

Leave-one-out approach to analyze high-d ERM's

Setting: Let $\mathcal{Z} = \{x_i, y_i\}_{i \in [n]}$ n iid samples.

We train a linear model:

$$\hat{w} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i \in [n]} \ell(\langle w, x_i \rangle, y_i) + \frac{\lambda}{2} \|w\|^2$$

(Ex) for today: Binary classification, Gaussian separated case
 $x_i \sim \mathcal{N}(0, \operatorname{Id})$, $y_i = \operatorname{sgn}(\langle \beta, x_i \rangle)$
 $\|\beta\| = 1.$

The test error is $\varepsilon_g(\hat{w}) = \mathbb{P}_{x_{\text{test}}, y_{\text{test}}} \left[y_{\text{test}} \neq \operatorname{sgn}(\langle \hat{w}, x_{\text{test}} \rangle) \right]$

Claim (e.g. El Karoui, 2013)

In the limit $\begin{cases} n, d \rightarrow \infty \\ n/d = \alpha = \Theta_n(1) \end{cases}$, $\varepsilon_g(\hat{w}) \xrightarrow{P} E(\alpha, \lambda)$.

where E is some complicated function.

→ Gives a typical case characterization in high-d.

In the following, we highlight (heuristically, a proof being typically pretty lengthy) how to reach this claim.

Leave-one-out

Idea: $\varepsilon_g(\hat{w}) = \frac{1}{\pi} \arccos \left[\frac{\langle \hat{w}, \beta \rangle}{\|\hat{w}\|} \right]$

Remember that since \hat{w} minimizes the empirical risk,

$$\hat{w} = -\frac{1}{\lambda n} \sum_i \ell'_i(\langle \hat{w}, x_i \rangle) x_i \equiv r_i$$

Thus:

$$\|\hat{w}\|^2 = -\frac{1}{\lambda n} \sum_i \ell'_i(\langle \hat{w}, x_i \rangle) r_i$$

Ideally, we would like
① Concentration
② Distribution of r_i
→ what follows

We would like to understand the distribution of

$$r_i = \langle \hat{w}, x_i \rangle$$

complicated correlations with x_i !

Idea: define the "leave-one-out" empirical risk

$$\text{L.o.o.} \begin{cases} \hat{R}_{12}(w) = \frac{1}{n} \sum_{j \neq i} \ell_j(\langle w, x_j \rangle) + \frac{\lambda}{2} \|w\|^2 \\ \hat{w}_{12} = \underset{w}{\operatorname{argmin}} \hat{R}_{12}(w) \end{cases}$$

Then the estimator

$$\hat{w} = \underset{w}{\operatorname{argmin}} \left[\hat{R}_{12}(w) + \frac{1}{n} \ell_i(\langle w, x_i \rangle) \right] \quad \text{Taylor exp}$$

should intuitively be close to

$$\tilde{w}_2 = \underset{w}{\operatorname{argmin}} \left[\hat{R}_{12}(\hat{w}_{12}) + \frac{1}{2} \langle w - \hat{w}_{12}, H_{12}(w - \hat{w}_{12}) \rangle + \frac{1}{n} \ell_i(\langle w, x_i \rangle) \right]$$

where we introduced the Hessian

$$H_{ii} = \frac{1}{n} \sum_{j \neq i} \ell_j''(\langle \hat{w}_{12}, x_j \rangle) x_j x_j^T + \lambda \text{Id}$$

$$\text{Thus, } \hat{w} \approx \tilde{w}_2 = \hat{w}_{12} - \frac{1}{n} \ell_i'(\langle \hat{w}_{12}, x_i \rangle) H_{ii}^{-1} x_i \quad (*)$$

$$\left(\text{Lemma } \sup_{i \in [n]} \|\hat{w} - \tilde{w}_2\| = O_p \left(\frac{\text{poly}(\ln n)}{n} \right) \right)$$

Taking $\langle x_i, (*) \rangle$:

$$\hat{r}_i = \underbrace{\langle \hat{w}_{12}, x_i \rangle}_{\sim \mathcal{N}(0, \|w\|)} - \underbrace{\ell_i'(\langle \hat{w}_{12}, x_i \rangle) \frac{x_i^T H_{ii}^{-1} x_i}{n}}_{\approx \frac{1}{n} \operatorname{tr}(H_{ii}^{-1}) \equiv V}$$

Classroom: we removed the complicated stat. dependencies!

Thus

$$\tilde{r}_i = \text{prox}_{V \ell(\cdot, y_i)} (\|w\| g)$$

This result allows to go back and "close" the equation on $\|w\|$

$$\begin{aligned} \|w\|^2 &\approx -\frac{1}{\lambda n} \sum_i \ell'_i(\tilde{r}_i) r_i \approx -\frac{1}{\lambda n} \sum_i \ell'_i(\tilde{r}_i) \tilde{r}_i \\ &\approx -\frac{1}{\lambda} \mathbb{E} \left[\ell'(\text{prox}_{V \ell(\cdot, y)} (\|w\| g), \text{prox}_{V \ell(\cdot, y)} (\|w\| g)) \right] \end{aligned}$$

Iterated ERM's:

ERM I $w_0 = \underset{w \in \mathbb{R}^d}{\text{argmin}} \sum_{i \in [n]} \ell_0(\langle w, x_i \rangle, y_i) + \frac{\lambda}{2} \|w\|^2$

ERM II $w^i = \underset{w \in \mathbb{R}^d}{\text{argmin}} \sum_{i \in [n]} \ell(\langle w, x_i \rangle, \underbrace{\langle \hat{w}_0, x_i \rangle, y_i}_{\downarrow}) + \frac{\lambda}{2} \|w\|^2$

Depends on all the $\{x_j\}_{j \in [n]}$ through \hat{w}_0 (see ERM I).

