



Tools for Uniform Convergence

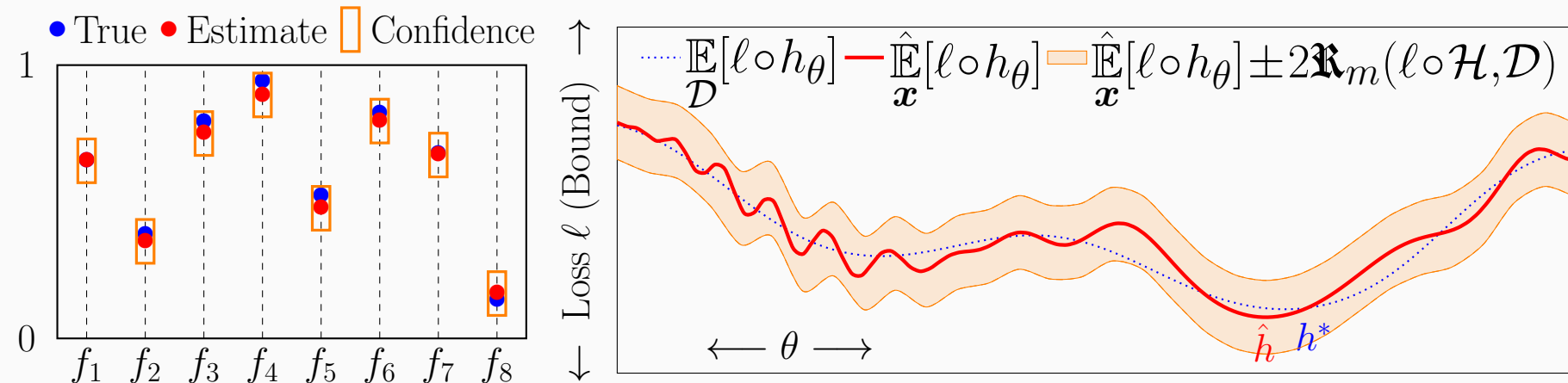
Family $\mathcal{F} \subseteq \mathcal{X} \rightarrow [0, r] \subseteq \mathbb{R}$ Sample $\mathbf{x} \sim \mathcal{D}^m$ Rademacher $\sigma \sim \mathcal{U}(\pm 1)$

Empirical Rade Avg $\hat{\mathfrak{R}}_m(\mathcal{F}, \mathbf{x}) \doteq \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i f(\mathbf{x}_i) \right| \right]$

Rademacher Avg $\mathfrak{R}_m(\mathcal{F}, \mathcal{D}) \doteq \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i f(\mathbf{x}_i) \right| \right] \right]$

The Wimpy Variance:

	Empirical	True
Raw	$\hat{v}_{\text{sup}}^{\text{raw}} \doteq \sup_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{x}} [f^2]$	$v_{\text{sup}}^{\text{raw}} \doteq \sup_{f \in \mathcal{F}} \mathbb{E} [f^2]$
Centralized	$\hat{v}_{\text{sup}} \doteq \sup_{f \in \mathcal{F}} \mathbb{V}_{\mathbf{x}} [f]$	$v_{\text{sup}} \doteq \sup_{f \in \mathcal{F}} \mathbb{V} [f]$



(Empirical) Centralization

The Symmetrization Inequality:

$$\frac{1}{2} \mathfrak{R}_m(\mathcal{F} - \mathbb{E}[\mathcal{F}], \mathcal{D}) \leq \mathbb{E}_{\mathbf{x}} \left[\underbrace{\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\mathcal{D}} [f] - \hat{\mathbb{E}}_{\mathbf{x}} [f] \right|}_{\text{SUPREMUM DEVIATION}} \right] \leq \underbrace{\frac{2 \mathfrak{R}_m(\mathcal{F} - \mathbb{E}[\mathcal{F}], \mathcal{D})}{2 \mathfrak{R}_m(\mathcal{F}, \mathcal{D})}}_{\text{CENTRALIZED / NONCENTRALIZED}}$$

