

Wednesday 07/28: TA session 1 (THEODOR MISIAKIEWICZ)
(STANFORD)

Complement to ANDREA MONTANARI's lectures

① Today: Implicit regularization

② Friday: "Benign overfitting in linear regression" ✓
[Bartlett, Long, Lugosi, Trigler, 2019]

③ Next Tuesday: Going beyond the linear regime (?)

① IMPLICIT REGULARIZATION (or BIAS)

ANDREA's lectures: linear regime

→ implicit regularization of GD: minimum ℓ_2 -norm interpolating solution

⇒ More generally: IR plays a crucial role in the generalization properties of overparametrized models

Setting:

- Supervised learning: $\{(y_i, \alpha_i)\}_{i \leq n}$ data

$$\alpha_i \in \mathbb{R}^d, \quad y_i \in \mathbb{R} \quad (\text{or } y_i \in \{\pm 1\})$$

- Parametric family of models: $\mathcal{F} = \{f(\cdot, \theta) : \theta \in \mathbb{R}^P\}$

$f(\cdot, \theta) : \mathbb{R}^d \rightarrow \mathbb{R}$

$x \mapsto f(x, \theta)$ parameters
 $\theta \in \mathbb{R}^P$
 input

- Learning: * loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

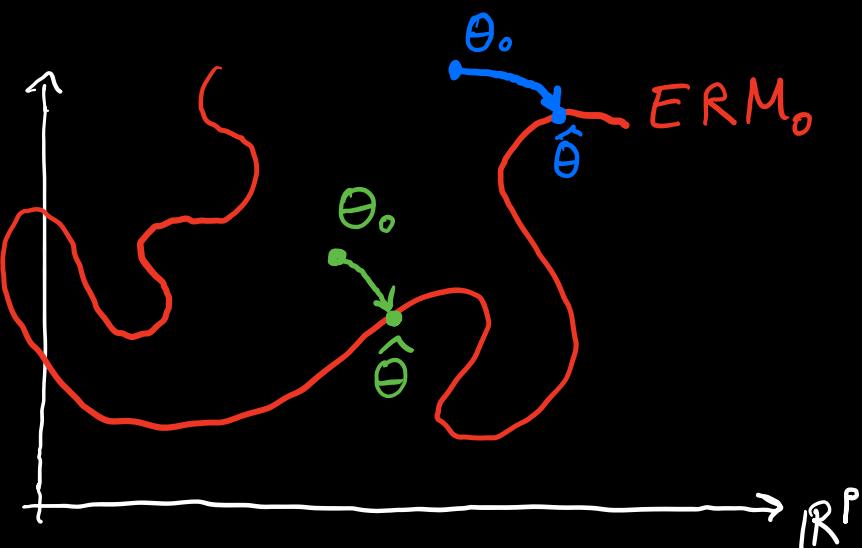
* Empirical Risk Minimization (ERM):

$$\text{minimize } \hat{R}_m(\theta) = \frac{1}{m} \sum_{i=1}^m \ell(y_i, f(x_i, \theta))$$

Modern approach: GD / SGD $\theta^t \leftarrow \sum_{i=1}^m (\langle x_i, \theta_* \rangle - \ell(x_i, \theta))^2$
 $\hat{R}_m(\theta^t) \rightarrow 0$ ("INTERPOLATION")
 $\|X(\theta_* - \theta)\|_2^2 + \lambda \|\theta\|_2^2$

Overparametrized models:

$\hat{f}(x_i, \theta) = y_i$ $\begin{matrix} \text{m} \times p \text{ matrix} \\ \text{m} \ll p \end{matrix} \quad \begin{matrix} \text{p vector} \\ \theta = \theta_* + \omega \end{matrix}$
 θ_* is large $p-m$ subspace



Empirical minimizers generalize more or less well

→ Opt algo introduce a bias in this choice

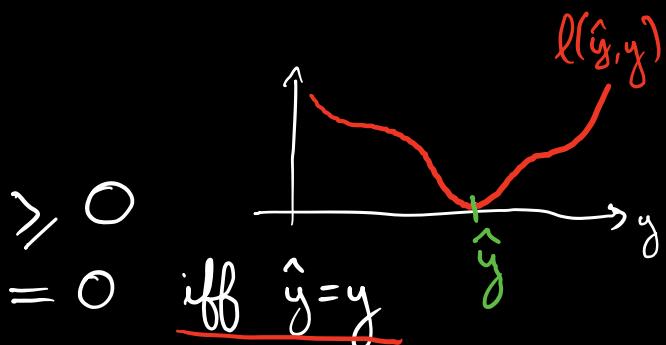
\Rightarrow Opt algorithm selects solution $\hat{\theta} \in ERM_0$

