

# A Discrepancy-Based Theory of Adaptation

Joint work with:

Pranjal Awasthi (Google Research), Corinna Cortes (Google Research), Andres Muñoz Medina (Google), Yishay Mansour (Google Research & Tel-Aviv), Afshin Rosamizadeh (Google Research)

MEHRYAR MOHRI    MOHRI@

GOOGLE RESEARCH & COURANT INSTITUTE

# Motivation: Domain Adaptation

- Distribution mismatch: in many real-world problems, source and target domains differ.
- Challenges: collecting labeled data for target domains is costly, generalization problem.
- Special instances: sample bias correction, covariate-shift problems, fine-tuning for LLMs.
- Real-world applications: healthcare, autonomous driving, speech recognition, best-effort fairness.
- Can we design a theoretical framework to guide adaptation methods?

# This Talk

- Discrepancy.
- Reweighting algorithms.
- Experimental results.

# Multiple-Source Adaptation

- Multiple-source adaptation problem: no labeled data.
  - Theoretical analysis (Mansour, MM, and Rostamizadeh, 2008, 2009).
  - Theory and algorithms (Hoffman, MM, and Zhang, 2021), (Cortes, MM, Suresh, Zhang, 2021).
- Learning with multiple source distribution: labeled data.
  - Theoretical analysis and algorithms, application to federated learning (MM, Sivek, and Suresh, 2019).
  - Boosting with multiple sources (Cortes, MM, Storcheus, Suresh, 2021).
  - Limited target data (Mansour, MM, Ro, Suresh, Wu, 2021).

# Adaptation Scenario

- Input space  $\mathcal{X}$ , output space  $\mathcal{Y}$ .
- Loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$ .
- Hypothesis set  $\mathcal{H}$  of functions mapping from  $\mathcal{X}$  to  $\mathcal{Y}$ .
- Learner receives:
  - Labeled sample from source domain, distribution  $\mathcal{Q}$ .
  - Labeled points from target domain, distribution  $\mathcal{P}$ :  
supervised adaptation (fair amount), weakly supervised (only some), unsupervised (none).
  - Typically large unlabeled sample from  $\mathcal{P}$ .

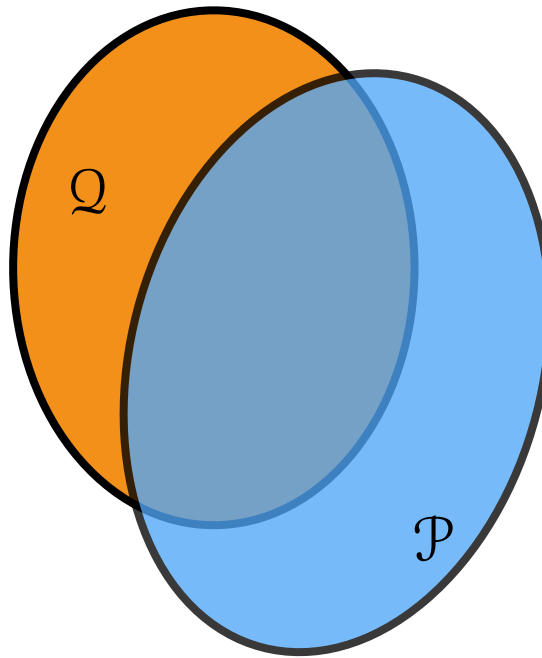
# Adaptation Problem

## ■ Learning problem:

- Use labeled samples from  $\mathcal{Q}$  and  $\mathcal{P}$  (different scenarios) as well as typically large unlabeled sample from  $\mathcal{P}$  to find hypothesis  $h \in \mathcal{H}$  with small target expected loss

$$\mathcal{L}(\mathcal{P}, h) = \mathbb{E}_{(x,y) \sim \mathcal{P}} [\ell(h(x), y)].$$

# Challenging Problem



- Which divergence between distributions should we use?

# Divergence

- Some key desiderata:
  - Tailored to adaptation problem.
  - Captures structure: loss function, hypothesis set.
  - Can be estimated from finite samples.
  - Can be leveraged algorithmically.



