

# Coherent risk measures

Convex Analysis  $f: \mathbb{R}^d \rightarrow \mathbb{R}$

Lower semi-continuity (lsc)

$$\liminf_{x \rightarrow x} f(x) \geq f(x)$$

FACT Any lsc convex function is the supremum of all affine functions minorizing  $f$   
 i.e.  $f(x) = \sup_{a \leq f} a(x)$  where sup is over all affine  $a$ .

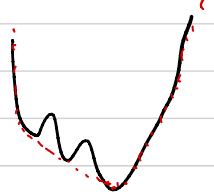


of. Proof of this uses the separating hyperplane theorem. Argues any point outside the epigraph has an affine minorant between it and the epigraph.

The Fenchel conjugate of  $f$  is  $f^*(s) := \sup_t \{ \langle s, t \rangle - f(t) \}$ .

The biconjugate is simply the conjugate of  $f^*$ , which we denote by  $f^{**}$ .

The biconjugate  $f^{**}$  is the largest lsc convex function minorizing  $f$ .



Lemma  $f^{**}(x) = \sup_{a \in \mathbb{R}^d, b \in \mathbb{R}} \{ \langle a, x \rangle - b : \langle a, t \rangle - b \leq f(t) \forall t \}$

Let  $A \subset \mathbb{R}^d \times \mathbb{R}$  be set of all  $(a, b)$  minorizing  $f$ :  $f(t) \geq \langle a, t \rangle - b \quad \forall t$ .

Then,  $(a, b) \in A \Leftrightarrow f(t) \geq \langle a, t \rangle - b \quad \forall t \Leftrightarrow b \geq \langle a, t \rangle - f(t) \quad \forall t \Leftrightarrow b \geq f^*(a), a \in \text{dom } f^*$   
 and  $\text{RHS} = \sup \{ \langle a, x \rangle - b : a \in \text{dom } f^*, b \leq f^*(a) \} = \sup \{ \langle a, x \rangle - f^*(a) \}$   $\square$

In particular, if  $f$  is convex, then  $f^{**} = f$ .  $\checkmark$  Fenchel-Moreau

## Coherence

Consider a risk measure that maps random losses to a risk value

$\mathcal{R}: \mathcal{L}^0(P) \rightarrow \mathbb{R}$ . The simplest such measure is  $\mathcal{R}(W) = \mathbb{E}W$ .

Coherence defines a class of "sensible" disutility functions.

Def  $\mathcal{R}$  is a coherent risk measure if

- 1) (Convexity)  $\mathcal{R}(\lambda W + (1-\lambda)W') \leq \lambda \mathcal{R}(W) + (1-\lambda) \mathcal{R}(W') \quad \forall \lambda \in [0, 1]$
- 2) (Monotonicity)  $W \leq W' \text{ P-a.s.} \Rightarrow \mathcal{R}(W) \leq \mathcal{R}(W')$
- 3) (Translation equivariance)  $\mathcal{R}(W+c) = \mathcal{R}(W) + c \quad \forall c \in \mathbb{R}$
- 4) (Positive homogeneity)  $\mathcal{R}(\lambda W) = \lambda \mathcal{R}(W) \quad \forall \lambda > 0$ .

Diminishing marginal returns  
 ↓ or caring more about high losses.

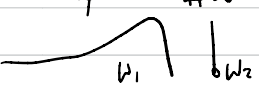
→ constant amt of deterioration in loss → same for risk

Example (semideviations)  $\mathcal{R}(W) = \mathbb{E}W + \left( \mathbb{E} (W - \mathbb{E}W)_+^2 \right)^{\frac{1}{2}}$  is coherent.

1, 3, 4 are obvious. To see 2,

$$\begin{aligned} \mathcal{R}(W) &\leq \mathbb{E}W + \left( \mathbb{E} (W' - \mathbb{E}W)_+^2 \right)^{\frac{1}{2}} = \mathbb{E}W + \left( \mathbb{E} (W' - \mathbb{E}W + \mathbb{E}W - \mathbb{E}W)_+^2 \right)^{\frac{1}{2}} \\ &= \mathbb{E}W + 2 \left( \mathbb{E} \left( \frac{1}{2} (W' - \mathbb{E}W) + \frac{1}{2} (\mathbb{E}W - \mathbb{E}W) \right)_+^2 \right)^{\frac{1}{2}} \leq \mathbb{E}W + \left( \mathbb{E} (W' - \mathbb{E}W)_+^2 \right)^{\frac{1}{2}} + \mathbb{E}W - \mathbb{E}W = \mathcal{R}(W'). \end{aligned}$$

cf.  $\mathbb{E}W + \left( \mathbb{E} (W - \mathbb{E}W)^2 \right)^{\frac{1}{2}}$  is NOT coherent since it also penalizes downward deviations of  $W$ .  
 then  $\mathcal{R}(W_1) \geq \mathcal{R}(W_2)$  may hold.



Rank  $W$  as DRD is not coherent. Why? (no monotonicity)

Example  $\mathcal{Q}(W) = \inf_{\gamma} \left\{ c \left( \mathbb{E}(W-\gamma)_+^{k_0} \right)^{\frac{1}{k_0}} + \gamma \right\}$  is coherent.

In fact, ANY DRO problem over convex  $\mathcal{Q} \subseteq \{Q: Q \ll P\}$ ,  $\mathcal{Q}(W) = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q W$  defines a coherent risk measure. This follows by the likelihood ratio representation

$$\mathcal{Q}(W) = \sup \left\{ \mathbb{E}_P L W : L = \frac{dQ}{dP}, Q \in \mathcal{Q} \right\}$$

We can show that the converse is also true. Recall the conjugate function

$$\mathcal{Q}^*(L) = \sup_{W \in \mathcal{L}^k(P)} \left\{ \mathbb{E}_P L W - \mathcal{Q}(W) \right\}.$$

Theorem For any coherent risk measure,  $\exists \mathcal{L} \subset \mathcal{L}^k(P)$  s.t.  $\mathcal{Q}(W) = \sup_{L \in \mathcal{L}} \mathbb{E}_P L W$ .

PF First note that since any finite-valued convex function is continuous,  $\mathcal{Q}^{**} = \mathcal{Q}$  by Fenchel-Moreau.

$$\text{So } \mathcal{Q}(W) = \sup_{L \in \mathcal{L}^k(P)} \left\{ \mathbb{E}_P L W - \mathcal{Q}^*(L) \right\} = \sup_{L \in \text{dom}(\mathcal{Q}^*)} \left\{ \mathbb{E}_P L W - \mathcal{Q}^*(L) \right\} \quad (*)$$

We now proceed in three parts.

① Property 2  $\Leftrightarrow \forall L \in \text{dom}(\mathcal{Q}^*), L \geq 0$  P-a.s.

