

Benign overfitting

Last lecture we showed that uniform convergence allows for arguments like, $\forall \epsilon > 1-\delta, \forall f \in \mathcal{F}, |\hat{L}(f) - L(f)| \leq \sqrt{\frac{\text{Complexity}(\mathcal{F})}{n}}$

i.e. we can guarantee small population error if empirical error is small & # samples is sufficiently large.

However, we also saw in Zhang et al's paper that interpolators ($\hat{L}(f) = 0$) can achieve good performance even on noisy problems ($L(f) \geq c > 0$). In particular, in many modern deep learning settings, we have

$$c \leq L(f) = |L(f) - \hat{L}(f)| \leq 2 \cdot \min_{f \in \mathcal{F}} L(f).$$

Clearly, such settings cannot have a simple uniform convergence argument.

We'll now show a result which provably allows for

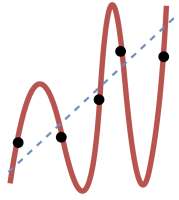
$$\min_{f \in \mathcal{F}} L(f) = c \leq L(f_n) = |L(f_n) - \hat{L}(f_n)| \leq \min_{f \in \mathcal{F}} L(f) + o_n(1).$$

This is called "benign overfitting":

- "overfitting" since $c \leq L(f) \leq \hat{L}(f) = 0 < c$.

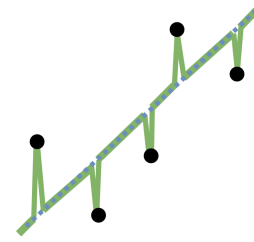
- "benign" since as $n \rightarrow \infty$, $L(f_n) \rightarrow \min_{f \in \mathcal{F}} L(f)$ while $\hat{L}(f_n) = 0 \forall n$.

"Classical" (catastrophic) form of overfitting:
think of using high degree polynomial to fit noisy linear model:



Catastrophic overfitting

• compare w/ another
• "overfitting" estimator:



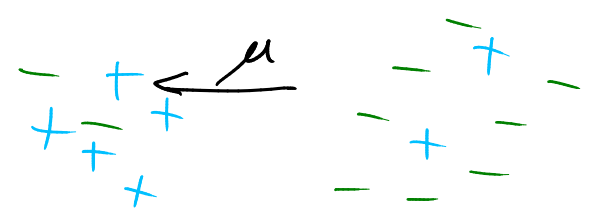
Benign overfitting

The setting: binary Gaussian mixture model. $\mu \in \mathbb{R}^d$, $p \in (0, \frac{1}{2})$

$\tilde{y} \sim \text{Unif}(\{\pm 1\})$, $x | \tilde{y} \sim \tilde{y}\mu + z$, $z \sim \mathcal{N}(0, \mathbb{I}_d)$;

$y = \begin{cases} \tilde{y}, & \text{w.p. } 1-p, \\ -\tilde{y}, & \text{w.p. } p. \end{cases}$

$(x, y) \sim \mathbb{P}$.



Suppose $(x_i, y_i) \stackrel{\text{iid}}{\sim} \mathbb{P}$.

Call $i \in \mathcal{N}$ if $y_i = -\tilde{y}_i$; $i \in \mathcal{C}$ if $y_i = \tilde{y}_i$. (Learner does not know if $i \in \mathcal{C}$ or $i \in \mathcal{N}$).

Then $\mathcal{N} \cup \mathcal{C} = [n]$.

We will show that $n\text{AVG} := \sum_{i=1}^n y_i x_i$ exhibits benign overfitting, under certain conditions.

Two things to show:

(1) $\hat{L}(n\text{AVG}) = 0$: for each $k \in [n]$, $\int_k \langle n\text{AVG}, x_k \rangle > 0$

(2) $p \leq L(n\text{AVG}) \leq p + o_n(1)$.