# Discrepancy

# Discrepancy

## ■ Labeled discrepancy:

$$\text{dis}(\mathcal{P}, \mathcal{Q}) = \sup_{h \in \mathcal{H}} \left\{ \mathbb{E}_{(x,y) \sim \mathcal{P}} [\ell(h(x), y)] - \mathbb{E}_{(x,y) \sim \mathcal{Q}} [\ell(h(x), y)] \right\}.$$

## ■ Unlabeled discrepancy:

$$\overline{\text{dis}}(\mathcal{P}, \mathcal{Q}) = \sup_{h, h' \in \mathcal{H}} \left\{ \mathbb{E}_{x \sim \mathcal{P}_X} [\ell(h(x), h'(x))] - \mathbb{E}_{x \sim \mathcal{Q}_X} [\ell(h(x), h'(x))] \right\}.$$

- also finer local labeled or unlabeled discrepancy (Cortes et al., 2019).
- unlabeled discrepancy coincides with  $d_A$ -distance of (Kifer et al., 2004), for zero-one loss.

# Discrepancy - Properties

- Takes into account hypothesis set and loss function.
- Can be accurately estimated from finite samples for a hypothesis set with favorable complexity:

$$\left| \text{dis}(\mathcal{P}, \mathcal{Q}) - \text{dis}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}) \right| = O\left(\frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}}\right).$$

- Triangle inequality, distance under some assumptions.
- Upper bounds in terms of  $\ell_1$ -distance, relative entropy, Wassertein distance.
- Coincides with  $d_A$ -distance of (Kifer et al., 2004), for zero-one loss.

# Discrepancy Estimation

## ■ Notation:

- $\hat{\mathcal{P}}$  empirical distribution for sample drawn from  $\mathcal{P}^n$ .
- $\hat{\mathcal{Q}}$  empirical distribution for sample drawn from  $\mathcal{Q}^n$ .

## ■ Theorem: With high probability, the following holds:

$$\left| \text{dis}(\mathcal{P}, \mathcal{Q}) - \text{dis}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}) \right| \leq 2\mathfrak{R}_n(\ell \circ \mathcal{H}) + 2\mathfrak{R}_m(\ell \circ \mathcal{H}) + O\left(\frac{1}{\sqrt{n}} + \frac{1}{\sqrt{m}}\right).$$

## ■ Proof:

$$\left| \text{dis}(\mathcal{P}, \mathcal{Q}) - \text{dis}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}) \right| \leq \sup_{h \in \mathcal{H}} \left| \left[ \mathcal{L}(\mathcal{P}, h) - \mathcal{L}(\hat{\mathcal{P}}, h) \right] - \left[ \mathcal{L}(\mathcal{Q}, h) - \mathcal{L}(\hat{\mathcal{Q}}, h) \right] \right|.$$

# Discrepancy - Upper Bounds

- Upper bounded by  $\ell_1$ -distance and relative entropy:

$$\begin{aligned} \text{dis}(\mathcal{P}, \mathcal{Q}) &= \sup_{h \in \mathcal{H}} \left\{ \mathbb{E}_{(x,y) \sim \mathcal{P}} [\ell(h(x), y)] - \mathbb{E}_{(x,y) \sim \mathcal{Q}} [\ell(h(x), y)] \right\} \\ &\leq \sup_{h \in \mathcal{H}} \iint_{\mathcal{X} \times \mathcal{Y}} |p(x, y) - q(x, y)| |\ell(h(x), y)| dx dy \leq \ell_1(\mathcal{P}, \mathcal{Q}). \end{aligned}$$

- Upper bounded via importance weights  $w(x, y) = \frac{p(x, y)}{q(x, y)}$ :

$$\begin{aligned} \text{dis}(\mathcal{P}, \mathcal{Q}) &= \sup_{h \in \mathcal{H}} \iint_{\mathcal{X} \times \mathcal{Y}} [w(x, y) - 1] q(x, y) \ell(h(x), y) dx dy \\ &= \sup_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{Q}} [\Delta w(x, y) \ell(h(x), y)] \\ &\leq \sqrt{\mathbb{E}_{(x,y) \sim \mathcal{Q}} [\Delta^2 w(x, y)] \sup_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{Q}} [\ell^2(h(x), y)]}. \end{aligned}$$