PF ' $\Rightarrow$ ' Assume  $\exists L \in \text{dom}(\mathcal{Q}^*)$  s.t.  $L < 0$  on some set  $S$  of positive measure. Fix any  $W \in \mathcal{L}^k(P)$ .

$$\text{Then, } \mathcal{Q}^*(L) \geq \sup_{\lambda > 0} \left\{ \mathbb{E}_P L(W - \lambda \mathbb{1}_S) - \mathcal{Q}(W - \lambda \mathbb{1}_S) \right\} \geq \sup_{\lambda > 0} \left\{ \mathbb{E}_P L(W - \lambda \mathbb{1}_S) - \mathcal{Q}(W) \right\} = \infty.$$

since  $\mathcal{Q}(W) \leq \mathcal{Q}(W - \lambda \mathbb{1}_S)$  by Property 2.

' $\Leftarrow$ ' For any  $W \in \mathcal{W}$  P-a.s.,  $\mathbb{E}_P L W \leq \mathbb{E}_P L W'$  so (\*) gives the result.

② Property 3  $\Leftrightarrow \forall L \in \text{dom}(\mathcal{Q}^*), \mathbb{E}_P L = 1$

PF ' $\Rightarrow$ '  $\mathcal{Q}^*(L) \geq \sup_{c \in \mathbb{R}} \left\{ \mathbb{E}_P L(W+c) - \mathcal{Q}(W+c) \right\} = \sup_{c \in \mathbb{R}} \left\{ c(\mathbb{E}_P L - 1) + \mathbb{E}_P L W - \mathcal{Q}(W) \right\} = \infty$  if  $\mathbb{E}_P L \neq 1$ .

' $\Leftarrow$ ' Follows from (\*).

③ Property 4  $\Leftrightarrow \mathcal{Q}(W) = \sup_{L \in \text{dom}(\mathcal{Q}^*)} \mathbb{E}_P L W$ .

PF ' $\Rightarrow$ '  $\mathcal{Q}^*(L) = \sup_{W \in \mathcal{L}^k(P)} \left\{ \mathbb{E}_P L W - \mathcal{Q}(W) \right\} = \sup_{\lambda > 0, W \in \mathcal{L}^k(P)} \left\{ \mathbb{E}_P L(\lambda W) - \mathcal{Q}(\lambda W) \right\}$

$$= \sup_{\lambda > 0} \lambda \mathcal{Q}^*(L). \quad \text{So either } \mathcal{Q}^*(L) = 0 \text{ or } \infty.$$

' $\Leftarrow$ ' trivial. □

So DRO  $\Leftrightarrow$  risk-aversion (coherence). i.e., good perf. under distributional drifts  
 $\Leftrightarrow$  good tail-performance.

Qs • Which  $\mathcal{Q}$  to use? Or equivalently, which  $\mathcal{Q}$ ?

- ① Desired notion of risk-aversion
- ② statistical efficiency
- ③ Computational efficiency

- Risk-aversion generally means less sample efficiency
- Open Qs:
  - Linking problem structure to a particular kind of distr shift
  - ↳ this gives rough guidelines on how to choose  $\mathcal{Q}$ .

• So far, we considered shifts in full distr of  $\mathcal{Z}$ .

What if we are interested in partial distr shifts, only w.r.t. marginal distr of  $\mathcal{P}_{\text{param}}$ .

- Outliers
- Statistics in high-dim
- Training deep nets with risk-aversion

## Transition

and solve

So far, we studied DRO as a population formulation.  
 In practice, we would formulate an empirical approximation to this problem

$$\underset{\theta \in \Theta}{\text{minimize}} \quad \sup_{\mathcal{Q}: \mathcal{P}(\mathcal{Q}, \hat{P}_n) \leq \rho} \mathbb{E}_{\mathcal{Q}} \ell(\theta; \mathcal{Z}),$$

where  $\hat{P}_n$  is the empirical distribution,  $\frac{1}{n}$  uniform weights on each data point.

For  $f$ -divergence DRO, duality gives

$$\underset{\theta \in \Theta, \lambda > 0, \eta \in \mathbb{R}}{\text{minimize}} \quad \left\{ \frac{1}{n} \sum_{i=1}^n \hat{p}_i \cdot f^* \left( \frac{\ell(\theta; z_i) - \eta}{\lambda} \right) + \lambda \rho + \eta \right\}.$$

It turns out that by setting the radius  $\rho_n \rightarrow 0$  as  $n \rightarrow \infty$  above, we can view these finite sample procedures as approximations to the **average-case optimization problem**

i.e.  $\min_{\theta \in \Theta} \sup_{\mathcal{Q}: \frac{1}{n} \sum_{i=1}^n q_i \cdot \ell(\theta; z_i) \leq \rho_n} \sum_{i=1}^n q_i \cdot \ell(\theta; z_i)$  is a good approximation to  $\mathbb{E} \ell(\theta; \mathcal{Z})$

if we choose  $\rho_n$  appropriately.

This is what we will show now.

# Subexponential RVs & Bernstein bounds

So far, we studied tail-bounds for subGaussian RVs  $X$  satisfying  $\mathbb{E} e^{\lambda(X-\mathbb{E}X)} \leq e^{\frac{\lambda^2 \sigma^2}{2}} \quad \forall \lambda \in \mathbb{R}$ . This puts a restrictive condition on the tails of  $X$ , so it's natural to consider relaxations.

Def A RV  $X$  is sub-exponential with parameters  $(\nu, \alpha)$  if  $\mathbb{E} e^{\lambda(X-\mathbb{E}X)} \leq e^{\frac{\lambda^2 \nu^2}{2}}$  for all  $\lambda$  s.t.  $|\lambda| < \frac{1}{\alpha}$ . ↳ nonnegative

Any subGaussian RV is sub-exponential, but the converse is not true.

Ex Let  $X \sim N(0,1)$ . Then,  $\mathbb{E} e^{\lambda(X^2-1)} = \frac{1}{\sqrt{2\pi}} e^{-\lambda} \leq e^{2\lambda^2} \quad \forall |\lambda| \leq \frac{1}{4}$

By the Chernoff bound, we can again derive a tail-bound for sub-exponential RVs

$$\mathbb{P}(X-\mathbb{E}X \geq t) \leq e^{-\lambda t} \mathbb{E} e^{\lambda(X-\mathbb{E}X)} \leq \exp\left(-\lambda t + \frac{\lambda^2 \nu^2}{2}\right) \quad \text{if } 0 < \lambda < \frac{1}{\alpha}.$$

Unconstrained min of  $-\lambda t + \frac{\lambda^2 \nu^2}{2}$ :  $\lambda^* = \frac{t}{\nu^2}$ . This gives RHS =  $-\frac{t^2}{2\nu^2}$  if  $\frac{t}{\nu^2} < \frac{1}{\alpha}$ .