Now, removing the i -th term from the sum no longer suffices to unravel all stat. dependences on x_i : \rightarrow nested L.O.O.

Thom (very informed) with Yue Lu, '26.

$$\left\{ \begin{aligned} \tilde{r}_{0,i} &= \text{prox}_{V_0 \ell_0(\cdot, y_i)} (\|w_0\| g_i) \rightarrow \text{same} \\ \tilde{r}_i &= \text{prox}_{V \ell(\cdot, \tilde{r}_{0,i}, y_i)} (\|w\| g_2 + \underbrace{\lambda \text{ du } \ell_0(\tilde{r}_{0,i})}_{\text{additional term}}) \end{aligned} \right.$$

Applications:

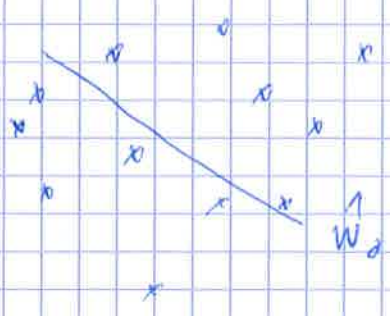
- Optimization: e.g. proximal descent.
- Sample reweighting, e.g. Ada Boost.

(E.x.)

Active learning

You have n unlabeled data pts $\{x_i\}_{i \in [n]}$, but only enough budget to label a fraction $\gamma \in (0, 1)$ of them. How to choose which points to label?

Idea (margin-based / uncertainty based selection)



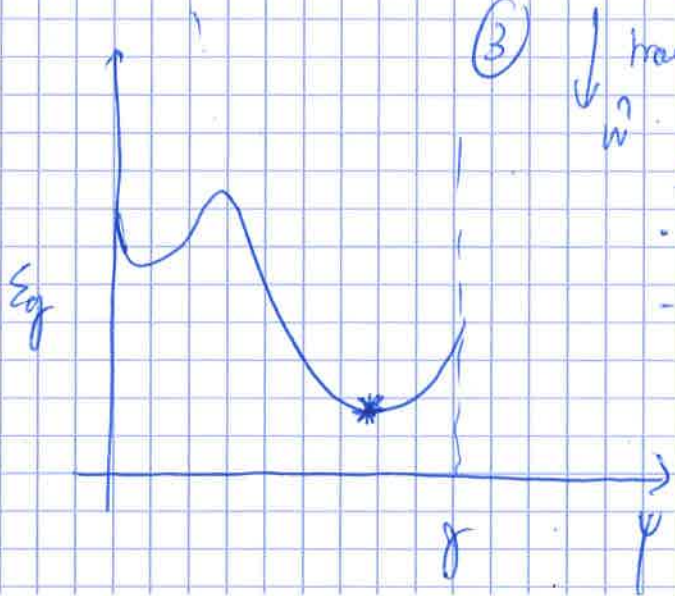
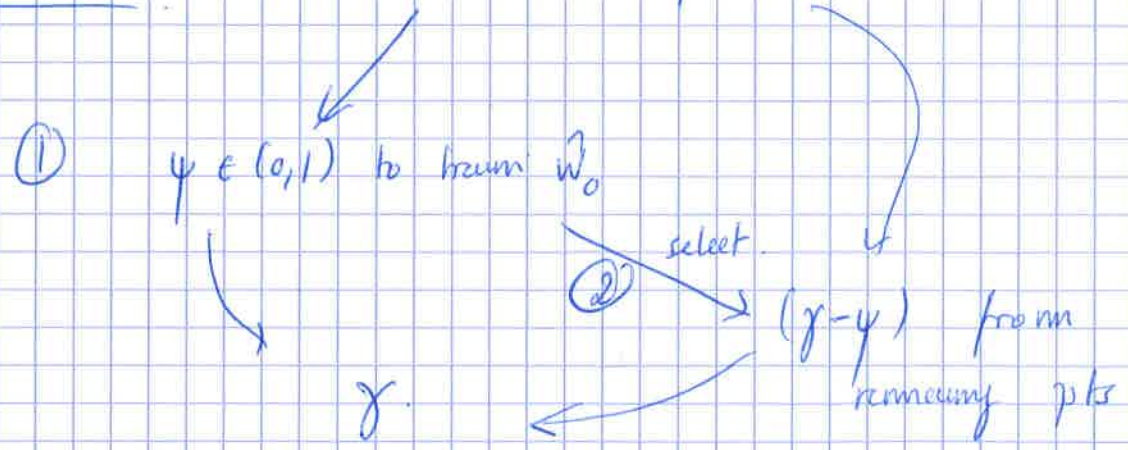
If a pre-estimate w_0 is available, label points closest to decision boundary (asymmetr. $|\langle w_0, x_i \rangle|$)

Many works in $n, d \rightarrow \infty, n/d = \Theta(1)$ regime [Koltassov et al, '23, Sorscher et al, '22]

assume $w_0 \perp$ dataset (oracle, or held-out data).

More realistic:

n samples



- some observations
- Tradeoff in choosing γ
 - Double descent.