- This bias is really understood in only a few models.
 - Here focus on linear models.
- "Characterizing implicit bias in terms of optimization geometry"
- Ganeshan, Lee, Soudy, Srebro (2018)

- Linear models: $f(\alpha, \Theta) = \langle \alpha, \Theta \rangle \quad \alpha, \Theta \in \mathbb{R}^P$
- ERM: $\hat{R}_m(\Theta) = \frac{1}{m} \sum_{i=1}^m l(y_i, \langle \alpha_i, \Theta \rangle)$

2 family of losses:

① Regression setting: $l(y, \hat{y}) \geq 0$



$$ERM_0 = \{ \Theta \in \mathbb{R}^P : \langle \alpha_i, \Theta \rangle = y_i \quad \forall i \leq m \}$$

$$= \{ \Theta \in \mathbb{R}^P : \underline{X\Theta = y} \}$$

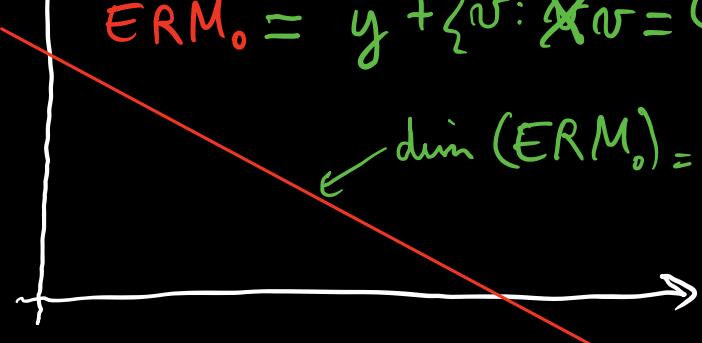
$$y = (y_1, \dots, y_m) \in \mathbb{R}^m$$

$$ERM_0 = y + \{ w : \cancel{Xw = 0} \}$$

$$X = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix} \in \mathbb{R}^{m \times P}$$

$$\dim(ERM_0) = p - m$$

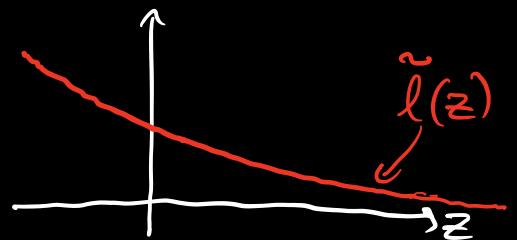
(if X full row rank)



② Classification setting: $y \in \{-1, +1\}$ $\ell(\hat{y}, y) = \tilde{\ell}(\hat{y}y)$

$$\tilde{\ell}(z) > 0, \quad \ell(z) \rightarrow 0 \text{ iff } z \rightarrow \infty$$

monotonically decreasing



e.g.: $\tilde{\ell}(z) = e^{-z}$ (exp loss)

$\tilde{\ell}(z) = \log(1 + e^{-z})$ (logistic)

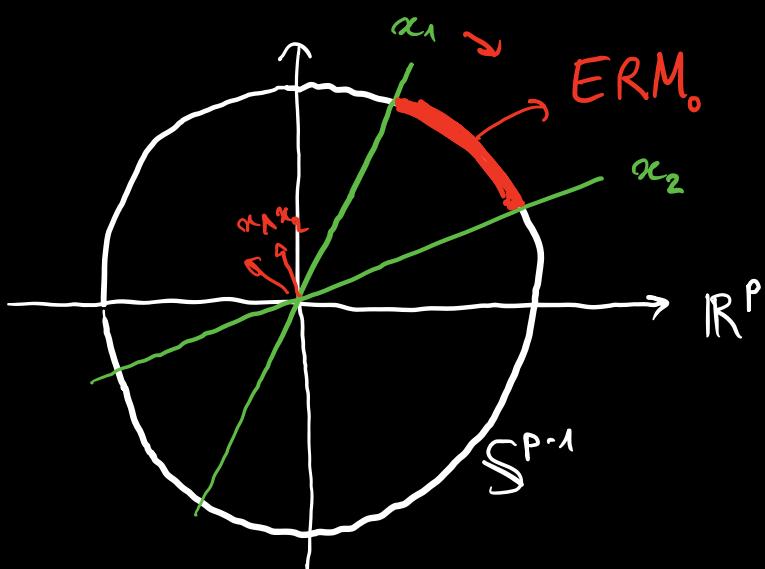
→ In this case: $\hat{R}_m(\theta) > 0 \quad \forall \theta \text{ finite}$