If  $\frac{t}{\nu^2} \geq \frac{1}{\alpha}$ ,  $\lambda \mapsto -\lambda t + \frac{\lambda^2 \nu^2}{2}$  is monotone, so  $\lambda^* = \frac{1}{\alpha}$  gives RHS =  $-\frac{t}{\alpha} + \frac{1}{2\alpha} \frac{\nu^2}{\alpha} \leq -\frac{t}{\alpha} + \frac{1}{2\alpha} t = -\frac{t}{2\alpha}$ .

Combining,  $\mathbb{P}(X-\mathbb{E}X \geq t) \leq \begin{cases} \exp(-t^2/2\nu^2) & \text{if } 0 \leq t \leq \frac{\nu^2}{\alpha} \\ \exp(-t/2\alpha) & \text{if } t > \frac{\nu^2}{\alpha} \end{cases}$

(Similarly, left tail can be controlled.)

(Bernstein's condition) Let  $\sigma^2 = \text{Var } X$ .  $\forall k \geq 2$ ,  $|\mathbb{E}(X-\mathbb{E}X)^k| \leq \frac{1}{2} k! \sigma^2 b^{k-2}$

e.g. If  $|X-\mu| \leq b$ , B's condition holds.

If  $X$  satisfies B's condition, it is sub-exponential.

$$\mathbb{E} e^{\lambda(X-\mathbb{E}X)} = \mathbb{E} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} (X-\mu)^k = 1 + \frac{\lambda^2 \sigma^2}{2} + \sum_{k=3}^{\infty} \lambda^k \frac{1}{k!} \mathbb{E}(X-\mathbb{E}X)^k \leq 1 + \frac{\lambda^2 \sigma^2}{2} + \sum_{k=3}^{\infty} \frac{\lambda^k \sigma^2}{2} \cdot (|\lambda| \cdot b)^{k-2}$$

For any  $|\lambda| < 1/b$ , RHS =  $1 + \frac{\lambda^2 \sigma^2}{2} \cdot \frac{1}{1-b|\lambda|} \leq \exp\left(\frac{\lambda^2 \sigma^2}{2} \cdot \frac{1}{1-b|\lambda|}\right)$  by  $1+t \leq e^t$ .  $\dots$  (\*)

So for  $|\lambda| < 1/2b$ ,  $\mathbb{E} e^{\lambda(X-\mathbb{E}X)} \leq \exp(\lambda^2 \sigma^2)$ .

i.e. B's condition  $\Rightarrow$  sub-exponential with  $(\sqrt{2}\sigma, 2b)$ .

Directly applying condition (\*) in the Chernoff bound,  
 $P(X - \mathbb{E}X \geq t) \leq e^{-\lambda t} \mathbb{E} e^{\lambda(X - \mathbb{E}X)} \leq \exp\left(-\lambda t + \frac{1}{1 - b|\lambda|} \cdot \frac{\lambda^2 \sigma^2}{2}\right)$ . Let  $\lambda = \frac{t}{bt + \sigma^2} \in (0, b^{-1})$ ,

$$\begin{aligned} \text{RHS} &= -\frac{t^2}{bt + \sigma^2} + \frac{1}{1 - \frac{bt}{bt + \sigma^2}} \cdot \frac{\sigma^2}{2} \cdot \frac{t^2}{(bt + \sigma^2)^2} \\ &= -\frac{t^2}{bt + \sigma^2} + \frac{bt + \sigma^2}{\sigma^2} \cdot \frac{\sigma^2}{2} \cdot \frac{t^2}{(bt + \sigma^2)^2} = -\frac{t^2}{2(bt + \sigma^2)} \end{aligned}$$

We have the following lemma.

Lemma If  $\mathbb{E} \exp(\lambda(X - \mathbb{E}X)) \leq \exp\left(\frac{\lambda^2 \sigma^2}{2} \cdot \frac{1}{1 - b|\lambda|}\right)$  for  $|\lambda| < 1/b$ , then  
 $P(X - \mathbb{E}X \geq t) \leq \exp\left(-\frac{t^2}{2(bt + \sigma^2)}\right)$ .

Example Recalling that if  $|X - \mathbb{E}X| \leq b$  then  $X - \mathbb{E}X$  is subG with param  $b^2$   
the tail-bound for subG RVs gives  $P(X - \mathbb{E}X \geq t) \leq \exp\left(-\frac{t^2}{2b^2}\right)$   
 $\equiv X - \mathbb{E}X \leq b\sqrt{2s}$  w.p.  $\geq 1 - e^{-s}$

Let us try using the above lemma. Set  $s = \frac{t^2}{2(bt + \sigma^2)}$ . Solving for  $t$ ,  $t = bs + \sqrt{b^2 s^2 + 2\sigma^2 s}$ .  
So  $X - \mathbb{E}X \leq \sqrt{2\sigma^2 s} + 2bs$  w.p.  $\geq 1 - e^{-s}$ , where we used  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$   $\forall a, b > 0$ .

When  $\sigma^2 \ll b^2$  which is often true when  $X$  occasionally take large values,  
the second bound is better.

Example Consider i.i.d.  $X_i \in [-M, M]$ , so that (\*) holds with  $b = M$ .

$$\begin{aligned} \mathbb{E} \exp\left(\lambda \left(\frac{1}{n} \sum X_i - \mathbb{E}X\right)\right) &\leq \prod_{i=1}^n \mathbb{E} \exp\left(\frac{\lambda}{n} (X_i - \mathbb{E}X)\right) \leq \prod_{i=1}^n \exp\left(\frac{\lambda^2 \sigma^2}{2n^2} \cdot \frac{1}{1 - M|\lambda|/n}\right) \quad \text{for } |\lambda| < \frac{n}{M} \\ &= \exp\left(\frac{\lambda^2 \sigma^2}{2n} \cdot \frac{1}{1 - M|\lambda|/n}\right). \end{aligned}$$

So lemma gives  $P\left(\frac{1}{n} \sum X_i - \mathbb{E}X \geq t\right) \leq \exp\left(-\frac{nt^2}{2(\sigma^2 + Mt)}\right)$ .

We arrive at the following bound:

$$\mathbb{E}X \leq \frac{1}{n} \sum X_i + \sqrt{\frac{2s \sigma^2}{n}} + 2M \cdot \frac{s}{n} \quad \text{w.p. } \geq 1 - e^{-s}$$

If  $\text{Var} X$  concentrates around the sample variance  $\hat{\sigma}_n^2$ , then

$$\mathbb{E}X \leq \frac{1}{n} \sum X_i + \sqrt{2s \cdot \hat{\sigma}_n^2 / n} + C \cdot M \frac{s}{n} \quad \text{w.h.p.}$$

# Variance regularization

Given Bernstein's inequality we just derived, we have

$$\mathbb{E} l(\theta; z) \leq \frac{1}{n} \sum l(\theta; z_i) + \sqrt{\frac{cs}{n} \widehat{\text{Var}}(l(\theta; z))} + \frac{cMs}{n} \quad \text{w.p.} \geq 1 - e^{-s}$$

for some numerical constant  $c > 0$ .

