

Implicit regularization

We saw last class a theorem (by Lyou-Li '19, Ji-Telgarsky '20) which shows implicit bias of G.D. towards ℓ^2 -margin maximization.

We will now prove a number of implicit bias results.

We'll first consider regression.

Let $X \in \mathbb{R}^{n \times d}$, $y \in \mathbb{R}^n$, assume $n < d$ (over-parameterized / high-dimensional)
 \hat{X} full rank

consider $\hat{L}(\beta) := \frac{1}{2} \|X\beta - y\|_2^2 = \frac{1}{2} \sum_{i=1}^n (\langle x_i, \beta \rangle - y_i)^2$.

Lemma Let X^+ denote pseudo-inverse of X . Then β is a global min of $\hat{L} \iff \beta = X^+ y + \xi$ for some $\xi \perp \text{span}\{x_1, \dots, x_n\}$.

Pf. Any $\beta \in \mathbb{R}^d$ can be represented as

$$\beta = X^+ y + \xi \text{ for some } \xi \in \mathbb{R}^d \quad (\xi = \beta - X^+ y)$$

$$\begin{aligned} \text{Since } X\beta &= X(X^+ y + \xi) = XX^+ y + X\xi \\ &= X \cdot X^T (XX^T)^{-1} y + X\xi \\ &= y + X\xi. \end{aligned}$$

Since X is full rank, $X^+ = X^T (XX^T)^{-1}$

Now, β is a global min of $\hat{L} \Leftrightarrow \|X\beta - y\|_2^2 = 0$.

$$\Leftrightarrow \|y + X\xi - y\|_2^2 = \|X\xi\|_2^2 = 0$$

$$\|X\xi\|_2^2 = 0 \Leftrightarrow \sum_{i=1}^n \langle x_i, \xi \rangle^2 = 0 \Leftrightarrow \langle x_i, \xi \rangle = 0 \quad \forall i=1, \dots, n.$$
$$\Leftrightarrow \xi \perp \text{span}\{x_1, \dots, x_n\}.$$

□

So every global min of squared loss lies in subspace spanned by data, and in that subspace is given by $X^+ y$.

Lemma Let $\beta^* := \arg \min \{ \|\beta\|_2^2 : \hat{L}(\beta) = 0 \}$.

Then $\beta^* = X^+ y$.

pf. Let β be s.t. $\hat{L}(\beta) = 0$. By prev lemma, $\exists \xi \perp \{x_1, \dots, x_n\}$ s.t.
 $\beta = X^+ y + \xi$. Thus, $\|X\xi\|_2 = 0$, and we have

$$\begin{aligned}
\|\beta\|_2^2 &= \|X^+y\|^2 + \|\xi\|^2 + 2\langle X^+y, \xi \rangle \\
&= \|X^+y\|^2 + \|\xi\|^2 + 2\langle X^T (XX^T)^{-1} y, \xi \rangle \\
&= \|X^+y\|^2 + \|\xi\|^2 + 2\langle (XX^T)^{-1} y, X\xi \rangle \\
&= \|X^+y\|^2 + \|\xi\|^2 \geq \|X^+y\|^2.
\end{aligned}$$

We get equality $\Leftrightarrow \xi = 0$. □

Now consider iterates of G.D. started from β_0 :

$\beta_{t+1} = \beta_t - \eta \nabla \hat{L}(\beta_t)$. \hat{L} is convex and 1-smooth, so we know $\beta_t \rightarrow$ global min of \hat{L} . We now show it converges to min. ℓ_2 norm solution, if $\beta_0 = 0$.

Theorem Let $\beta_0 = 0$, $\{\beta_t\}$ G.D. iterates. Suppose $\beta_t \rightarrow \hat{\beta}$ with $\hat{L}(\hat{\beta}) = 0$. Then $\hat{\beta} = \beta^* = \arg \min \{\|p\|_2 : \hat{L}(p) = 0\}$.

PF First, note that $\beta_t \in \text{span}\{x_1, \dots, x_n\} \forall t$: clear @ $t=0$.

$$-\nabla \hat{L}(\beta) = \frac{1}{n} \sum_i (y_i - \langle w, x_i \rangle) \cdot x_i \in \text{span}\{x_1, \dots, x_n\},$$

So an induction arg clearly shows $\beta_t \in \text{span}\{x_1, \dots, x_n\} \forall t$.