→ Prediction: $\hat{f}(\alpha, \theta) = \text{sign}(\langle \theta, \alpha \rangle) = \text{sign}\left(\underbrace{\langle \frac{\theta}{\|\theta\|_2}, \alpha \rangle}_{=\beta}\right)$

$\beta \in \mathbb{S}^{p-1}$: unit sphere in p dimension

$$\hat{f}(\alpha_i, \theta) y_i \geq 0$$

$$\text{ERM}_0 := \left\{ \beta \in \mathbb{S}^{p-1} : \langle \beta, \alpha_i \rangle y_i \geq 0 \quad \forall i \leq m \right\}$$



① Regression setting: (assume ℓ is differentiable)

Mirror descent (MD):

- potential: $\Psi: \mathbb{R}^P \rightarrow \mathbb{R}$ differentiable + strictly convex
- Bregman divergence w.r.t. Ψ :

$$D_\Psi(\theta, \theta_0) := \Psi(\theta) - \Psi(\theta_0) - \underbrace{\langle \nabla \Psi(\theta_0), \theta - \theta_0 \rangle}_{> 0}$$

$D_\Psi(\cdot, \theta_0)$: strictly convex with unique minimizer at $\theta = \theta_0$.

- MD algo:

- initialization $\theta_0 \leftarrow$
- step size $\eta_t \leftarrow$
- update:

$$\theta^{t+1} = \underset{\theta \in \mathbb{R}^P}{\operatorname{argmin}} \left\{ \eta_t \underbrace{\langle \theta, \nabla \hat{R}_m(\theta_t) \rangle}_{+ D_\Psi(\theta, \theta^t)} + D_\Psi(\theta, \theta^t) \right\}$$

Example: 1) $\Psi(\theta) = \frac{1}{2} \|\theta\|_2^2$ $D_\Psi(\theta, \theta_0) = \frac{1}{2} \|\theta - \theta_0\|_2^2$

$$\theta^{t+1} = \theta^t - \eta_t \nabla \hat{R}_m(\theta^t) \quad \text{Gradient Descent}$$

2) $\theta \in \mathbb{R}_{>0}^P$ (positive orthant)

$$\Psi(\theta) = \sum_{i=1}^f \theta_i \log \theta_i \leftarrow \alpha \log n$$

$$\Delta = \left\{ \alpha \in \mathbb{R}_{>0}^P \mid \alpha^\top \alpha = 0 \right\}$$

- * Lots of work on conditions such that MD converges to global optimum
- * Question here: which point do we converge to?
 \exists interpolator

Proposition: Assume $\{\theta : X\theta = y\} \neq \emptyset$.

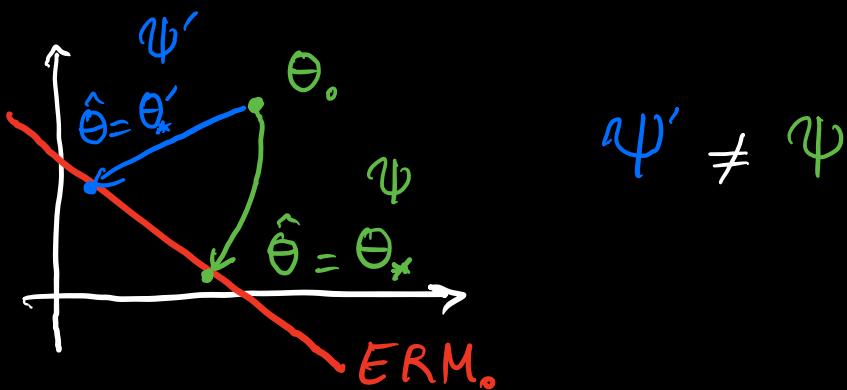
Consider MD with initialization θ^0 .

Assume $\hat{R}_n(\theta^t) \xrightarrow{t \rightarrow \infty} 0$ and $(\theta^t)_{t \geq 1}$ remains bounded

Then:

$$\lim_{t \rightarrow \infty} \theta^t = \theta_* := \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ D_\psi(\theta, \theta_0) : \hat{R}_n(\theta) = 0 \right\}$$

Remark:



Proof: KKT conditions of $\theta_* = \underset{\theta}{\operatorname{argmin}} \left\{ D_\psi(\theta, \theta_0) : \hat{R}_n(\theta) = 0 \right\}$