So given this upper bound on the population loss, it is natural to optimize

$$\underset{\theta \in \Theta}{\text{minimize}} \left\{ \frac{1}{n} \sum l(\theta; z_i) + c \cdot \sqrt{\frac{\widehat{\text{Var}}(l(\theta; z_i))}{n}} \right\}$$

↳ We call this the variance regularized problem.

Problem  $\theta \mapsto \sqrt{\widehat{\text{Var}} l(\theta; z)}$  is nonconvex even when  $\theta \mapsto l(\theta; z)$  is convex.

Consider the empirical  $f$ -divergence DRO formulation with  $f(t) = \frac{1}{2}(t-1)^2$ , and  $\rho_n > 0$ .

$$\begin{aligned} & \max_Q \left\{ \mathbb{E}_Q l(\theta; z) : D_f(Q, \hat{P}_n) = \mathbb{E}_{\hat{P}_n} \frac{1}{2} \left( \frac{dQ}{d\hat{P}_n} - 1 \right)^2 = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \left( \frac{q_i}{1/n} - 1 \right)^2 \leq \rho_n \right\} \\ & = \max_{q \geq 0} \left\{ \sum_{i=1}^n q_i l(\theta; z_i) : \frac{1}{2n} \sum_{i=1}^n (nq_i - 1)^2 \leq \rho_n, q^T \mathbf{1} = 1 \right\}. \end{aligned}$$

Taking  $\rho_n = \frac{c}{n}$ , we show that is actually equivalent to Variance Reg.

Fix  $\theta \in \Theta$ , and let  $W = l(\theta; z)$ ,  $\hat{\mu}_n := \frac{1}{n} \sum_{i=1}^n w_i$ ,  $\hat{\sigma}_n^2 := \frac{1}{n} \sum_{i=1}^n w_i^2 - \left( \frac{1}{n} \sum_{i=1}^n w_i \right)^2$ .

Define  $Q_n(c) := \left\{ q \in \mathbb{R}_+^n : q^T \mathbf{1} = 1, \frac{1}{2} \sum (nq_i - 1)^2 = \frac{1}{2} \|nq - \mathbf{1}\|_2^2 \leq c \right\}$

Doing a change of variables  $u = nq - \mathbf{1}$ , and denoting  $\vec{w}_n := [w_1, \dots, w_n] \in \mathbb{R}^n$ ,

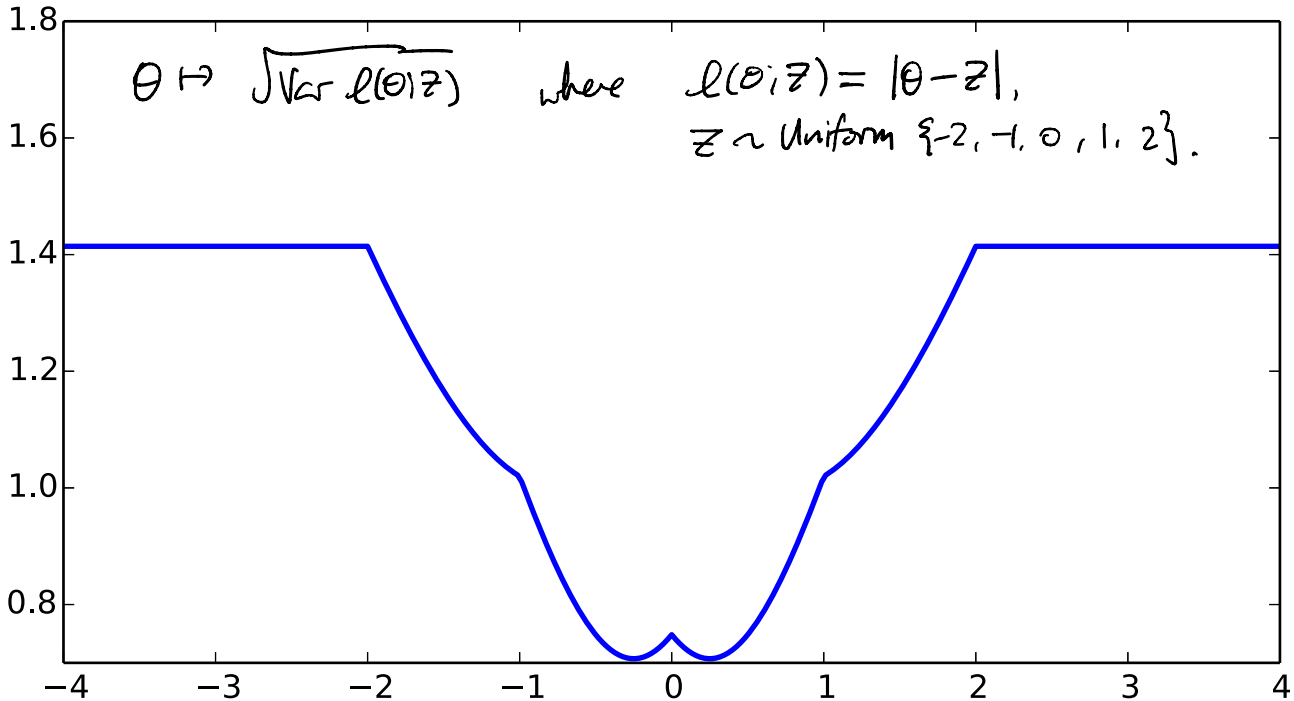
$$\sup_{q \in Q_n} q^T \vec{w}_n = \hat{\mu}_n + \sup \left\{ u^T (\vec{w}_n - \hat{\mu}_n \mathbf{1}) : u \geq -\frac{1}{n} \mathbf{1}, u^T \mathbf{1} = 0, \|u\|_2^2 \leq \frac{2c}{n^2} \right\}.$$

From Cauchy-Schwarz,  $u^T (\vec{w}_n - \hat{\mu}_n \mathbf{1}) \leq \|u\|_2 \|\vec{w}_n - \hat{\mu}_n \mathbf{1}\|_2 \leq \frac{\sqrt{2c}}{n} \|\vec{w}_n - \hat{\mu}_n \mathbf{1}\|_2 = \sqrt{\frac{2c \hat{\sigma}_n^2}{n}}$

Equality is attained iff  $u \propto \vec{w}_n - \hat{\mu}_n \mathbf{1}$ , or  $u_i^* = \frac{\sqrt{2c} (w_i - \hat{\mu}_n)}{n \|\vec{w}_n - \hat{\mu}_n \mathbf{1}\|_2} = \frac{\sqrt{2c} (w_i - \hat{\mu}_n)}{n \sqrt{\hat{\sigma}_n^2}}$ .