Lemma If training data is p -orthogonal in the sense that for some $p \geq 2$ and for $R^2 = \max_{i \neq j} \frac{\|x_i\|^2}{\|x_j\|^2} < \infty$, we have

$\|x_k\|^2 \geq pR^2 n \cdot \max_{i \neq j} |\langle x_i, x_j \rangle|$ for $k=1, \dots, n$, then for all $k \in [n]$, $y_k \langle x_k, \text{AVG} \rangle > 0$.

$$\begin{aligned} \text{Pf. } \langle y_k x_k, \text{AVG} \rangle &= \left\langle y_k x_k, \sum_{i=1}^n y_i x_i \right\rangle \\ &= \|x_k\|^2 + \sum_{i \neq k} \langle y_k x_k, y_i x_i \rangle \\ &\geq \|x_k\|^2 - n \cdot \max_{i \neq j} |\langle x_i, x_k \rangle| \\ &\geq \frac{1}{2} \|x_k\|^2 > 0, \quad \text{since } R^2 \geq 1 \text{ \& } p \geq 2. \quad \square \end{aligned}$$

Thus, AVG interpolates the training data: $(\hat{L}(f) = 0)$ if the training data is p -orthogonal. We now establish sufficient conditions for this.

Note: $\|x_k\|^2 = \|\tilde{y}_k + z_k\|^2 = \|\tilde{y}_k\|^2 + \|z_k\|^2 + 2\langle \tilde{y}_k, z_k \rangle$, so suffices to control both $\|z_k\|^2$, $\langle z_k, u \rangle$ for fixed vector u .

Lemma There is a $C_0 > 1$ s.t. for $\delta \in (0, \frac{1}{2})$, if $d > C_0^3 \log(12n/\delta)$, then w.p. $> 1-\delta$, we have:

$$\textcircled{1} \forall k, \left| \frac{\|z_k\|}{\sqrt{d}} - 1 \right| \leq C_0 \sqrt{\frac{\log(12n/\delta)}{d}} \Leftrightarrow \left| \|z_k\| - \sqrt{d} \right| \leq C_0 \sqrt{\log\left(\frac{12n}{\delta}\right)}$$

$$\textcircled{2} \forall i \neq j, |\langle z_i, z_j \rangle| \leq C_0 \sqrt{d} \log(12n^2/\delta)$$

Pf Part $\textcircled{1}$ & $\textcircled{2}$ were HW.

Lemma There is $C_1 > 1$ s.t. for $\delta \in (0, \frac{1}{2})$, if $d > C_0^3 \log(12n/\delta)$, and if $d \geq C_0 (n \|\mu\|^2 \vee n^2 \log(12n/\delta))$, $\|\mu\| \geq 3C_0$, then for some absolute $C_1 > 1$,

$$\text{w.p. } > 1-\delta, \quad \|X_{k2}\|^2 \geq C_1 n \cdot \frac{\max_{i,j} \|x_i\|^2}{\|x_j\|^2} \cdot \max_{i \neq j} |\langle x_i, x_j \rangle|.$$

$$\begin{aligned} \text{Pf } |\langle x_i, x_j \rangle| &= |\langle \tilde{y}_i \mu + z_i, \tilde{y}_j \mu + z_j \rangle| \\ &\leq \|\mu\|^2 + |\langle \mu, z_i \rangle| + |\langle \mu, z_j \rangle| + |\langle z_i, z_j \rangle| \\ &\leq \|\mu\|^2 + 2C_0 \|\mu\| \log(12n/\delta) + C_0 \sqrt{d} \log(12n^2/\delta). \\ &\leq \|\mu\|^2 + 2C_0 (\|\mu\| \sqrt{d}) \log(12n^2/\delta). \end{aligned}$$

$$\|X_{k2}\|^2 = \|\tilde{y}_k \mu + z_k\|^2 = \|\mu\|^2 + 2\tilde{y}_k \langle \mu, z_k \rangle + \|z_k\|^2. \quad \|\mu\|^2 \geq \frac{d}{C_0 n}, \quad \|\mu\| \leq \sqrt{\frac{d}{C_0 n}}.$$