Lagrangian: $L(\theta, u) = \underbrace{D_\psi(\theta, \theta_0)}_{\text{green}} - \langle u, \underbrace{X\theta - y}_{\text{green}} \rangle$

$$D_\psi(\theta, \theta_0) = \underbrace{\psi(\theta) - \psi(\theta_0)}_{\text{green}} - \underbrace{\langle \nabla \psi(\theta_0), \theta - \theta_0 \rangle}_{\text{green}}$$

$$\rightarrow \nabla_\theta L(\theta, u) = \underbrace{\nabla \psi(\theta)}_{\text{green}} - \underbrace{\nabla \psi(\theta_0)}_{\text{green}} - \underbrace{X^\top u}_{\text{green}}$$

$$\begin{aligned} \text{RHS: } & \quad \left\{ \begin{array}{l} (\nabla \Psi(\theta^*) - \nabla \Psi(\theta_0) = X^T u \\ X \theta^* = y \end{array} \right. \\ \text{KKT: } & \quad \left\{ \begin{array}{l} \end{array} \right. \end{aligned} \quad (I)$$

$$\text{MD: } \theta^{t+1} = \underset{\theta}{\operatorname{argmin}} \left\{ \eta_t \langle \theta, \nabla \hat{R}_m(\theta^t) \rangle + D_\psi(\theta, \theta^t) \right\}$$

$$\text{KKT: } \nabla \Psi(\theta^{t+1}) - \nabla \Psi(\theta^t) + \eta_t \nabla \hat{R}_m(\theta^t) = 0$$

$$\nabla \hat{R}_m(\theta) = \frac{1}{m} \sum_{i=1}^m l'(y_i, \langle \theta, \alpha_i \rangle) \alpha_i = X^T n(\theta)$$

$$n(\theta) = (n_1(\theta), \dots, n_m(\theta)) \quad n_i(\theta) = \frac{1}{m} l'(y_i, \langle \alpha_i, \theta \rangle)$$

$$\text{So we can write: } \nabla \Psi(\theta^{t+1}) - \nabla \Psi(\theta^t) = -\eta_t X^T n(\theta^t)$$

$$\rightarrow \nabla \Psi(\theta^t) - \nabla \Psi(\theta^0) = X^T \left(- \sum_{s=0}^{t-1} \eta_s n(\theta^s) \right) \quad (*)$$

By compactness: θ^t w along a subsequence $\theta^{t_k} \rightarrow \underline{\theta^s} \in \mathbb{R}^P$

By differentiability: $\nabla \Psi(\theta^{t_k}) \rightarrow \nabla \Psi(\underline{\theta^s})$

By (*): $v^{t_k} \rightarrow v^s$ + $v^{t_k} \in \operatorname{Im}(X^T)$ closed
 $\Rightarrow v^s \in \operatorname{Im}(X^T)$

$$v^{t_k} \rightarrow X^T u^s \quad v^s = X^T u^s$$

$$\left\{ \begin{array}{l} \underline{\nabla \Psi(\Theta^s) - \nabla \Psi(\Theta^o) = X^T u^s} \\ \underline{X \Theta^s = y} \quad (\text{because } R(\Theta^s) = 0) \end{array} \right.$$

$R(\Theta_t) \rightarrow 0$
assumption

\equiv KKT condition (I) (of Θ_*)

D_Ψ strictly convex \Rightarrow unique solution $\Theta^s = \Theta_*$

Every subsequence $\underline{\Theta^{t_n}}$ wrt to Θ_* $\Rightarrow \underline{\underline{\Theta^t}} \rightarrow \underline{\underline{\Theta_*}}$ □

Summary: if $\hat{R}_m(\Theta^t) \rightarrow 0 + \{\Theta^t\}$ bounded then

$$\Theta^t \xrightarrow[t \rightarrow \infty]{} \Theta_* = \underset{\Theta}{\operatorname{argmin}} \left\{ D_\Psi(\Theta, \Theta^o) : \hat{R}_m(\Theta) = 0 \right\}$$

"Minimum Bregman divergence interpolating solution"

Hence $\hat{\Theta}^*$ depends on :

- ① Algorithm (choice of Ψ)
- ② Initialization (choice of Θ^o)
- ③ Parametrization ϵ

Parametrization example: $F: \mathbb{R}^P \rightarrow \mathbb{R}^P$ $A \in \mathbb{R}^{P \times P}$ invertible
 $\theta \mapsto A\theta$

→ approach 1: GD on θ init $\theta^0 = 0$

→ approach 2: GD on $\tilde{\theta} = A\theta$, $\tilde{\theta}^0 = 0$

$$\hat{\theta}^{(1)} = \operatorname{argmin} \left\{ \|\theta\|_2^2 : X\theta = y \right\} \quad (1) \quad \begin{matrix} X^\top \\ AX^\top \end{matrix}$$

$$\hat{\tilde{\theta}}^{(2)} = \operatorname{argmin} \left\{ \|\tilde{\theta}\|_2^2 : XA^{-1}(\tilde{\theta}) = y \right\} \quad \begin{matrix} X \\ XA^{-1} \end{matrix}$$

$$\begin{aligned} \hat{\theta}^{(2)} &= \underbrace{A^{-1}(\hat{\theta}^{(1)})}_{=} \\ &= \operatorname{argmin} \left\{ \|A\theta\|_2^2 : X\theta = y \right\} \quad (2) \end{aligned}$$