# Discrepancy - Upper Bounds

- Upper-bound in terms of Wasserstein distance

$$\mathcal{W}(\mathcal{P}, \mathcal{Q}) = \sup_{\|f\|_{\text{Lip}} \leq 1} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_X} [f(x)] - \mathbb{E}_{\mathbf{x} \sim \mathcal{Q}_X} [f(x)] \right\}.$$

- For  $\ell$   $\mu_\ell$ -Lipschitz,  $\ell(h(x), h'(x)) \leq \mu_\ell |h(x) - h'(x)|$ , and for a hypothesis set  $\mathcal{H}$  of  $\mu_{\mathcal{H}}$ -Lipschitz functions,

$$|h(x') - h'(x')| - |h(x) - h'(x)| \leq 2\mu_{\mathcal{H}} |x' - x|,$$

$$\begin{aligned} \overline{\text{dis}}(\mathcal{P}, \mathcal{Q}) &= \sup_{h, h' \in \mathcal{H}} \left\{ \mathbb{E}_{x \in \mathcal{P}_X} [\ell(h(x), h'(x))] - \mathbb{E}_{x \in \mathcal{Q}_X} [\ell(h(x), h'(x))] \right\} \\ &\leq 2\mu_\ell \mu_{\mathcal{H}} \mathcal{W}(\mathcal{P}, \mathcal{Q}). \end{aligned}$$

# Discrepancy-Based Guarantee

■ Notation:

$$\mathcal{L}(\mathcal{P}, h) = \mathbb{E}_{(x,y) \sim \mathcal{P}} [\ell(h(x), y)] \quad \mathcal{L}(\mathcal{P}, h, h') = \mathbb{E}_{x \sim \mathcal{P}_x} [\ell(h(x), h'(x))].$$

- **Theorem:** Assume that  $\ell$  verifies the triangle inequality. Then, the following inequality holds for all  $h \in \mathcal{H}$ :

$$\mathcal{L}(\mathcal{P}, h) \leq \inf_{\substack{(h_{\mathcal{Q}}, h_{\mathcal{P}}) \in \mathcal{H}_{\text{all}} \times \mathcal{H} \\ \vee (h_{\mathcal{Q}}, h_{\mathcal{P}}) \in \mathcal{H} \times \mathcal{H}_{\text{all}}}} \left\{ \mathcal{L}(\mathcal{Q}, h, h_{\mathcal{Q}}) + \overline{\text{dis}}(\mathcal{P}, \mathcal{Q}) + \mathcal{L}(\mathcal{P}, h_{\mathcal{P}}) + \min \{ \mathcal{L}(\mathcal{Q}, h_{\mathcal{Q}}, h_{\mathcal{P}}), \mathcal{L}(\mathcal{P}, h_{\mathcal{Q}}, h_{\mathcal{P}}) \} \right\}.$$

# Discrepancy-Based Guarantee

## ■ Properties:

- always tighter than bound of (Ben-David et al., 2010):

$$\mathcal{L}(\mathcal{P}, h) \leq \mathcal{L}(\mathcal{Q}, h) + \overline{\text{dis}}(\mathcal{P}, \mathcal{Q}) + \min_{h' \in \mathcal{H}} \{\mathcal{L}(\mathcal{Q}, h') + \mathcal{L}(\mathcal{P}, h')\}.$$

- for same best-in class hypotheses  $h_{\mathcal{Q}}^* = h_{\mathcal{P}}^* = h^*$ , bound becomes:

$$\mathcal{L}(\mathcal{P}, h) \leq \mathcal{L}(\mathcal{Q}, h, h^*) + \overline{\text{dis}}(\mathcal{P}, \mathcal{Q}) + \mathcal{L}(\mathcal{P}, h^*).$$

- for consistent case, bound becomes:

$$\mathcal{L}(\mathcal{P}, h) \leq \mathcal{L}(\mathcal{Q}, h, f_{\mathcal{P}}) + \overline{\text{dis}}(\mathcal{P}, \mathcal{Q}).$$



