

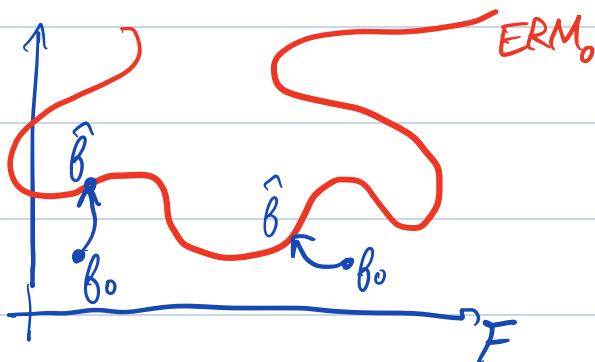
①

Lecture 3: [Implicit / Algorithmic bias]

Setting: $\hat{R}_m(f) = \frac{1}{m} \sum_{i=1}^m l(y_i, f(x_i))$

$F = \{f(\cdot; \theta) : \theta \in \mathbb{R}^P\}$ overparametrized models
 $P \gg m$

$\text{ERM}_0 := \{f \in F : \hat{R}_m(f) = 0\}$



→ ERM minimizer generalizes more or less well

→ Opt algo introduce a bias in this choice

⇒ Opt algorithm selects a solut° $f \in \text{ERM}_0$

Is that a fair description of what NNs do in practice?

- For LLMs $\text{data} = \infty$ (for now) : doesn't train until interpolation regime
- Previously : not exactly trained until interpolation but still much closer to interpolation learning than UC
- Regime to "magnify" new properties of NNs : see [Misha Balakin : Fit without fear review paper 2021]

(2)

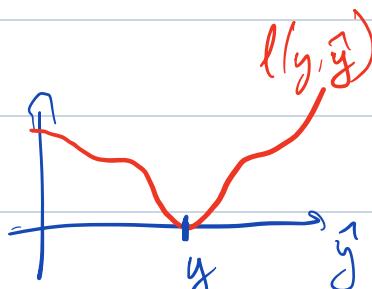
- * The bias is really understood in only a few models
- * Here focus on linear models

"Characterizing implicit bias in terms of optimization geometry"

- Gunasekar, Lee, Soudry, Srebro (2018)

- Linear models: $f(x, \theta) = \langle x, \theta \rangle$ $x, \theta \in \mathbb{R}^P$

- ERM: $\hat{R}_m(\theta) = \frac{1}{m} \sum_{i=1}^m l(y_i, \langle x_i, \theta \rangle)$



2 Family of losses:

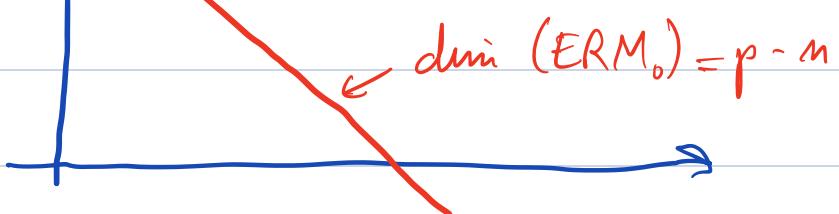
① Regression setting: $y \in \mathbb{R}$ $l(y, \hat{y}) \geq 0$

$$= 0 \text{ iff } \hat{y} = y$$

$$\text{ERM}_0 = \{\theta \in \mathbb{R}^P : \langle x_i, \theta \rangle = y_i \forall i \in [m]\} = \{\theta \in \mathbb{R}^P : X\theta = y\}$$

$$\text{ERM}_0 = y + \{v : Xv = 0\}$$

$$X = \begin{bmatrix} -x_1 \\ -x_2 \\ \vdots \\ -x_m \\ \end{bmatrix}$$

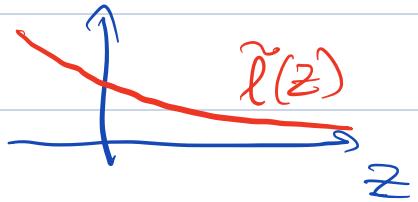


(if X is full rank)

(3)

② Classification setting $y \in \{+1, -1\}$

$\ell(\hat{y}, y) = \tilde{\ell}(\hat{y}y)$ $\tilde{\ell}(z) > 0$ $\tilde{\ell}(z) \rightarrow 0$ iff $z \rightarrow \infty$
 monotonically decreasing