Unless A orthogonal $\hat{\theta}^{(2)} \neq \hat{\theta}^{(1)}$ (1) \neq (2)

Other example: $F(\theta) = (\sqrt{|\theta_1|}, \dots, \sqrt{|\theta_p|})$

→ $\theta^0 = \underbrace{\alpha \mathbf{1}}_{\alpha > 0}$ was shown: $\theta^t \rightarrow \operatorname{argmin} \left\{ \|\theta\|_1 : X\theta = y \right\}$ (GD)

[Ganesh et al, 2017]

Kernel method

Remark: ANDREA's lectures

Linear regime:

$$\begin{aligned} \alpha \in \mathbb{R}^P &\mapsto \alpha \in L^2(\mathbb{R}^d) \\ f_{lin}(\alpha, \alpha) &= \int \sigma(\langle \alpha, \theta \rangle) \alpha(\theta) d\theta \\ f(\alpha, \theta) &\simeq f_{lin}(\alpha, \alpha, \theta_0) = f(\alpha, \theta_0) + \langle \nabla f(\alpha, \theta_0), \underline{\alpha} \rangle \\ \alpha &= \theta - \theta_0 \end{aligned}$$

$$\text{Denote: } y'_i = y_i - f(x_i, \Theta_0)$$

$$\Phi = \begin{bmatrix} \nabla f(x_1, \Theta_0) \\ \vdots \\ \nabla f(x_m, \Theta_0) \end{bmatrix} \in \mathbb{R}^{m \times p}$$

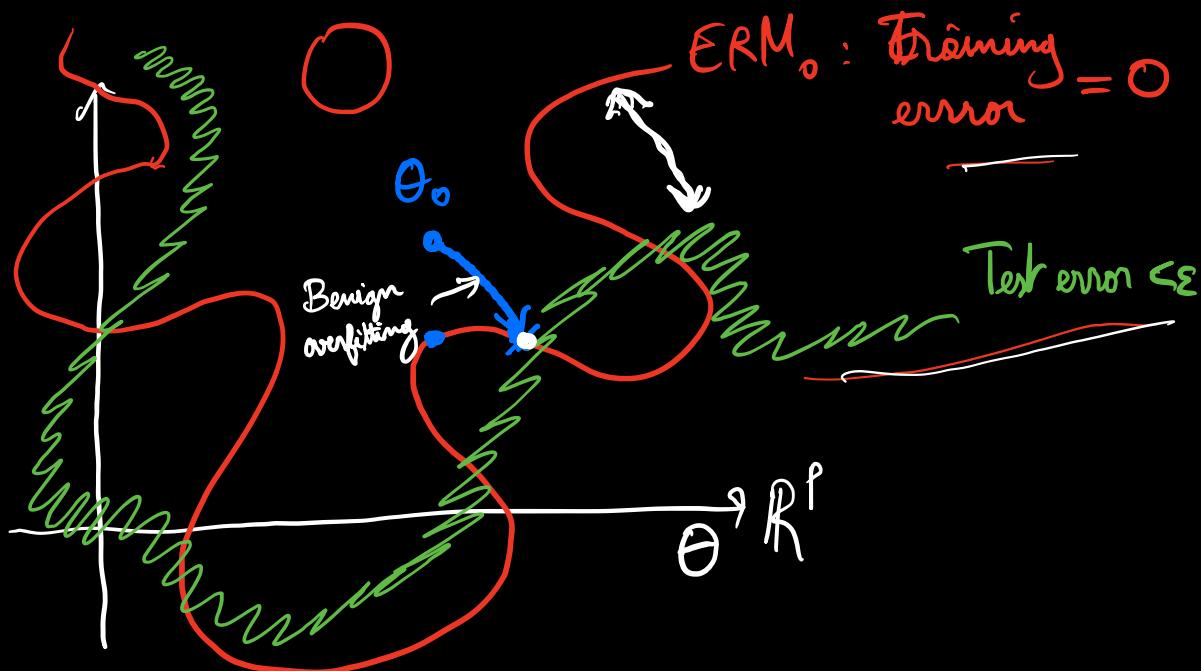
$\leftarrow P \rightarrow$

$$GD \text{ w.r.t. } \hat{a} = \arg\min \left\{ \|a\|_2^2 : y' = \Phi a \right\}$$