# Discrepancy-Based Guarantee

- **Proof:** By definition of the triangle inequality and the discrepancy,

$$\mathcal{L}(\mathcal{P}, h) \leq \inf_{h_{\mathcal{P}} \in H} \left\{ \mathcal{L}(\mathcal{P}, h, h_{\mathcal{P}}) + \mathcal{L}(\mathcal{P}, h_{\mathcal{P}}) \right\} \quad (\text{triangle ineq.})$$

$$\leq \inf_{h_{\mathcal{P}} \in H} \left\{ \mathcal{L}(\mathcal{Q}, h, h_{\mathcal{P}}) + \overline{\text{dis}}(\mathcal{P}, \mathcal{Q}) + \mathcal{L}(\mathcal{P}, h_{\mathcal{P}}) \right\} \quad (\text{def. of discrepancy})$$

$$\leq \inf_{h_{\mathcal{Q}} \in \mathcal{H}_{\text{all}}, h_{\mathcal{P}} \in H} \left\{ \mathcal{L}(\mathcal{Q}, h, h_{\mathcal{Q}}) + \mathcal{L}(\mathcal{Q}, h_{\mathcal{Q}}, h_{\mathcal{P}}) + \overline{\text{dis}}(\mathcal{P}, \mathcal{Q}) + \mathcal{L}(\mathcal{P}, h_{\mathcal{P}}) \right\}. \quad (\text{triangle ineq.})$$

- Combining similar inequalities yields:

$$\mathcal{L}(\mathcal{P}, h) \leq$$

$$\inf_{\substack{(h_{\mathcal{Q}}, h_{\mathcal{P}}) \in \mathcal{H}_{\text{all}} \times \mathcal{H} \\ \forall (h_{\mathcal{Q}}, h_{\mathcal{P}}) \in \mathcal{H} \times \mathcal{H}_{\text{all}}}} \left\{ \mathcal{L}(\mathcal{Q}, h, h_{\mathcal{Q}}) + \overline{\text{dis}}(\mathcal{P}, \mathcal{Q}) + \mathcal{L}(\mathcal{P}, h_{\mathcal{P}}) + \min\{\mathcal{L}(\mathcal{Q}, h_{\mathcal{Q}}, h_{\mathcal{P}}), \mathcal{L}(\mathcal{P}, h_{\mathcal{Q}}, h_{\mathcal{P}})\} \right\}.$$

# Reweighting Algorithms

# Reweighting Algorithms

## ■ Ideas:

- Sample weights to reduce empirical discrepancy.
- Weights can affect weighted empirical loss.
- Select weights and predictor jointly.

## ■ General class of adaptation algorithm:

- KMM (Huang et al., 2006).
- KLIEP (Sugiyama et al., 2007).
- Importance weighting (analysis by (Cortes et al., 2010)).
- Discrepancy minimization (Cortes & MM, 2014).
- Generalized disc. minimization (Cortes et al., 2019).

# Learning Setup

## ■ General supervised adaptation scenario:

- Labeled sample  $S = ((x_1, y_m), \dots, (x_m, y_m)) \sim \mathcal{Q}^m$ .
- Labeled sample  $S' = ((x_{m+1}, y_{m+1}), \dots, (x_{m+n}, y_{m+n})) \sim \mathcal{P}^n$ .
- Non-negative weight vector  $q \in [0, 1]^{m+n}$ .
- Total weight of first  $m$  samples:  $\bar{q} = \sum_{i=1}^m q_i$ .

## ■ Problem:

- Find weights  $q \in [0, 1]^{m+n}$  and  $h \in \mathcal{H}$  to achieve small target domain expected loss  $\mathcal{L}(\mathcal{P}, h)$ .

# Weighted Rademacher Comp.

■ For  $\mathbf{q} \in [0, 1]^{[m+n]}$ ,

$$\mathfrak{R}_{\mathbf{q}}(\ell \circ \mathcal{H}) = \mathbb{E}_{S, S', \sigma} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^{m+n} \sigma_i \mathbf{q}_i \ell(h(x_i), y_i) \right].$$

- By Talagrand's contraction lemma,

$$\mathfrak{R}_{\mathbf{q}}(\ell \circ \mathcal{H}) \leq \|\mathbf{q}\|_{\infty} (m+n) \mathfrak{R}_{m+n}(\ell \circ \mathcal{H}).$$