The *Khintchine* and *Massart* inequalities tell us that

Noncentralized: $\sqrt{\frac{\hat{v}_{\text{sup}}^{\text{raw}}}{2m}} \leq \hat{\mathfrak{R}}_m(\mathcal{F}, \mathbf{x}) \leq \sqrt{\frac{2 \hat{v}_{\text{sup}}^{\text{raw}} \ln(2|\mathcal{F}|)}{m}}$

Centralized: $\sqrt{\frac{\hat{v}_{\text{sup}}}{2m}} \leq \hat{\mathfrak{R}}_m(\mathcal{F} - \mathbb{E}[\mathcal{F}], \mathbf{x}) \lesssim \sqrt{\frac{2 \hat{v}_{\text{sup}} \ln(2|\mathcal{F}|)}{m}}$

Problem: need to know $\mathbb{E}_{\mathcal{D}}[\mathcal{F}]$ to centralize!

Idea: Use ERA of empirically centralized family $\hat{\mathfrak{R}}_m(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x})$

Lemma: Empirical Centralization Bias

$$\frac{\mathbb{E}_{\mathbf{x}} [\hat{\mathfrak{R}}_m(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x})]}{1 + \Theta(1/\sqrt{m})} \leq \mathfrak{R}_m(\mathcal{F} - \mathbb{E}[\mathcal{F}], \mathcal{D}) \leq \frac{\mathbb{E}_{\mathbf{x}} [\hat{\mathfrak{R}}_m(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x})]}{1 - \Theta(1/\sqrt{m})}$$

Perceived Problems with the Rademacher Average

- Common criticisms of Rademacher averages
 - “Highly theoretical tool”
 - “Not useful in practice”
 - “Fail to explain” generalization of popular ML methods
- Where do these claims come from?
 - Is the symmetrization inequality tight?
 - Do we have sharp tail bounds on the empirical SD?
 - Are data-dependent bounds on $\hat{\mathfrak{R}}_m(\mathcal{F}, \mathbf{x})$ tight?
- In standard practice: **no, no, and no!**
 - To “fix” the bounds, we must repair *every part* of them
 - Can we match lower-bounds at every step?
 - * Minimax mean-estimation bound: $m \in \Omega\left(\frac{v_{\text{sup}} \ln \frac{1}{\delta}}{2\epsilon^2}\right)$
 - The loosest bound is the bottleneck!

Standard McDiarmid Bounds

$$\underbrace{\mathbb{E}_{\mathbf{x}, \sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i f(\mathbf{x}_i) \right| \right]}_{\text{RA}} \leq \underbrace{\mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i f(\mathbf{x}_i) \right| \right]}_{\text{EMPIRICAL RA}} + r \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}$$

$$\underbrace{\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\mathcal{D}} [f] - \hat{\mathbb{E}}_{\mathbf{x}} [f] \right|}_{\text{SUPREMUM DEVIATION}} \leq 2 \underbrace{\mathbb{E}_{\mathbf{x}, \sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i f(\mathbf{x}_i) \right| \right]}_{\text{RA}} + r \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}$$

Can we break the

McDiarmid bottleneck?



Target: Bousquet’s Inequality

We want dependence on v_{sup} , not $v_{\text{sup}}^{\text{raw}}$ or r^2

In particular, we want to match *Bousquet’s inequality*:

$$\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\mathcal{D}} [f] - \hat{\mathbb{E}}_{\mathbf{x}} [f] \right| \leq 2 \mathfrak{R}_m(\mathcal{F}, \mathcal{D}) + \frac{r \ln \frac{1}{\delta}}{3m} + \sqrt{\frac{2(v_{\text{sup}} + 4r \mathfrak{R}_m(\mathcal{F}, \mathcal{D})) \ln \frac{1}{\delta}}{m}}$$