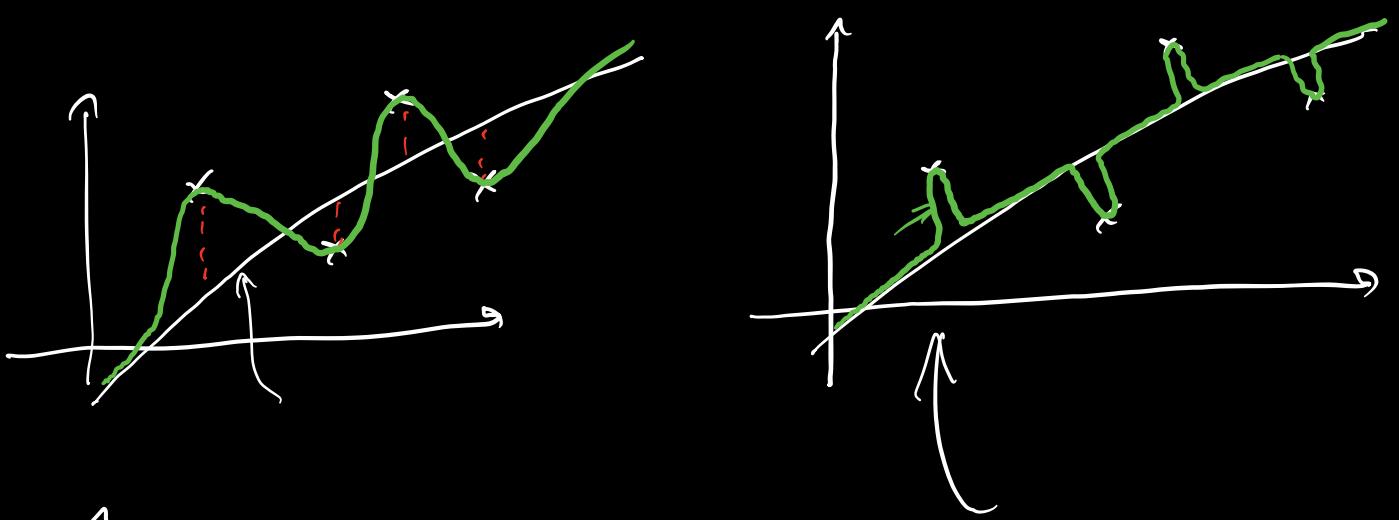
"Ridgeless" regression solution $\hat{a}(0^+)$

$$\text{where } \hat{a}(\lambda) = \arg\min_{a \in \mathbb{R}^p} \left\{ \|y' - \Phi a\|_2^2 + \lambda \|a\|_2^2 \right\}$$

$$\lambda > 0^+$$



$$\hat{\Theta} \rightarrow R_{\text{test}}(\hat{\Theta}) = \mathbb{E}[l(y_{\text{new}}, f(x_{\text{new}}, \hat{\Theta}))]$$



$$\hat{f}(\alpha, \theta) = \langle \alpha, \theta \rangle$$

$$f^*(\alpha, \theta_*) = \langle \alpha, \theta_* \rangle$$

$$y_i = f^*(\alpha_i, \theta_*) + \varepsilon_i$$

$$\hat{f}(\alpha, \hat{\theta}) = \langle \alpha, \hat{\theta} \rangle$$

$$+ \Delta(\alpha)$$

↑
small spikes

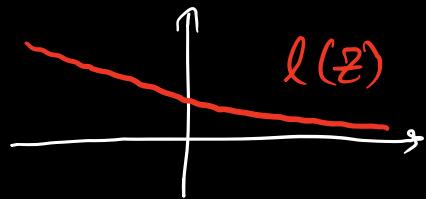
$$\mathbb{E}[\Delta(\alpha)^2] \leq 1$$

② Classification setting:

Setting: $\{(y_i, \alpha_i)\}_{i \in \mathcal{N}} \quad y_i \in \{-1, +1\}$

$$\hat{f}(\alpha, \Theta) = \text{sign}(\langle \Theta, \alpha \rangle)$$

$$\hat{R}_m(\Theta) = \sum_{i=1}^m \ell(\langle \Theta, \alpha_i \rangle y_i)$$



Steepest descent (SD):

- $\|\cdot\|$ norm but not necessarily ℓ_2 -norm

- SD algo:

- initialization Θ^0

- step-sizes η_t

- update:

$$\Theta^{t+1} = \underset{\Theta \in \mathbb{R}^P}{\operatorname{argmin}} \left\{ \eta_t \langle \nabla \hat{R}_m(\Theta^t), \Theta \rangle + \frac{1}{2} \|\Theta - \Theta^t\|^2 \right\}$$