# Reweighting Learning Bound

- **Theorem:** fix weights  $\mathbf{q} \in [0, 1]^{[m+n]}$ . Then, with probability at least  $1 - \delta$  over the draw of a sample  $S \sim \mathcal{Q}^m$  from the source domain and  $S' \sim \mathcal{P}^n$ , for any  $h \in \mathcal{H}$ ,

$$\mathcal{L}(\mathcal{P}, h) \leq \sum_{i=1}^{m+n} \mathbf{q}_i \ell(h(x_i), y_i) + \underbrace{\text{dis}\left(\left[(1 - \|\mathbf{q}\|_1) + \bar{\mathbf{q}}\right] \mathcal{P}, \bar{\mathbf{q}} \mathcal{Q}\right)}_{\substack{\mathbf{q} \text{ distribution} \\ \bar{\mathbf{q}} \text{dis}(\mathcal{P}, \mathcal{Q})}} + 2\mathfrak{R}_{\mathbf{q}}(\ell \circ \mathcal{H}) + \|\mathbf{q}\|_2 \sqrt{\frac{\log \frac{1}{\delta}}{2}}.$$

# Reweighting Lower Bound

- **Theorem:** fix distribution  $q \in \Delta_{m+n}$ . Then, for any  $\epsilon > 0$ , there exists  $h \in \mathcal{H}$  such that for any  $\delta > 0$ , with probability at least  $1 - \delta$  over the draw of a sample  $S \sim \mathcal{Q}^m$  from the source domain and  $S' \sim \mathcal{P}^n$ ,

$$\mathcal{L}(\mathcal{P}, h) \geq \sum_{i=1}^{m+n} q_i \ell(h(x_i), y_i) + \bar{q} \text{dis}(\mathcal{P}, \mathcal{Q}) + \Omega\left(\frac{1}{\sqrt{m+n}}\right).$$

- for  $\|q\|_2$ ,  $\mathfrak{R}_q(\ell \circ \mathcal{H}) \in O\left(\frac{1}{\sqrt{m+n}}\right)$ .

# Reweighting Uniform Bound

- **Theorem:** For any  $\delta > 0$ , with probability at least  $1 - \delta$  over the draw of a sample  $S \sim \mathcal{Q}^m$  and sample  $S' \sim \mathcal{P}^n$ , the following holds for all  $h \in \mathcal{H}$  and  $\mathbf{q} \in B_1(\mathbf{p}_0, 1)$ ,

$$\begin{aligned} \mathcal{L}(\mathcal{P}, h) \leq & \sum_{i=1}^{m+n} \mathbf{q}_i \ell(h(x_i), y_i) + \bar{\mathbf{q}} \text{dis}(\mathcal{P}, \mathcal{Q}) + \text{dis}(\mathbf{p}^0, \mathbf{q}) + 2\mathfrak{R}_{\mathbf{q}}(\ell \circ \mathcal{H}) \\ & + 8\|\mathbf{q} - \mathbf{p}^0\|_1 + [\|\mathbf{q}\|_2 + 2\|\mathbf{q} - \mathbf{p}^0\|_1] \left[ \sqrt{\log \log_2 \frac{2}{1 - \|\mathbf{q} - \mathbf{p}^0\|_1}} + \sqrt{\frac{\log \frac{2}{\delta}}{2}} \right]. \end{aligned}$$

- $\mathbf{p}_0$ : reference weights.



# Empirical Disc. Estimation

- Optimization problem:

$$\hat{d} = \max_{h \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=m+1}^{m+n} \ell(h(x_i), y_i) - \frac{1}{m} \sum_{i=1}^m \ell(h(x_i), y_i) \right\}.$$

- for a convex loss, can be cast as DC-programming problem and solved via DCA ([Tao and An, 1988](#)).
- for squared loss, global optimum convergence guarantee.

$$h_{t+1} \in \operatorname{argmin}_{h \in \mathcal{H}} \left\{ \frac{1}{m} \sum_{i=1}^m \ell(h(x_i), y_i) - \frac{1}{n} \sum_{i=m+1}^{m+n} \nabla \ell(h_t(x_i), y_i) \cdot (h - h_t) \right\}.$$

# Algorithm

- **Optimization problem:** SBEST algorithm.

$$\min_{h \in \mathcal{H}, \mathbf{q} \in [0,1]^{m+n}} \sum_{i=1}^{m+n} \mathbf{q}_i \ell(h(x_i), y_i) + \bar{\mathbf{q}} \text{dis}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}) + \lambda_\infty \|\mathbf{q}\|_\infty \|h\|^2 + \lambda_1 \|\mathbf{q} - \mathbf{p}^0\|_1 + \lambda_2 \|\mathbf{q}\|_2^2.$$

- Alternate minimization solution.
- For squared loss with linear predictors, convex optimization problem.
- Empirical discrepancy estimation via DC-programming.
- Extension to **weakly** or **unsupervised** adaptation.