Such choice is possible if  $u_i^* \geq -\frac{1}{n} \forall i$ .

$\theta \mapsto \sqrt{\text{Var } \ell(\theta; Z)}$  where  $\ell(\theta; Z) = |\theta - Z|$ ,  
 $Z \sim \text{Uniform } \{-2, -1, 0, 1, 2\}$ .



So as long as  $\sqrt{2c} \min_{1 \leq i \leq n} (W_i - \hat{\mu}_n) \geq -\sqrt{n} \hat{\sigma}_n^2 \equiv \hat{\sigma}_n \geq \sqrt{\frac{2c}{n}} \cdot \max_{1 \leq i \leq n} (\hat{\mu}_n - W_i)$ ,

$$\sup_{q \in Q_n(c)} q^T \vec{W}_n = \hat{\mu}_n + \sqrt{\frac{2c}{n}} \hat{\sigma}_n^2$$

By arguing that  $\hat{\sigma}_n^2$  is large enough with high probability, we arrive at the following.

Theorem Let  $W \in [0, M]$  a.s., and  $\sigma^2 = \text{Var} W > 0$

$$\left( \sqrt{\frac{2c}{n}} \hat{\sigma}_n^2 - \frac{2Mc}{n} \right)_+ \leq \sup_{q \in Q_n(c)} q^T W - \hat{\mu}_n \leq \sqrt{\frac{2c}{n}} \hat{\sigma}_n^2$$

and for  $n \geq \max(2, \frac{M^2}{\sigma^2} \max(80, 44))$ , w.p.  $\geq 1 - \exp(-\frac{3nb^2}{5M^2})$

$$\sup_{q \in Q_n(c)} q^T W = \hat{\mu}_n + \sqrt{\frac{2c}{n}} \hat{\sigma}_n^2$$

That is, we have shown

$$\frac{1}{n} \sum l(\theta_i; z_i) + c \cdot \sqrt{\frac{\hat{\text{var}}(l(\theta_i; z_i))}{n}} \stackrel{*}{=} \sup_{q \in Q_n(c)} \sum_i q_i l(\theta_i; z_i) \quad \text{w.h.p.}$$

LHS: nonconvex, not coherent, but a natural quantity from a learning theoretic perspective

RHS: convex, coherent, "robust" w.r.t. reweighting.

A computationally tractable way of regularizing by variance.

Deep connections to EL

- Remark
- 1) The guarantee  $*$  can be made uniform in  $\theta \in \Theta$
  - 2) A variant can be shown for any smooth  $f$ -divergence
  - 3) An asymptotic version can be shown for any fast mixing seq of RV

So what does this give us?

upweight harder examples

$$\hat{\theta}_n^{\text{rob}} = \underset{\theta \in \Theta}{\text{argmin}} \sup_{q \in Q_n(c)} \sum_i q_i l(\theta_i; z_i)$$

$$\hat{\theta}_n^{\text{ERM}} = \underset{\theta \in \Theta}{\text{argmin}} \frac{1}{n} \sum_i l(\theta_i; z_i)$$



Theorem Let  $c = s + C \cdot \mathbb{P}_n \{ l(\theta; z) : \theta \in \Theta \}$ . If  $l(\theta; z) \in [0, M]$  a.s.,  
 $\mathbb{E}_{z \sim P} l(\hat{\theta}_n^{\text{rob}}; z) \leq \inf_{\theta \in \Theta} \mathbb{E}_P l(\theta; z) + 2 \sqrt{\frac{2c}{n} \text{Var} l(\theta; z)} + \frac{cM}{n} \cdot a$   
 w.p.  $\geq 1 - e^{-s}$ , for some numerical const  $a > 0$

Optimal bias-variance trade-off

Using a uniform version of Bernstein's inequality, w.p.  $\geq 1 - e^{-s}$ ,  
 $\mathbb{E}_{z \sim P} l(\hat{\theta}_n^{\text{ERM}}; z) \leq \inf_{\theta \in \Theta} \mathbb{E}_P l(\theta; z) + \sqrt{\frac{2cM}{n} \inf_{\theta \in \Theta} \mathbb{E}_P l(\theta; z)} + \frac{dMc}{n}$

If  $\text{Var} l(\theta; z) \ll M \cdot \mathbb{E}_P l(\theta; z)$ , then bound for  $\hat{\theta}_n^{\text{rob}}$  is tighter.

of. We can also construct an explicit (contrived) example where

$$\mathbb{E}_{z \sim P} l(\hat{\theta}_n^{\text{rob}}; z) \leq \inf_{\theta \in \Theta} \mathbb{E}_P l(\theta; z) + \frac{a}{n} \quad \text{but} \quad \mathbb{E}_{z \sim P} l(\hat{\theta}_n^{\text{ERM}}; z) \geq \inf_{\theta \in \Theta} \mathbb{E}_P l(\theta; z) + \frac{a'}{\sqrt{n}}$$

## Wasserstein DRO & Regularization

By choosing certain cost functions, we can show that Wasserstein DRO is equivalent to classical regularizers.

Proposition (Regression) Consider  $c((x, y), (x', y')) = \| (x, y) - (x', y') \|_k^2$  for some  $k \in (1, \infty]$ .

$$\sup_{Q: W_k(Q, \hat{P}) \leq \rho} \mathbb{E}_Q (Y - \theta^T X)^2 = \left( \left( \frac{1}{n} \sum_{i=1}^n (Y_i - \theta^T X_i)^2 \right)^{\frac{1}{2}} + \sqrt{\rho} \| [\theta, 1] \|_k \right)^2$$

where  $k_* = \frac{k}{k-1}$  so  $\frac{1}{k} + \frac{1}{k_*} = 1$ .