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

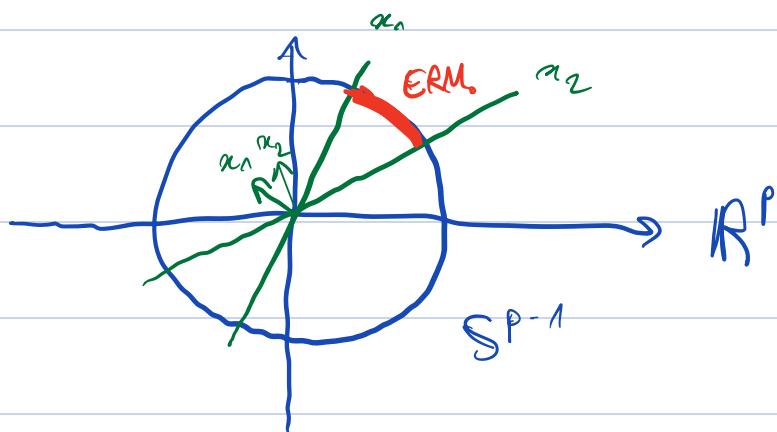
$\tilde{\ell}(z) = \log(1 + e^{-z})$ (logistic)

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

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

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

$\widetilde{\text{ERM}}_{\alpha} := \{ \exists \in \mathbb{S}^{p-1} : \langle \exists, \alpha_i \rangle y_i \geq 0 \quad \forall i \leq m \}$



interpolating
here means
 $\text{sign}(\langle \theta, \alpha_i \rangle) = y_i$

①

Regression setting

④

(Assume l is differentiable)

Morion descent (MD):

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

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

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

- * MD algo:
 - initialization θ_0
 - step size η_t
 - update:

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

(S)

Examples:

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

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

Gradient Descent

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

$$\Psi(\theta) = \sum_{i=1}^p \theta_i \log \theta_i$$

Rmk: * MD can adapt to the geometry of \mathbb{H}

* lots of work on conditions such that MD cr to global optimum

* Question here:

which point do we converge to?

Theorem [Gomeskar et al., '18]

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

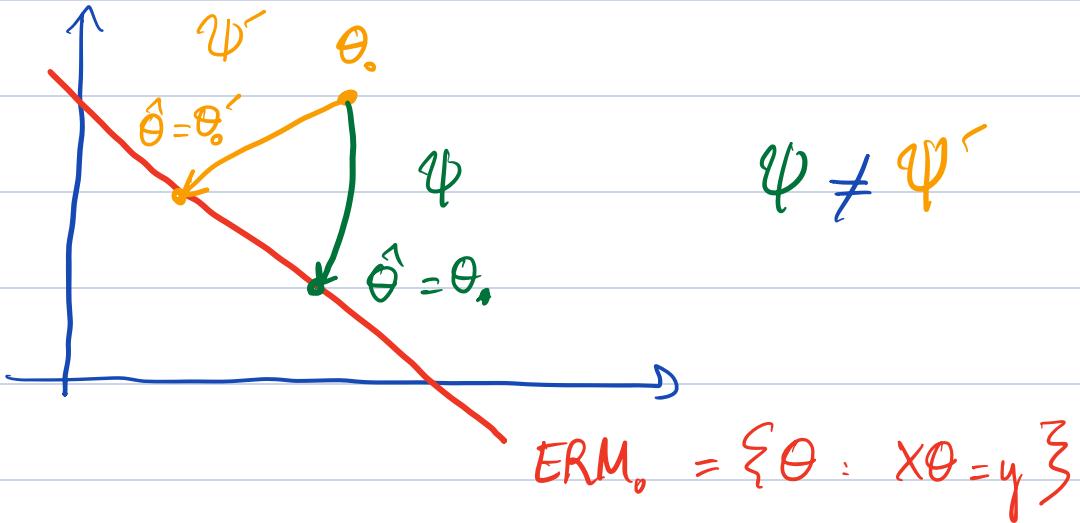
Consider MD with initialization θ^0

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

Then

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

(6)

Rmk:

Proof: KKT conditions of $\theta_* = \arg\min \{D_{\psi}(\theta, \theta_0) : \hat{R}_m(\theta) = 0\}$

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

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

$$\rightarrow \nabla_{\theta} \mathcal{L}(\theta, u) = \nabla \psi(\theta) - \nabla \psi(\theta_0) - X^T u$$

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

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

(7)

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

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

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

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

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

$\underbrace{\qquad\qquad\qquad}_{= v^t}$

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

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

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

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

$$v^s = X^T u^s$$

$$\left\{ \begin{array}{l} \nabla \mathcal{D}(\theta^s) - \nabla \mathcal{D}(\theta^*) = X^T u^s \\ X \theta^s = y \end{array} \right.$$

$R(\theta^*) \rightarrow 0$ by
 assumption


(because $R(\theta^s) = 0$)

\equiv KKT condition (I) (of θ_*)

\mathcal{D}_ψ strictly convex \Rightarrow unique solution $\theta^s = \theta_*$

Every subsequence $\theta^{t_k} \subset v \rightarrow \theta_*$ $\Rightarrow \theta^* \rightarrow \theta_*$