# Labeled Disc. Upper Bounds

- **Theorem:** for squared loss, for any  $h_0 \in \mathcal{H}$ ,

$$\text{dis}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}) \leq \overline{\text{dis}}_{\mathcal{H} \times \{h_0\}}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}) + 2\delta_{\mathcal{H}, h_0}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}).$$

- where:

$$\delta_{\mathcal{H}, h_0}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}) = \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{(x, y) \sim \hat{\mathcal{P}}} [h(x)(y - h_0(x))] - \mathbb{E}_{(x, y) \sim \hat{\mathcal{Q}}} [h(x)(y - h_0(x))] \right|.$$

- favorable when  $h_0$  can be chosen so that  $|y - h_0(x)|$  is relatively small for both samples.
- note that  $\delta_{\mathcal{H}, h_0}(\hat{\mathcal{P}}, \hat{\mathcal{Q}}) = 0$  for  $\hat{\mathcal{P}} = \hat{\mathcal{Q}}$ .

# Experimental Results

# Classification Tasks

| Dataset      | Train source $\mathcal{Q}$ | Train target $\mathcal{P}$ | KMM              | gapBoost         | SBEST                              |
|--------------|----------------------------|----------------------------|------------------|------------------|------------------------------------|
| Adult        | $82.72 \pm 0.10$           | $81.61 \pm 0.42$           | $81.24 \pm 0.01$ | $83.1 \pm 0.02$  | $83.30 \pm 0.28$                   |
| German       | $68.24 \pm 0.21$           | $69.87 \pm 0.27$           | $65.7 \pm 0.01$  | $69.8 \pm 0.03$  | <b><math>71.26 \pm 0.11</math></b> |
| Accent       | $27.20 \pm 0.26$           | $81.64 \pm 0.22$           | $53.1 \pm 0.03$  | $81.2 \pm 0.04$  | <b><math>84.15 \pm 0.30</math></b> |
| comp vs sci  | $83.2 \pm 0.004$           | $89.4 \pm 0.03$            | $83.1 \pm 0.004$ | $92.08 \pm 0.01$ | <b><math>94.4 \pm 0.01</math></b>  |
| rec vs sci   | $79.2 \pm 0.007$           | $91.3 \pm 0.02$            | $79.7 \pm 0.004$ | $92.2 \pm 0.01$  | <b><math>92.4 \pm 0.004</math></b> |
| comp vs talk | $71.4 \pm 0.002$           | $89.9 \pm 0.02$            | $71 \pm 0.006$   | $90.6 \pm 0.01$  | <b><math>91 \pm 0.02</math></b>    |
| comp vs rec  | $65.4 \pm 0.007$           | $85.2 \pm 0.01$            | $67.7 \pm 0.007$ | $85.9 \pm 0.01$  | <b><math>88 \pm 0.01</math></b>    |
| rec vs talk  | $81.3 \pm 0.004$           | $88 \pm 0.02$              | $81.2 \pm 0.005$ | $89.2 \pm 0.01$  | <b><math>92.3 \pm 0.03</math></b>  |
| sci vs talk  | $88.2 \pm 0.005$           | $93.3 \pm 0.008$           | $88.5 \pm 0.003$ | $94.6 \pm 0.01$  | $94.6 \pm 0.02$                    |

- Details of experimental results in (Awasthi, Cortes, and MM, 2024).