Proof of main proposition To ease notation, define  $z = (x, y)$ ,  $\bar{\theta} = [\theta, 1] \in \mathbb{R}^{d+1}$

From the duality result for Wasserstein DRO,

$$\sup_{Q: W_k(Q, \hat{P}) \leq \rho} \mathbb{E}_Q (Y - \theta^T X)^2 = \inf_{\lambda \geq 0} \left\{ \lambda \ell + \mathbb{E}_{\hat{P}} \sup_{z'} \left\{ (\bar{\theta}^T z')^2 - \lambda \| z - z' \|_k^2 \right\} \right\}$$

We just simplify the robust surrogate loss  $\sup_{z'} \left\{ (\bar{\theta}^T z')^2 - \lambda \| z - z' \|_k^2 \right\} = \phi_\lambda(\bar{\theta}; z)$

Doing a change of variable  $\Delta = z - z'$ . since sup should be attained when signs match.

$$\phi_\lambda(\bar{\theta}; z) = \sup_{\Delta} \left\{ (\bar{\theta}^\top z + \bar{\theta}^\top \Delta)^2 - \lambda \|\Delta\|_k^2 \right\} = \sup_{\Delta} \left\{ (\bar{\theta}^\top z + \text{sgn}(\bar{\theta}^\top z) \cdot \bar{\theta}^\top \Delta)^2 - \lambda \|\Delta\|_k^2 \right\}$$

$$= \sup_{\Delta} \left\{ (\bar{\theta}^\top z + \text{sgn}(\bar{\theta}^\top z) \cdot \|\bar{\theta}\|_{k_*} \|\Delta\|_k)^2 - \lambda \|\Delta\|_k^2 \right\} \quad \text{since Holder's inequality is tight for some choice of } \Delta$$

$$= (\bar{\theta}^\top z)^2 + \sup_{\Delta} \left\{ -(\lambda - \|\bar{\theta}\|_{k_*}) \cdot \|\Delta\|_k^2 + 2|\bar{\theta}^\top z| \|\bar{\theta}\|_{k_*} \|\Delta\|_k \right\}$$

$$= \begin{cases} \frac{\lambda}{\lambda - \|\bar{\theta}\|_{k_*}^2} (\bar{\theta}^\top z)^2 & \text{if } \lambda > \|\bar{\theta}\|_{k_*}^2 \\ \infty & \text{o/w} \end{cases}$$

So we conclude  $\sup_{\mathcal{Q}: W_0(\mathcal{Q}, \hat{\theta}) \leq \rho} \mathbb{E}_{\mathcal{Q}} (Y - \theta^\top X)^2 = \inf_{\lambda > \|\bar{\theta}\|_{k_*}^2} \left\{ \lambda \rho + \frac{\lambda}{\lambda - \|\bar{\theta}\|_{k_*}^2} \frac{1}{n} \sum_{i=1}^n (\bar{\theta}^\top z_i)^2 \right\}$ .

Noting that the optimum is achieved at  $\lambda^* = \|\bar{\theta}\|_{k_*}^2 + \left( \frac{\|\bar{\theta}\|_{k_*}^2}{\rho} \frac{1}{n} \sum_{i=1}^n (\bar{\theta}^\top z_i)^2 \right)^{\frac{1}{2}}$ , we have the result.  $\square$

Similarly, we can consider only perturbing the feature vector by setting

$$c((x, y), (x', y')) = \begin{cases} \|x - x'\|_k^2 & \text{if } y = y' \\ \infty & \text{if } y \neq y' \end{cases} \quad \text{“covariate shift”}$$

$$\sup_{\mathcal{Q}: W_0(\mathcal{Q}, \hat{\theta}) \leq \rho} \mathbb{E}_{\mathcal{Q}} (Y - \theta^\top X)^2 = \left( \left( \frac{1}{n} \sum_{i=1}^n (Y_i - \theta^\top X_i)^2 \right)^{\frac{1}{2}} + \sqrt{\rho} \|\theta\|_{k_*} \right)^2$$

We can show a similar equivalence for linear classification models.  $Y \in \{\pm 1\}$

$$c((x, y), (x', y')) = \begin{cases} \|x - x'\|_k & \text{if } y = y' \\ \infty & \text{if } y \neq y' \end{cases}$$

$$\sup_{\mathcal{Q}: W_0(\mathcal{Q}, \hat{\theta}) \leq \rho} \mathbb{E}_{\mathcal{Q}} \log(1 + e^{-Y \theta^\top X}) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-Y_i \theta^\top X_i}) + \rho \|\theta\|_{k_*}$$

$$\sup_{\mathcal{Q}: W_0(\mathcal{Q}, \hat{\theta}) \leq \rho} \mathbb{E}_{\mathcal{Q}} (1 - Y \theta^\top X)_+ = \frac{1}{n} \sum_{i=1}^n (1 - Y_i \theta^\top X_i)_+ + \rho \|\theta\|_{k_*}$$

## Another basic connection

Consider the cost function  $c(z, z') = \frac{1}{2} \|z - z'\|_2^2$ , and the corresponding robust surrogate function  $\phi_\lambda(\theta; z) = \sup_{z'} \{ \ell(\theta; z') - \frac{\lambda}{2} \|z - z'\|_2^2 \}$ .

Plugging the first order approximation into the robust surrogate

$$\# = \sup_{z'} \{ \ell(\theta; z) + \nabla_z \ell(\theta; z)^T (z' - z) - \frac{\lambda}{2} \|z - z'\|_2^2 \} \quad \text{gradient ascent step for inc. loss}$$

The max is attained at  $\nabla_z \ell(\theta; z) = \lambda(z - z') \Rightarrow z' = z + \frac{1}{\lambda} \nabla_z \ell(\theta; z)$ ,

and  $\# = \ell(\theta; z) + \frac{1}{2\lambda} \|\nabla_z \ell(\theta; z)\|_2^2$ .

Plugging this first order approximation in to the dual, we get

$$\inf_{\lambda \geq 0} \left\{ \lambda \rho + \frac{1}{n} \sum_{i=1}^n \sup_{z'} \left\{ \ell(\theta; z_i) + \nabla_z \ell(\theta; z_i)^T (z' - z_i) - \frac{\lambda}{2} \|z_i - z'\|_2^2 \right\} \right\}$$

$$= \frac{1}{n} \sum \ell(\theta; z_i) + \inf_{\lambda \geq 0} \left\{ \lambda \rho + \frac{1}{2\lambda} \frac{1}{n} \sum \|\nabla_z \ell(\theta; z_i)\|_2^2 \right\}$$

$$= \frac{1}{n} \sum \ell(\theta; z_i) + \sqrt{\rho \cdot \left( \mathbb{E}_{\mathbb{P}_n} \|\nabla_z \ell(\theta; z)\|_2^2 \right)^{\frac{1}{2}}}$$

↳ Regularize to make the loss more stable against data perturbations

For smooth losses, Wasserstein DRO regularizes to make  $\|\nabla_z \ell(\theta; z)\|$  small, up to first order.