In particular, $\beta_t \rightarrow \hat{\beta} \in \text{span}\{x_1, \dots, x_n\}$.

So $\hat{\beta} = X^T \nu$ for some $\nu \in \mathbb{R}^n$.

Since by assumption $\hat{L}(\hat{\beta}) = 0$, $0 = \|X\hat{\beta} - y\|_2^2$

implies $0 = X\hat{\beta} - y = XX^T \nu - y$, so

$$\nu = (XX^T)^{-1} y \Rightarrow \hat{\beta} = X^T \nu = X^T (XX^T)^{-1} y = X^+ y = \beta^*$$

So, G.P. on squared loss has implicit bias towards minimum l^2 -norm solution (an implicit regularization effect).

We'll now look at classification setting. Ref: [T], ch. 10.

Def Data $\{x_i, y_i\}$, $x_i \in \mathbb{R}^d$, $y_i \in \{\pm 1\}$, is linearly separable if, $\exists w \in \mathbb{R}^d$ st $\min_i y_i \langle w, x_i \rangle > 0$. For linearly sep. data,

The l_2 -max margin predictor is $\bar{w} := \arg \min \{ \|w\|_2^2 : y_i \langle w, x_i \rangle \geq 1 \forall i \}$.

Equivalently, $\bar{u} := \arg \max_{\{ \min_i y_i \langle w, x_i \rangle : \|w\|_2 = 1 \}}$.

Exercise Show that if the l_2 -max-margin predictor exists, it is unique.

Prop Suppose $f(x; \theta)$ is L -homog. in θ , l is exp-loss, and

$$\exists \hat{\theta} \text{ s.t. } \hat{L}(\hat{\theta}) = \sum_1^n l(y_i f(x_i; \hat{\theta})) < \frac{l(0)}{n}$$

Then $\inf_{\theta} \hat{L}(\theta) = 0$, and the inf is not attained.

Pf. Let $m_i(\theta) := y_i f(x_i; \theta)$ (margin of ex. i).

Since l is decreasing

$$y = \exp(-x) \\ \log y = -x \rightarrow x = -\log y$$

$$l(\min_i m_i(\theta)) = \max_i l(m_i(\theta)) \leq \sum_{i=1}^n l(m_i(\theta)) = \hat{L}(\theta) < l(0)/n \leq l(0)$$

Since $l^{-1}(t) = -\log t$ is ^{strictly} decreasing, $l^{-1}(l(\min_i m_i(\theta))) = \min_i m_i(\theta) > l^{-1}(l(0)) = 0$.

Thus: $0 \leq \inf_{\theta} \hat{L}(\theta) \leq \limsup_{c \rightarrow \infty} \hat{L}(c\hat{\theta}) = \limsup_{c \rightarrow \infty} \frac{1}{c} \sum_1^n l(c \cdot m_i(\hat{\theta}))$

$\hat{L}(\hat{\theta}) = \sum_1^n l(\cdot)$
not $\frac{1}{n}$.



Since $m_i(\theta) > 0$, l decreasing, $c > 0$,

$$l(c \cdot m_i(\theta)) \rightarrow 0 \text{ as } c \rightarrow \infty.$$

$$0 \leq \inf_{\theta} \hat{L}(\theta) \leq \limsup_{c \rightarrow \infty} \frac{1}{n} \sum_1^n l(c \cdot m_i(\theta)) \leq \frac{1}{n} \sum_1^n \limsup_{c \rightarrow \infty} l(c \cdot m_i(\theta)) = 0.$$

Thus $\inf_{\theta} \hat{L}(\theta) = 0$. Since $l(t) > 0 \forall t$, impossible to have $\hat{L}(\theta) = 0$.

□

So we cannot "find" an "optimum": solutions are off at ∞ .

TO compare predictors; first note

$$\min_i m_i(\theta) = \|\theta\|_2^L \cdot \min_i m_i\left(\frac{\theta}{\|\theta\|_2}\right) \quad \text{by homogeneity.}$$

→ we can compare by normalized margin $\theta / \|\theta\|_2$.