# Fine-Tuning Tasks

| Fine-tuning           | Train on $\mathcal{P}$ | gapBoost       | SBEST                             |
|-----------------------|------------------------|----------------|-----------------------------------|
| Last layer (CIFAR-10) | 88.61 $\pm$ .43        | 87.1 $\pm$ .01 | <b>89.62 <math>\pm</math> .32</b> |
| Full model (CIFAR-10) | 90.18 $\pm$ .31        | 90.8 $\pm$ .02 | <b>92.30 <math>\pm</math> .24</b> |
| Last layer (Civil)    | 63.1 $\pm$ .12         | 64.7 $\pm$ .11 | <b>65.8 <math>\pm</math> .12</b>  |
| Full model (Civil)    | 65.8 $\pm$ .01         | 67.2 $\pm$ .01 | <b>68.3 <math>\pm</math> .14</b>  |

# Regression Tasks

| Dataset | KMM             | DM              | SBEST                              |
|---------|-----------------|-----------------|------------------------------------|
| Wind    | $1.2 \pm 0.04$  | $1.14 \pm 0.03$ | <b><math>0.97 \pm 0.02</math></b>  |
| Airline | $2.4 \pm 0.09$  | $1.72 \pm 0.1$  | <b><math>0.952 \pm 0.03</math></b> |
| Gas     | $0.41 \pm 0.01$ | $0.39 \pm 0.01$ | <b><math>0.38 \pm 0.02</math></b>  |
| News    | $1.08 \pm 0.01$ | $1.1 \pm 0.01$  | <b><math>0.99 \pm 0.01</math></b>  |
| Traffic | $2.1 \pm 0.1$   | $2.08 \pm 0.08$ | <b><math>0.99 \pm 0.002</math></b> |

# Sentiment Analysis

| $\mathcal{Q}$ | $\mathcal{P}$ | GDM             | DM              | KMM             | Train on $\mathcal{Q}$ |
|---------------|---------------|-----------------|-----------------|-----------------|------------------------|
| books         | <b>dvd</b>    | $1.25 \pm 0.01$ | $1.26 \pm 0.11$ | $1.43 \pm 0.08$ | $2.34 \pm 0.19$        |
|               | elec          | $0.88 \pm 0.01$ | $0.89 \pm 0.03$ | $1.50 \pm 0.05$ | $2.13 \pm 0.13$        |
|               | <b>ktchn</b>  | $1.06 \pm 0.03$ | $1.08 \pm 0.04$ | $1.47 \pm 0.01$ | $1.55 \pm 0.01$        |
| dvd           | <b>books</b>  | $1.14 \pm 0.02$ | $1.17 \pm 0.10$ | $1.64 \pm 0.14$ | $2.18 \pm 0.18$        |
|               | elec          | $1.08 \pm 0.01$ | $1.10 \pm 0.12$ | $2.40 \pm 0.05$ | $3.26 \pm 0.07$        |
|               | <b>ktchn</b>  | $1.1 \pm 0.03$  | $1.12 \pm 0.02$ | $1.10 \pm 0.02$ | $2.34 \pm 0.05$        |
| elec          | books         | $0.98 \pm 0.01$ | $1.00 \pm 0.01$ | $1.33 \pm 0.06$ | $1.34 \pm 0.04$        |
|               | <b>dvd</b>    | $0.98 \pm 0.02$ | $1.00 \pm 0.06$ | $1.00 \pm 0.06$ | $1.04 \pm 0.08$        |
|               | ktchn         | $0.96 \pm 0.01$ | $0.98 \pm 0.06$ | $1.04 \pm 0.01$ | $1.14 \pm 0.01$        |
| ktchn         | <b>books</b>  | $1.00 \pm 0.03$ | $1.04 \pm 0.07$ | $1.27 \pm 0.09$ | $1.12 \pm 0.08$        |
|               | <b>dvd</b>    | $1.2 \pm 0.002$ | $1.33 \pm 0.03$ | $1.32 \pm 0.03$ | $1.42 \pm 0.04$        |
|               | elec          | $1.64 \pm 0.02$ | $1.67 \pm 0.54$ | $1.87 \pm 0.56$ | $1.89 \pm 0.56$        |



# Conclusion

- Multiple-source adaptation problems.
- Discrepancy-based analysis of drifting (MM & Muñoz Medina, 2019; Awasthi, Cortes, and Mohri, 2022).
- Time series prediction and algorithms (MM & Kuznetsov, 2020).
- Differentially private adaptation from public to private domains or vice-versa (Bassily, Cortes, Mao, MM, 2024).
- Active learning (de Mathelin et al., 2022; Zhang et al., 2019, 2020).
- PAC-Bayesian analysis of adaptation (Germain et al., 2013).