By prev lemma, $-C_0 \sqrt{\log(12n/\delta)} \leq \|z_k\| - \sqrt{d} \leq C_0 \sqrt{\log(12n/\delta)}$. For $d > C_0^3 \log(12n/\delta)$,

$$\sqrt{d} = \frac{\sqrt{d}}{2} + \frac{\sqrt{d}}{2} = \frac{C_0^{3/2}}{2} \sqrt{\log(12n/\delta)} + \frac{\sqrt{d}}{2} \Rightarrow \|z_k\| \geq \frac{\sqrt{d}}{2} + \left(\frac{C_0^{3/2}}{2} - C_0\right) \sqrt{\log(12n/\delta)} \geq \frac{\sqrt{d}}{2} \text{ for } C_0 \text{ large enough.}$$

$$\rightarrow \|x_{\text{all}}\|^2 \leq \|\mu\|^2 + 2C_0 \|\mu\| \log(12n/\delta) + d \left(1 + C_0^2 \cdot \log(12n/\delta)\right)$$

$$\|x_{\text{all}}\|^2 \geq \|\mu\|^2 - 2C_0 \|\mu\| \log(12n/\delta) + d/4.$$

For $d \geq C_0 \|\mu\|^2$, $\|\mu\| \geq 3C_0$, $\& \ d \geq C_0^3 \log(12n/\delta)$, we get:

$$\begin{aligned} \|x_{\text{all}}\|^2 &\leq d \left(\frac{1}{C_0} + \frac{2}{C_0} \frac{\log(12n/\delta)}{n} + 1 + \frac{1}{C_0} \right), & \left\{ \begin{array}{l} R^2 \leq \frac{(1 + \frac{1}{C_0})}{(1 - \frac{1}{C_0})} \leq (1 + 2C_0)^2 \leq 1.01 \\ \text{for } C_0 \text{ large.} \end{array} \right. \\ \|x_{\text{all}}\|^2 &\geq \frac{d}{4} \left(1 - \frac{2C_0 \|\mu\| \log(12n/\delta)}{d} \right) \geq \frac{d}{5} \text{ for } C_0 \text{ large.} \end{aligned}$$

$$\Rightarrow \frac{\|x_{\text{all}}\|^2}{R^2 \max_{i,j} | \langle x_i, x_j \rangle |} \geq \frac{d/5}{1.01 \cdot (\|\mu\|^2 + 2C_0 (\|\mu\| \sqrt{d}) \log(12n^2/\delta))} = \frac{d/5.05}{\|\mu\|^2 + 2C_0 \sqrt{d} \log(12n^2/\delta)}.$$

In order for near-orthogonality, this must be $\Omega(n)$:

$$\frac{d}{\|\mu\|^2 + \sqrt{d} \log(12n^2/\delta)} = \Omega(n) \text{ holds if } \frac{d}{\|\mu\|^2} = \Omega(n) \quad \& \quad \frac{\sqrt{d}}{\log(12n^2/\delta)} = \Omega(n)$$

These are precisely the assumptions in the lemma. \square

Now we move to generalization. We want to show $p \leq \mathbb{P}(y \neq \text{sgn}(\langle w, x \rangle)) \leq p + o(d)$.

Lemma Suppose $w \in \mathbb{R}^d$. Then

$$\mathbb{P}(y \neq \text{sgn}(\langle w, x \rangle)) \leq p + \Phi(-\langle w/\|w\|, \mu \rangle), \text{ where } \Phi \text{ is normal CDF.}$$

Remark: If $\langle w, \mu \rangle \leq 0$, vacuous; otherwise, decays exp-fast in $\langle \frac{w}{\|w\|}, \mu \rangle$.