Rewrite: $\Theta^{t+1} = \Theta^t + \eta_t v^t$

where $v^t = \underset{v \in \mathbb{R}^P}{\operatorname{argmin}} \left\{ \langle \nabla \hat{R}_m(\Theta^t), v \rangle + \frac{1}{2} \|v\|^2 \right\}$

Dual norm: $\|\alpha\|_* := \sup_{\|v\| \leq 1} \langle \alpha, v \rangle \quad u_t = -\nabla \hat{R}_m(\Theta^t)$

$$\Rightarrow \min_v \left\{ \frac{1}{2} \|v\|^2 - \langle u_t, v \rangle \right\} = \min_{t \geq 0} \min_{\|v\|=t} \left\{ \frac{1}{2} t^2 - \langle v, u_t \rangle \right\}$$

$$= \min_t \left\{ \frac{1}{2} t^2 - t \|u_t\|_* \right\} = -\frac{1}{2} \|\ \|_*^2 \text{ iff } t = \|u_t\|_*$$

$$\Rightarrow v^t = \underset{v}{\operatorname{argmax}} \left\{ \langle v, u_t \rangle : \|v\| \leq \|u_t\|_* \right\}$$

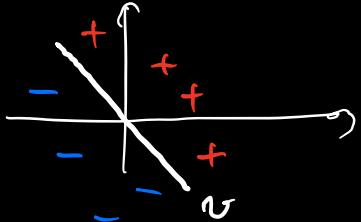
SD update:

$$\begin{cases} \theta^{t+1} = \theta^t - \eta_t v^t \\ v^t = \underset{\|v\| \leq \|\nabla \hat{R}_m(\theta^t)\|_*}{\operatorname{argmax}} \left\{ \langle v, \nabla \hat{R}_m(\theta^t) \rangle \right\} \end{cases}$$

Setting:

- Data is separable:

$$\exists v \in \mathbb{R}^d, y_i \langle \alpha_i, v \rangle > 0 \quad \forall i$$



- $l(z) = e^{-z}$

Thm: Assume separable data + $\|\alpha_i\|_* \leq B \quad \forall i \leq m$

$$\eta_t \leq \left[\frac{1}{B^2 R(\theta^t)} \wedge \eta_{\max} \right]$$

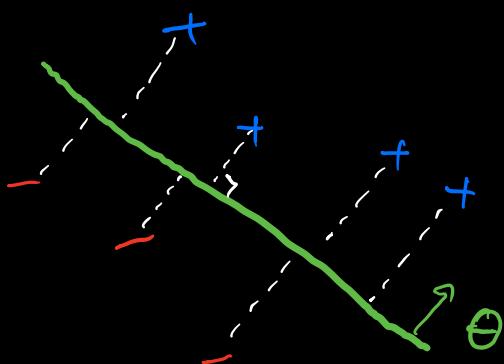
Then

$$\lim_{t \rightarrow \infty} \min_{i \leq m} \frac{y_i \langle \theta^t, \alpha_i \rangle}{\|\theta^t\|} = \max_{\|\theta\| \leq 1} \min_{i \leq m} y_i \langle \theta, \alpha_i \rangle$$

Further, if max on RHS is uniquely achieved at θ_* ,

then $\theta^t \rightarrow \theta_*$ as $t \rightarrow \infty$

Remark: * if $\|\cdot\| = \|\cdot\|_2$ then RMS = max margin classification



* Margin: $\gamma = \min_{\|\theta\|=1} \max_{i \leq m} \langle \theta, \tilde{x}_i \rangle$ ($\tilde{x}_i = y_i x_i$)
 (exercise) $= \min_{n \in \Delta_{m-1}} \|X^T n\|_*$

where $\Delta_{m-1} = \{n \in \mathbb{R}^m : n \geq 0, \langle 1, n \rangle = 1\}$

e.g.: $\hat{R}_m(\theta) = \sum_{i=1}^m e^{-y_i \langle x_i, \theta \rangle} = \sum_{i=1}^m n_i(\theta)$
 $n_i(\theta) = e^{-y_i \langle x_i, \theta \rangle}$

$$\nabla \hat{R}_m(\theta) = - \tilde{X}^T n(\theta)$$

$$\frac{\nabla \hat{R}_m(\theta)}{\hat{R}_m(\theta)} = - \frac{\tilde{X}^T n(\theta)}{\langle 1, n(\theta) \rangle}$$

$$\Rightarrow \frac{\|\nabla \hat{R}_m(\theta)\|_*}{\hat{R}_m(\theta)} \geq \gamma$$