We show that \hat{v}_{sup} is sufficient to *sharply bound* v_{sup}

With High Probability

Theorem: Concentration of Empirically Centralized RAs

With probability at least $1 - \delta$ over choice of \mathbf{x} :

$$\mathbb{E}_{\mathbf{x}} [\hat{\mathfrak{R}}_m(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x})] \leq \hat{\mathfrak{R}}_m(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x}) + \mathcal{O}\left(\frac{r \ln \frac{1}{\delta}}{m} + \sqrt{\frac{4r(\hat{\mathfrak{R}}_m(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x}) + r/\sqrt{m}) \ln \frac{1}{\delta}}{m}}\right)$$

The *Monte-Carlo Method*: Given $\sigma \in (\pm 1)^{n \times m}$, define

$$\hat{\mathfrak{R}}_m^n(\mathcal{F}, \mathbf{x}, \sigma) \doteq \frac{1}{n} \sum_{j=1}^n \sup_{f \in \mathcal{F}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_{j,i} f(\mathbf{x}_i) \right|$$

Theorem: Monte-Carlo Error Bounds

With probability at least $1 - \delta$ over choice of $\sigma \sim \mathcal{U}^{n \times m}(\pm 1)$:

$$\hat{\mathfrak{R}}_m(\mathcal{F}, \mathbf{x}) \leq \hat{\mathfrak{R}}_m^n(\mathcal{F}, \mathbf{x}, \sigma) + \frac{2r \ln \frac{1}{\delta}}{3nm} + \sqrt{\frac{4 \hat{v}_{\text{sup}}^{\text{raw}} \ln \frac{1}{\delta}}{nm}}$$

$$\hat{\mathfrak{R}}_m(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x}) \leq \hat{\mathfrak{R}}_m^n(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x}, \sigma) + \frac{4r \ln \frac{1}{\delta}}{3nm} + \sqrt{\frac{4 \hat{v}_{\text{sup}} \ln \frac{1}{\delta}}{nm}}$$

Comparative Analysis of Tail Bounds

- Prior Work: with probability at least $1 - \delta$, we may bound the SD as

$$2 \hat{\mathfrak{R}}_m(\mathcal{F}, \mathbf{x}) + \mathcal{O}\left(\frac{r \ln \frac{1}{\delta}}{m} + \sqrt{\frac{(v_{\text{sup}} + r \mathfrak{R}_m(\mathcal{F}, \mathcal{D})) \ln \frac{1}{\delta}}{m}}\right)$$
 - Implicit dependence of $\Omega\left(\sqrt{\frac{v_{\text{sup}}^{\text{raw}}}{m}}\right)$
 - Prior work: Monte-Carlo requires $n \in \Omega\left(\frac{r^2}{v_{\text{sup}}^{\text{raw}}}\right)$ trials
 - This work: $n = 1$ is *asymptotically optimal*
- This Work: With probability at least $1 - \delta$, we may bound the SD as

$$\frac{2 \hat{\mathfrak{R}}_m(\mathcal{F} - \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{F}], \mathbf{x})}{1 - \Theta\left(\sqrt{\frac{1}{m}}\right)} + \mathcal{O}\left(\frac{r \ln \frac{1}{\delta}}{m} + \sqrt{\frac{(v_{\text{sup}} + r \mathfrak{R}_m(\mathcal{F} - \mathbb{E}_{\mathcal{D}}[\mathcal{F}], \mathcal{D}) + \frac{r^2}{\sqrt{m}}) \ln \frac{1}{\delta}}{m}}\right)$$
 - All $v_{\text{sup}}^{\text{raw}}$ dependence moved to v_{sup} dependence
 - Can substitute \hat{v}_{sup} and Monte-Carlo ERAs
- Versus Localization
 - Complementary methods: *first* and *second* moment corrections
 - Localization often conflates *raw* and *centralized* variances