$$\begin{aligned}
 \text{Pf } \mathbb{P}(y \neq \text{sgn}(\langle w, x \rangle)) &= \mathbb{P}(y \langle w, x \rangle < 0) \\
 &= \mathbb{P}(y \langle w, x \rangle < 0, y = \tilde{y}) + \mathbb{P}(y \langle w, x \rangle < 0, y = -\tilde{y}) \\
 &\leq p + \mathbb{P}(y \langle w, x \rangle < 0, y = \tilde{y}).
 \end{aligned}$$

$$\begin{aligned}
 \mathbb{P}(y \langle w, x \rangle < 0, y = \tilde{y}) &= \mathbb{P}(\langle w, \tilde{y} x \rangle < 0) = \mathbb{P}(\langle w, \mu + \tilde{y} z \rangle < 0) \\
 &= \mathbb{P}(\langle w, \tilde{y} z \rangle < -\langle w, \mu \rangle) \\
 &= \mathbb{P}(N(0, 1) < -\langle \frac{w}{\|w\|}, \mu \rangle) \\
 &= \Phi(-\langle \frac{w}{\|w\|}, \mu \rangle). \quad \square
 \end{aligned}$$

So it suffices to show $\langle n \text{AVG}, \mu \rangle$ is large

Recall: $i \in R$ means $y_i = -\tilde{y}_i$, so $(x_i, y_i) = (\tilde{y}_i \mu + z_i, -\tilde{y}_i)$; $i \in C$: $(x_i, y_i) = (\tilde{y}_i \mu + z_i, \tilde{y}_i)$.

$$\langle n \text{AVG}, \mu \rangle = \langle \sum_{i=1}^n y_i x_i, \mu \rangle$$

$$= \langle \sum_{i \in C} \tilde{y}_i (\tilde{y}_i \mu + z_i), \mu \rangle + \langle \sum_{i \in R} -\tilde{y}_i (\tilde{y}_i \mu + z_i), \mu \rangle$$

$$= (|C| - |R|) \cdot \| \mu \|^2 + \sum_{i \in C} \langle \tilde{y}_i z_i, \mu \rangle - \sum_{i \in R} \langle \tilde{y}_i z_i, \mu \rangle.$$

$$= (n - 2|R|) \cdot \| \mu \|^2 + \langle \sum_{i=1}^n y'_i z_i, \mu \rangle, \quad \text{where}$$

$y_i' = \begin{cases} \tilde{y}_i: c \in \mathcal{C}, \\ -\tilde{y}_i: c \in \mathcal{V}. \end{cases}$
 Note that $y_i' \stackrel{\text{iid}}{\sim} \text{Unif}(\{\pm 1\}^d)$, since $\{\tilde{y}_i \in \mathcal{C}\}$ is indep. of \tilde{y}_i .

Thus it suffices to prove two things:

(1) an upper bound on $|M|$;

(2) an upper bound on $|\langle \sum_{i=1}^n y_i' z_i, \mu \rangle|$.

In homework, you will need to derive bounds on each of these.

But intuitively, $|M| \approx pn \pm O(\sqrt{n})$;

$\sum_{i=1}^n y_i' z_i \sim N(0, nI_d)$ by independence, so

$\langle \sum_{i=1}^n y_i' z_i, \mu \rangle \sim N(0, n\|\mu\|^2)$, so $|\langle \sum_{i=1}^n y_i' z_i, \mu \rangle| \lesssim \sqrt{n} \|\mu\|$;

thus

$$\langle n\text{AVG}, \mu \rangle \gtrsim (1-2p)n\|\mu\|^2 - \sqrt{n}\|\mu\|$$

$$\gtrsim n\|\mu\|^2 \text{ if } 1-2p < \text{const}, \|\mu\| \geq \text{const}.$$

Then need to bound $\|n\text{AVG}\|^2 = \|\sum_{i=1}^n y_i x_i\|^2$;

key here is to use near-orthogonality.

If $\frac{n\|\mu\|^4}{d} = \omega_d(1)$ then should get $\mathbb{P}(y \neq \text{sgn}(\langle \text{AVG}, x \rangle)) \leq p + o_d(1)$.