ERM Algorithm (FedAvg):

$$\min_w \frac{1}{n} \sum_{i=1}^n F_i(w)$$

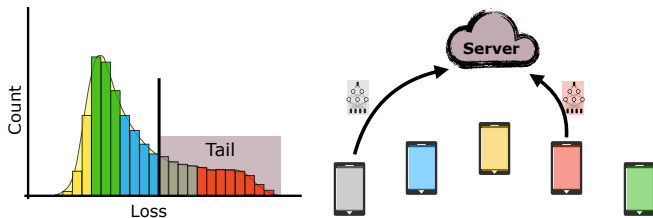
Step 3 of 3: Aggregate updates
contributed by **all clients**



Simplicial-FL Algorithm:

$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

Step 3 of 3: Aggregate updates
contributed by **tail clients** only



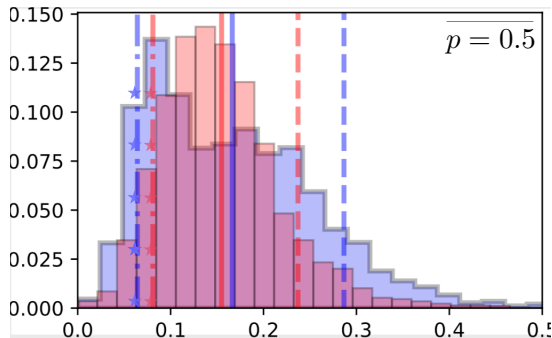
- Compatible with secure aggregation and differential privacy
- Analysis of the entropy-regularized version (both cvx and non-cvx)

Illustration

Classification task – ConvNet with EMNIST dataset (1730 users, 179 images/users)

Histogram over users of test misclassification error: **standard** vs. **DRO**

(dashed lines: 10%/90% -percentiles)



DRO reshapes test histograms

Conclusion

Main take-aways

- Distributionally robust optimization DRO is rich, active topic and has real-life applications, as in federated learning
- WDRO has nice generalization properties
- smoothed WDRO has nice properties
(general duality, approximation results, worst-case distribution, generalization)

On-going work

- Show that WDRO is not just a nice theory
- Further investigate applications... (in fairness?)

thank you all !