Moreover, for exp loss, we have: $l^{-1}(b) = -\log t$ (decreasing)

$$\frac{l^{-1}(\hat{L}(\theta))}{\|\theta\|_2^L} = \frac{l^{-1}\left(\frac{1}{n} \sum_1^n l(m_i(\theta))\right)}{\|\theta\|_2^L} \leq \frac{l^{-1}\left(\max_i l(m_i(\theta))\right)}{\|\theta\|_2^L} = \frac{\min_i m_i(\theta)}{\|\theta\|_2^L}$$

$$\frac{\ell^{-1}(\hat{L}(\theta))}{\|\theta\|_2^L} + \frac{\log(n)}{\|\theta\|_2^L} = \frac{\ell^{-1}\left(\frac{1}{n} \sum_i \ell(m_i(\theta))\right)}{\|\theta\|_2^L}$$

$$\geq \frac{\ell^{-1}(\max_i \ell(m_i(\theta)))}{\|\theta\|_2^L} \quad \text{since } \ell^{-1} \text{ is decreasing}$$

i.e. avg \leq max

$$= \frac{\min_i m_i(\theta)}{\|\theta\|_2^L} \quad \text{since } \ell \text{ is decreasing}$$

$$\geq \frac{\ell^{-1}\left(\sum_{n=1}^n \ell(m_i(\theta))\right)}{\|\theta\|_2^L}$$

since ℓ^{-1} is decr.
i.e. $\max_i \ell(m_i(\theta)) \leq \sum_i \ell(m_i(\theta))$

$$\Rightarrow \frac{\min_i m_i(\theta)}{\|\theta\|_2^L} \in \left[\frac{\ell^{-1}(\hat{L}(\theta))}{\|\theta\|_2^L}, \frac{\ell^{-1}(\hat{L}(\theta))}{\|\theta\|_2^L} + \frac{\log n}{\|\theta\|_2^L} \right].$$

We thus can approximate the normalized margin $\frac{\min_i m_i(\theta)}{\|\theta\|_2^L}$ by $\frac{\ell^{-1}(\hat{L}(\theta))}{\|\theta\|_2^L}$.

Def. Call data m -separable if $\exists \theta$ st $\min_i m_i(\theta) > 0$.

$$\gamma(\theta) := \min_i m_i(\theta) / \|\theta\|_2; \quad \bar{\gamma} := \max_{\|\theta\|_2=1} \gamma(\theta); \quad \tilde{\gamma}(\theta) := \frac{\ell^{-1}(\hat{L}(\theta))}{\|\theta\|_2^L}$$

(normalized) margin max margin smooth (normalized) margin.

Prop Suppose data is m -separable. Then:

(1) $\bar{\gamma} := \max_{\|\theta\|_2=1} \gamma(\theta) > 0$ is well-defined

(2) For any $\theta \neq 0$, $\lim_{c \rightarrow \infty} \tilde{\gamma}(c\theta) = \gamma(\theta)$.

Pf (1) We want to show that the maximum is attained.

Note that by assumption, $\exists \hat{\theta}$ st $\gamma(\hat{\theta}) > 0$. Since $m_i(\theta)$ is homogeneous, we know that $m_i(\alpha\theta) = \alpha^L m_i(\theta)$ for $\alpha > 0$. Thus $m_i(0) = 0$, hence $\hat{\theta} \neq 0$, so consider $\theta := \hat{\theta} / \|\hat{\theta}\|_2$. Then $\|\theta\|_2 = 1$. Thus,

$$\gamma(\theta) = \min_i m_i(\theta / \|\theta\|_2) = \|\theta\|_2^{-L} \min_i m_i(\theta) > 0.$$

Since $m_i(\theta)$ is continuous, i -wise min of cts fns is cts, $\gamma(\theta)$ is continuous in θ , and is strictly positive on at least one point on the domain $\{\|\theta\|_2 = 1\} = \mathcal{D}$

Thus its maximum must be strictly positive, and attained on \mathcal{D} due to compactness.

(2) Recall that $\tilde{\gamma}(\theta) := \frac{\ell^{-1}(\hat{\ell}(\theta))}{\|\theta\|_2^L} \leq \gamma(\theta) = \frac{\min_i m_i(\theta)}{\|\theta\|_2^L} \leq \frac{\ell^{-1}(\hat{\ell}(\theta))}{\|\theta\|_2^L} + \frac{\log n}{\|\theta\|_2^L} = \tilde{\gamma}(\theta) + \frac{\log n}{\|\theta\|_2^L}$.

$$\Rightarrow \tilde{\gamma}(c\theta) \leq \gamma(c\theta) = \gamma(\theta) \leq \tilde{\gamma}(c\theta) + \frac{\log n}{c^L \|\theta\|_2^L}$$

$$\Rightarrow \limsup_{c \rightarrow \infty} \tilde{\gamma}(c\theta) \leq \gamma(\theta), \text{ and } \liminf_{c \rightarrow \infty} \left\{ \tilde{\gamma}(c\theta) + \frac{\log n}{c^L \|\theta\|_2^L} \right\} = \liminf_{c \rightarrow \infty} \tilde{\gamma}(c\theta) \geq \gamma(\theta).$$

□

Gradient flow maximizes the margin of linear predictors.