# Experiment 2: Fine-grained recognition

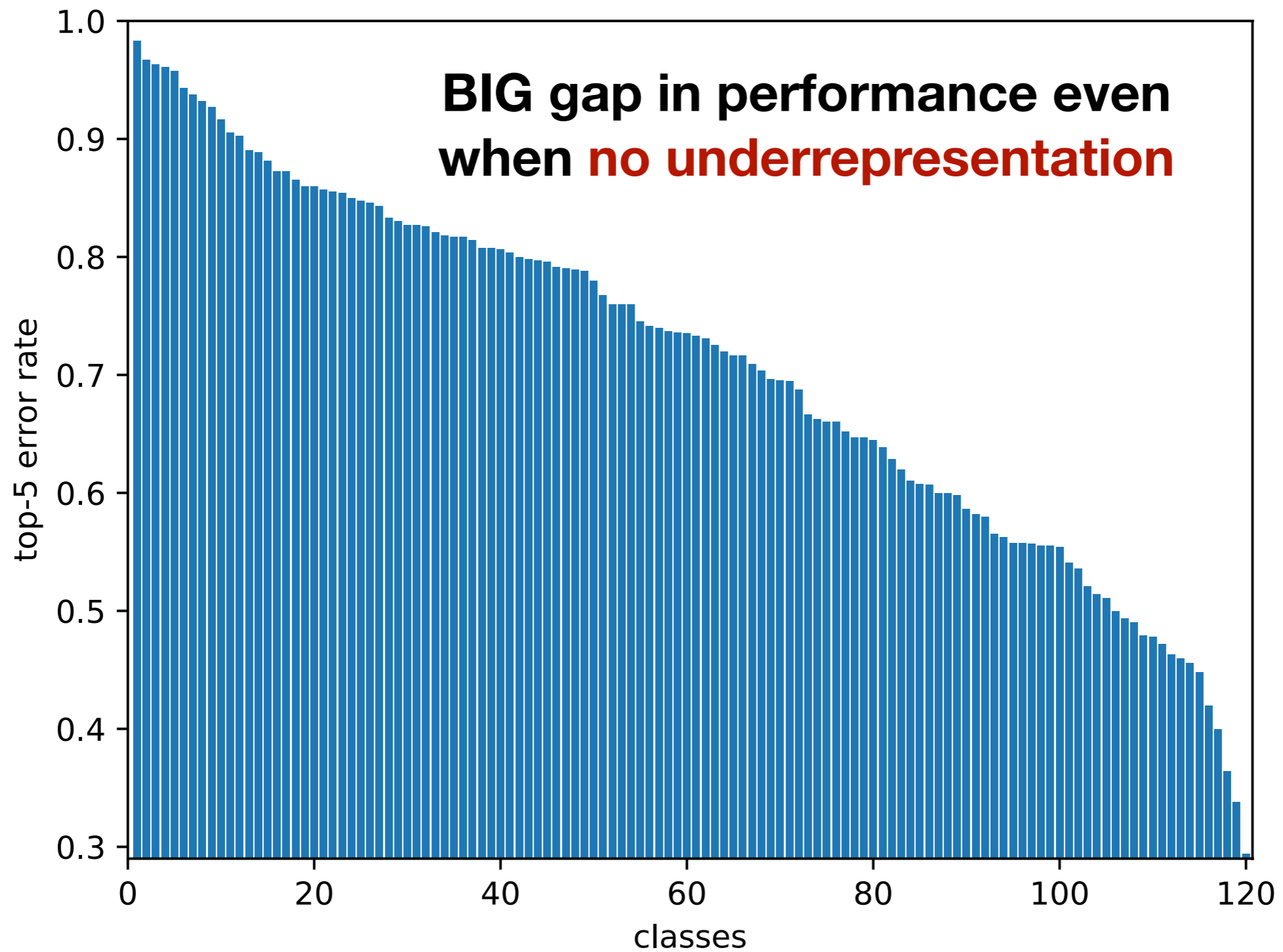
- Task: classify image of dog to breed (120 classes)
- Kernel features



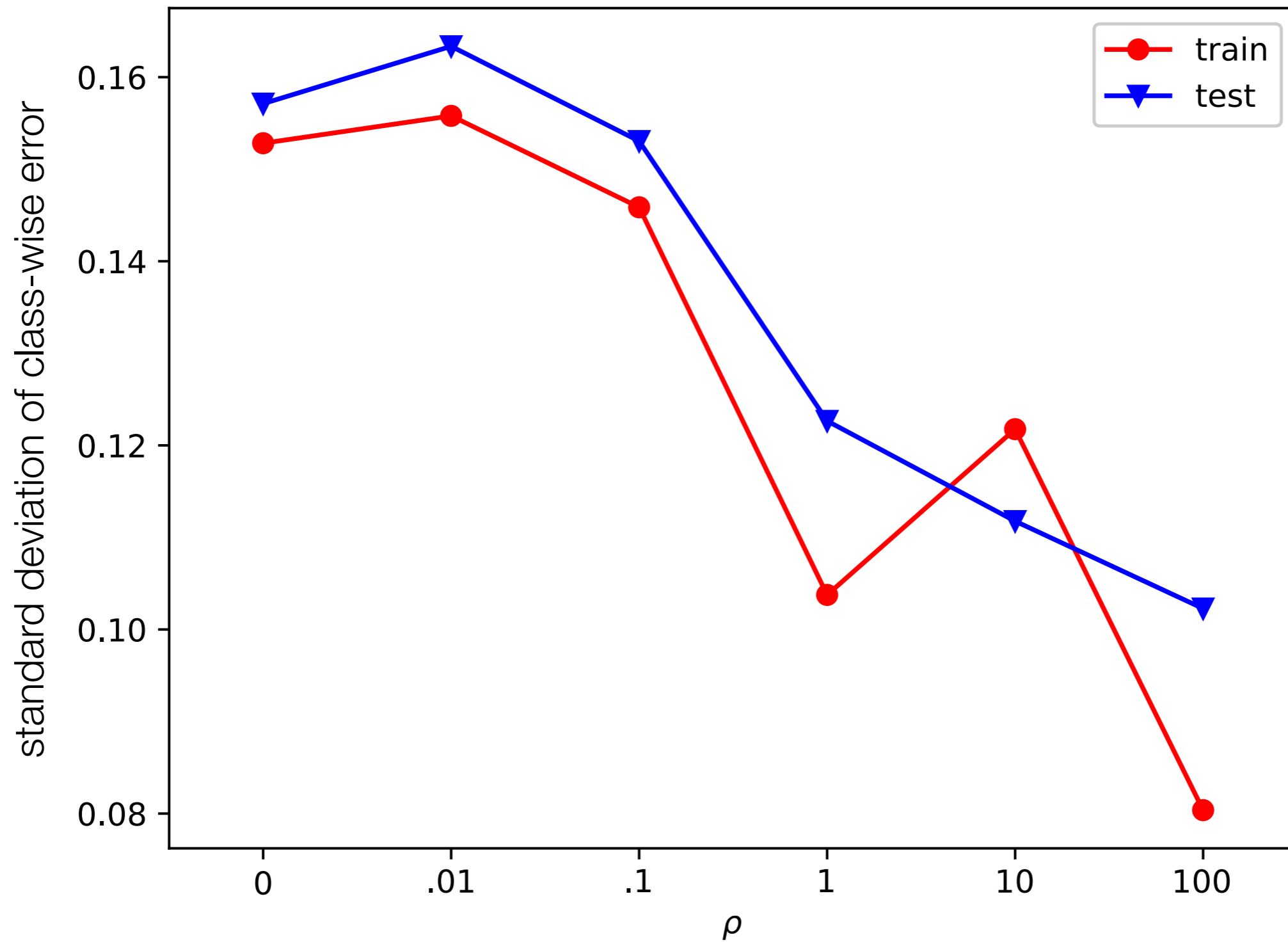
Stanford Dogs Dataset [Khosla et al. '11]