- Proof:
- WLOG: $y_i = +1$
 - Denote for simplicity $R(\theta) := \hat{R}_m(\theta) = \sum_{i=1}^m e^{-\langle \alpha_i, \theta \rangle}$

We will consider continuous time gradient flow:

$$\dot{\theta}_t = v_t : \quad \langle v_t, -\nabla R(\theta_t) \rangle = \|v_t\|^2 = \|\nabla R(\theta_t)\|_*^2$$

Lemma:

- (1) $\int_0^\infty \|\nabla R(\theta_t)\|_*^2 dt < \infty$
- (2) $R(\theta_t) \rightarrow 0$
- (3) $\int_0^\infty \|\nabla R(\theta_t)\|_* dt = \infty$

Proof: (1) $\frac{d}{dt} R(\theta_t) = \langle \nabla R(\theta_t), v_t \rangle = -\|\nabla R(\theta_t)\|_*^2$

$$\int_0^\infty \|\nabla R(\theta_t)\|_*^2 dt \leq R(\theta_0) < \infty$$

(2) (1) $\Rightarrow \|\nabla R(\theta_t)\|_* \rightarrow 0$

+ separability assumption: $\exists v : \langle v, \alpha_i \rangle > 0 \ \forall i$

$$\langle v, R(\theta_t) \rangle = \sum_{i=1}^m e^{-\langle \alpha_i, \theta_t \rangle} \langle \alpha_i, v \rangle \geq R(\theta^*) \cdot \min_{i \leq m} \langle v, \alpha_i \rangle$$

$$\nabla R(\theta_t) \rightarrow 0 \Rightarrow R(\theta_t) \rightarrow 0 \quad (\text{i.e. } \min_i \langle \alpha_i, \theta_t \rangle \rightarrow \infty)$$

$$(3) \quad \|\Theta^t\| \leq \|\Theta^0\| + \int_0^t \|v_s\| ds \\ = \|\Theta^0\| + \int_0^t \|\nabla R(\Theta_s)\|_* ds$$

$$\lim_{t \rightarrow \infty} \|\Theta^t\| = \infty \quad \Rightarrow \quad \int_0^\infty \|\nabla R(\Theta_s)\|_* ds = \infty \quad \square$$

Proof of thm: $\frac{d}{dt} R(\Theta^t) = - \|\nabla R(\Theta^t)\|_*^2$

Hence $- \frac{d}{dt} \log R(\Theta^t) = \frac{\|\nabla R(\Theta^t)\|_*^2}{R(\Theta^t)}$

$$\Rightarrow -\log R(\Theta^T) = -\log R(\Theta^0) + \int_0^T \frac{\|\nabla R(\Theta^t)\|_*^2}{R(\Theta^t)} dt$$

with $-\log \left(\sum_{i=1}^m e^{-\langle \alpha_i, \Theta^T \rangle} \right) \leq -\log (e^{-\langle \alpha_i, \Theta^T \rangle}) = \langle \alpha_i, \Theta^T \rangle$

$$\Rightarrow \min_{i \leq m} \langle \alpha_i, \Theta^T \rangle \geq -\log R(\Theta^0) + \int_0^T \frac{\|\nabla R(\Theta^t)\|_*^2}{R(\Theta^t)} dt$$

$$\|\Theta^T\| \leq \|\Theta^0\| + \int_0^T \|\nabla R(\Theta^t)\|_* dt$$

$$\Rightarrow \min_{i \leq m} \frac{\langle \alpha_i, \Theta^T \rangle}{\|\Theta^T\|} \geq \frac{-\log R(\Theta^0) + \int_0^T \frac{\|\nabla R(\Theta^t)\|_*^2}{R(\Theta^t)} dt}{\|\Theta^0\| + \underbrace{\int_0^T \|\nabla R(\Theta^t)\|_* dt}_{\substack{\rightarrow \infty \\ T \rightarrow \infty}}$$

$$\geq \frac{\int_0^T \|\nabla R(\Theta^t)\|_* \underbrace{\frac{\|\nabla R(\Theta^t)\|_*}{R(\Theta^t)}}_{\geq \delta} dt}{\int_0^T \|\nabla R(\Theta^t)\|_* dt} - o_T(1)$$

$$\geq \gamma - o_T(1)$$

and $\min_{i \leq m} \frac{\langle \alpha_i, \theta^T \rangle}{\|\theta^T\|} \leq \max_{\|\theta\| \leq 1} \min_{i \leq m} \langle \alpha_i, \theta \rangle = \gamma$

Hence $\lim_{T \rightarrow \infty} \left[\min_{i \leq m} \frac{\langle \alpha_i, \theta^T \rangle}{\|\theta^T\|} \right] = \gamma \quad \square$