Let $\hat{L}(\theta) = \sum_{i=1}^n \ell(y_i, \langle \theta, x_i \rangle)$. Gradient flow:

$$\frac{d\theta}{dt} = -\nabla \hat{L}(\theta(t)), \quad \text{with (assume) } \theta(0) = 0.$$

First, note that G.F. is always decreasing (even nonconvex):

$$\begin{aligned} \hat{L}(\theta(t)) - \hat{L}(\theta(0)) &= \int_0^t \langle \nabla \hat{L}(\theta(s)), \frac{d}{ds} \theta(s) \rangle ds \\ &= -\int_0^t \|\nabla \hat{L}(\theta(s))\|^2 ds \leq 0. \end{aligned}$$

Thus, $\min_{s \in [0, t]} \hat{L}(\theta(s)) = \hat{L}(\theta(t))$ for $t > 0$.

Theorem For any $z \in \mathbb{R}^d$, G.F. satisfies

$$t \hat{L}(\theta(t)) + \frac{1}{2} \|\theta(t) - z\|_2^2 \leq t \hat{L}(z) + \frac{1}{2} \|\theta(0) - z\|_2^2.$$

Pf. for any z ,

$$\frac{1}{2} \|\theta(t) - z\|_2^2 - \frac{1}{2} \|\theta(0) - z\|_2^2 = \frac{1}{2} \int_0^t \frac{d}{ds} \|\theta(s) - z\|_2^2 ds$$

$$= \int_0^t \left\langle \frac{d\theta}{ds}, \theta(s) - z \right\rangle ds$$

$$= \int_0^t \langle \nabla \hat{L}(\theta(s)), \theta(s) - z \rangle ds$$

$$\begin{aligned}
&= \int_0^t (\hat{L}(z) - \hat{L}(\theta(s))) ds \\
&= t \hat{L}(z) - \int_0^t \hat{L}(\theta(s)) ds
\end{aligned}$$

Now, $\int_0^t \hat{L}(\theta(s)) ds \geq \int_0^t \min_{s \in [0, t]} \hat{L}(\theta(s)) ds$

$$\begin{aligned}
&= t \cdot \min_{s \in [0, t]} \hat{L}(\theta(s)) \\
&= t \cdot \hat{L}(\theta(t)) \text{ by } \bullet.
\end{aligned}$$

$$\Rightarrow \frac{1}{2} \|\theta(t) - z\|_2^2 - \frac{1}{2} \|\theta(0) - z\|_2^2 \leq t \hat{L}(z) - t \hat{L}(\theta(t)). \quad \square$$

Lemma For linearly sep. data w/ $y_i \langle \bar{u}, x_i \rangle \geq \gamma$, $\|\bar{u}\| = 1$, $\|x_i\| \leq 1$,

$$\hat{L}(\theta(t)) \leq \frac{1 + \log(2t\gamma^2)}{2t\gamma^2} \quad \Leftrightarrow \quad \|\theta(t)\| \geq \log(2t\gamma^2) - \log(1 + \log(2t\gamma^2))$$

Pf take $z = \log(c) \bar{u} \gamma^{-1}$ (for some $c > 0$ to be determined)

in the preceding theorem to get:

$$\hat{L}(\theta(t)) \leq \hat{L}(z) + \frac{1}{2t} (\|z\|^2 - \|\theta(t) - z\|^2)$$

$$\leq \sum_i \ell(m_i(z)) + \frac{\|z\|^2}{2t}$$

$$\leq \sum_i \exp(-\log(c)) + \frac{\log^2(c)}{2t\gamma^2}$$

$$= \frac{n}{c} + \frac{\log^2(c)}{2t\gamma^2}$$

Take $c := 2t n \gamma^2$.