**No underrepresentation:**  
same number of images per class

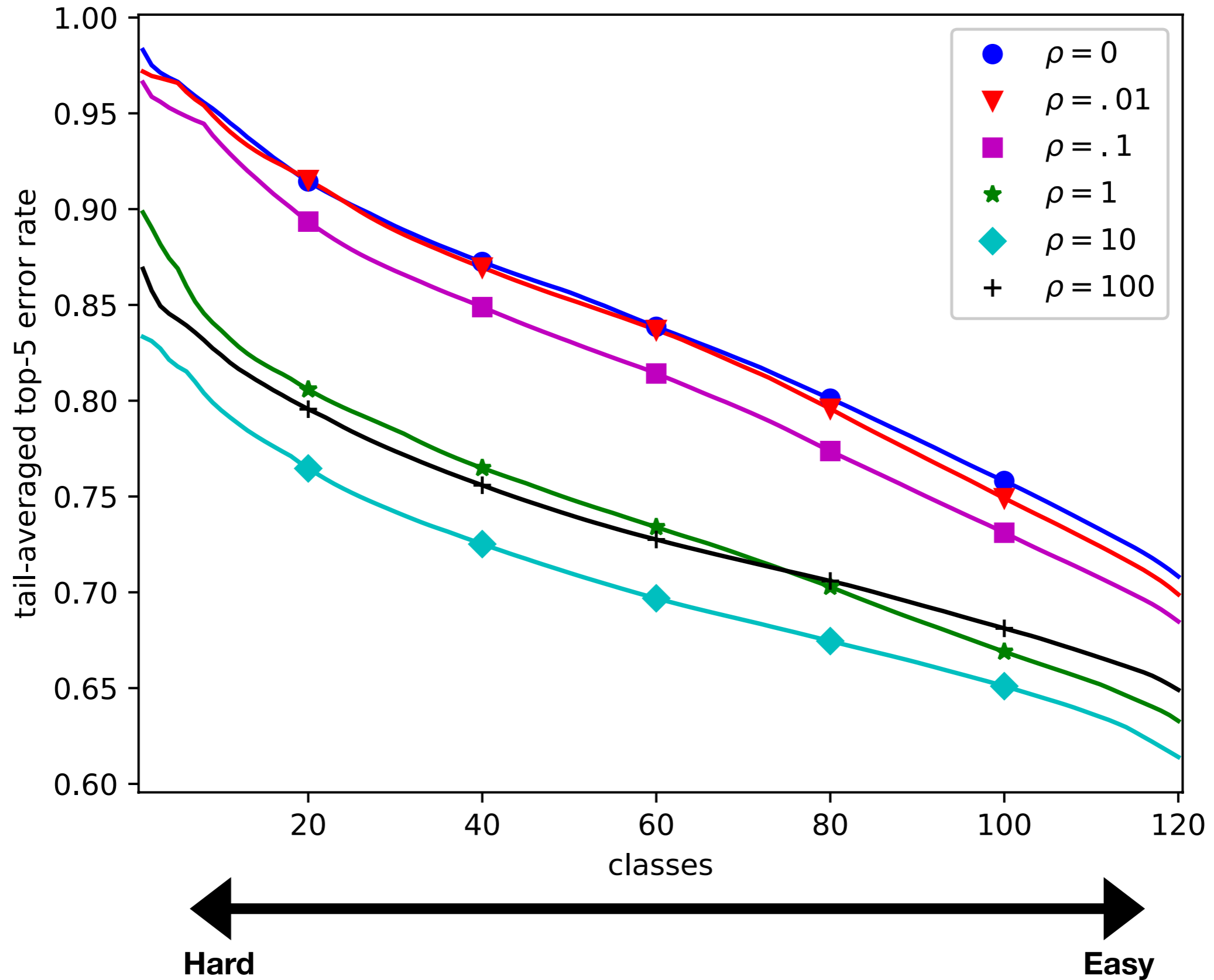
# ERM error rate



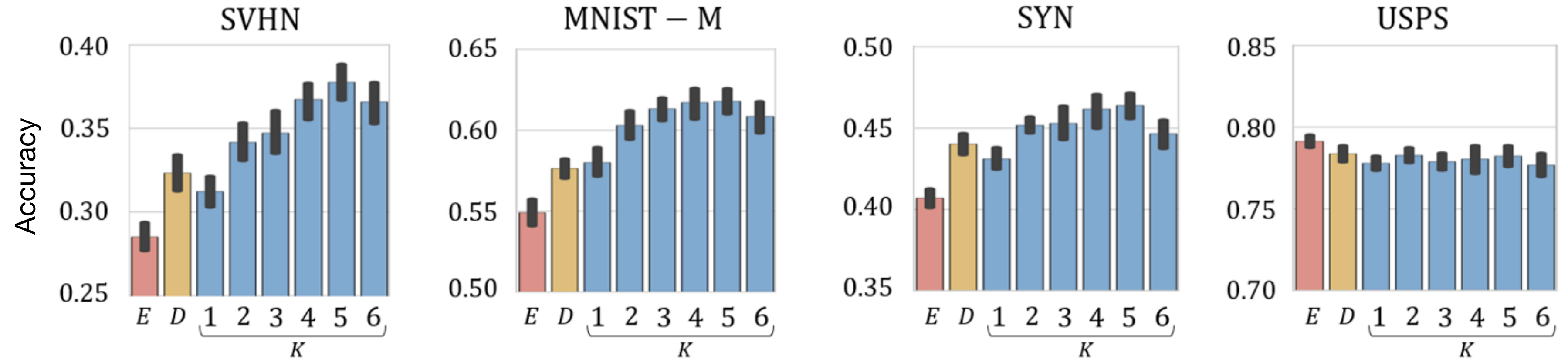
# Variation in error over 120 class



# Worst x-classes



# 30 seconds demo of Wasserstein DRO



E = ERM with L2 regularization

D = Dropout regularization

K = number of Wasserstein DRO gradient ascent steps

$$\begin{aligned} \phi_{\lambda}(\theta; NN(x)) \\ = \sup_{x'} \{ \ell(\theta; y, NN(x')) - \lambda \|NN(x) - NN(x')\|_2^2 \} \end{aligned}$$

Trained using lambda = 1.0, and an adaptive cost function defined on last hidden layer outputs of the neural network