$$\begin{aligned} m_i(z) &= \log(c) \gamma^{-1} y_i \langle \bar{u}, x_i \rangle \\ &\geq \log(c) \end{aligned}$$

as \uparrow

For lower bound on $\|\theta(t)\|$, note that

$$\|\theta(t)\| \geq y_i \langle \theta, x_i \rangle = m_i(\theta) \text{ for any } i \text{ by C-S } \hat{=} \|x_i\| \leq 1.$$

$$\begin{aligned} \Rightarrow \ell(\|\theta(t)\|) &\leq \ell(\max_i m_i(\theta(t))) \\ &= \min_i \ell(m_i(\theta(t))) \end{aligned}$$

$$\leq \frac{1}{n} \hat{L}(\theta(t)) \leq \frac{1 + \log^2(2t n \gamma^2)}{2t n \gamma^2} \quad \text{by } \cdot$$

Now take ℓ^{-1} of both sides $\hat{=}$ use ℓ^{-1} decreasing.

Thus:

$$\textcircled{1} \hat{L}(\theta_t) \rightarrow 0$$

$$\textcircled{2} \|\theta_t\| \rightarrow \infty.$$

Now we show margin maximiz.

Theorem Consider linearly sep. data w/ exp loss \hat{L} $\|x_i\| \leq 1$. Then:

$$\gamma(\theta_t) \approx \tilde{\gamma}(\theta_t) \geq \bar{\gamma} - \frac{\log n}{\log t + \log(2n\gamma^2) - 2\log\log(2n\gamma^2)}$$

Pf: Let $u(t) := \ell^{-1}(\hat{L}(\theta_t))$, $v(t) := \|\theta_t\|$.

Thus,

$$\tilde{\gamma}(\theta_t) = \frac{u(t)}{v(t)} = \frac{u(0) + \int_0^t \frac{du(s)}{ds} ds}{v(t)}$$

Want: $u(t)$ grows fast, $v(t)$ not too large.

Since $-\ell' = \ell$,

$$\frac{du}{dt} = \frac{d}{dt} \ell(\hat{L}(\theta_t)) = \left\langle \frac{-\nabla \hat{L}(\theta_t)}{\hat{L}(\theta_t)}, \frac{d\theta}{dt} \right\rangle = \frac{\|d\theta/dt\|^2}{\hat{L}(\theta_t)}$$

$$\begin{aligned} \left\| \frac{d\theta}{dt} \right\| &\geq \left\langle \frac{d\theta}{dt}, \bar{u} \right\rangle = \left\langle \sum_{i=1}^n -\ell'(m_i(\theta)) y_i x_i, \bar{u} \right\rangle \\ &= \left\langle \sum_{i=1}^n \ell(m_i(\theta)) y_i x_i, \bar{u} \right\rangle \\ &\geq \gamma \sum_{i=1}^n \ell(m_i(\theta)) = \gamma \hat{L}(\theta). \end{aligned}$$

$$0 \leq t \leq T, \quad v(t) = \|\theta(t)\| = \|\theta(t) - \theta(0)\| = \left\| \int_0^t \frac{d\theta(s)}{ds} ds \right\| \leq \int_0^t \left\| \frac{d\theta(s)}{ds} \right\| ds.$$

$$\Rightarrow \frac{\int_0^t \frac{dv}{ds} ds}{v(t)} \geq \frac{\int_0^t \frac{\|d\theta/ds\|^2}{\hat{L}(\theta(s))} ds}{\int_0^t \|d\theta/ds\| ds} = \frac{\int_0^t \|d\theta/ds\| \cdot \frac{\|d\theta/ds\|}{\hat{L}(\theta(s))} ds}{\int_0^t \|d\theta/ds\| ds}$$

$$\stackrel{\text{Cauchy-Schwarz}}{\geq} \frac{\int_0^t \|d\theta/ds\| \cdot \gamma ds}{\int_0^t \|d\theta/ds\| ds} = \gamma.$$

$$\text{Now: } \frac{v(0)}{v(t)} = \frac{-\log(\hat{L}(\theta(0)))}{\|\theta(t)\|} = \frac{-\log(\eta)}{\|\theta(t)\|} \stackrel{\text{prev. lemma}}{\geq} \frac{-\log \eta}{\log t + \log(2\eta r^2) - \log \log(2\eta e r^2)}$$

Putting into γ we get claimed thm. \square

This shows smooth normalized margin \rightarrow max margin as $t \rightarrow \infty$.

Since $\gamma(\theta) \geq \tilde{\gamma}(\theta)$, this implies $\theta(t)$ achieves max margin as $t \rightarrow \infty$.