Summary: If $\hat{R}_n(\theta^*) \rightarrow 0$ & $\{\theta^*\}$ bounded

Then $\boxed{\theta^* \xrightarrow[t \rightarrow \infty]{} \theta_* = \operatorname{argmin}_\theta \left\{ \mathcal{D}_\psi(\theta, \theta^*) : \hat{R}_n(\theta) = 0 \right\}}$

"Minimum Bregman divergence interpolating solution"

Hence: $\hat{\theta} = \theta_*$ depends on ① Alg (choice of ψ)

② Initializat[°] (choice of θ^*)

③ Parametrization

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

3

→ approach 1: GD on Θ mit $\Theta^* = 0$

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

$$\hat{\theta}^{(1)} = \operatorname{arg\min} \left\{ \|\theta\|_2^2 : X\theta = y \right\}$$

$$\hat{\vec{\theta}}^{(2)} = \arg\min \left\{ \|\vec{\theta}\|_2^2 : X A^{-1} \vec{\theta} = y \right\}$$

$$\hat{\theta}^{(2)} = A^{-1}(\hat{\theta}^{(1)}) = \underset{\theta}{\operatorname{arg\,min}} \left\{ \|A\theta\|_2 : X\theta = y \right\}$$

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

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

$\Theta^* = \alpha I$ was shown $\Theta^* \rightarrow \operatorname{argmin} \{ \| \Theta \|_1 ; X\Theta = y \}$

[Gunesekar et al., 2017]

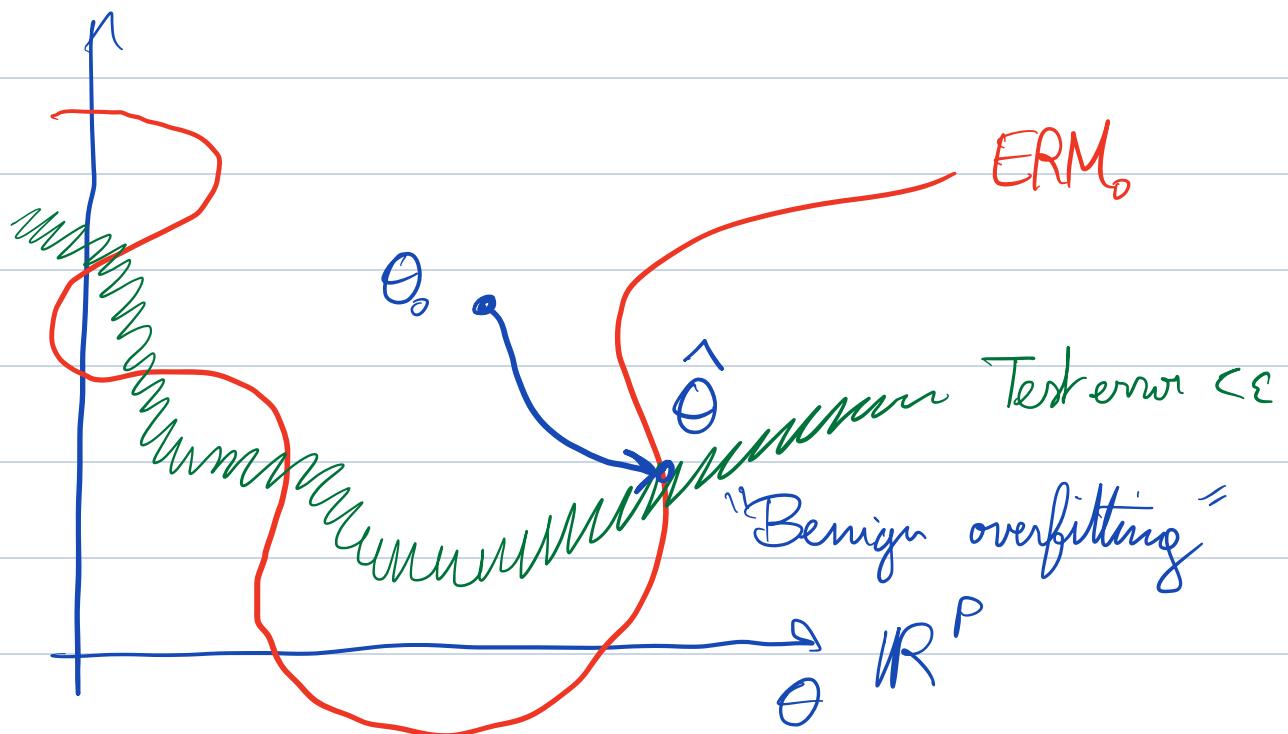
Next few lectures: $f(x; \theta) = \langle \theta, \Phi(x) \rangle$

GD cv to $\hat{\theta} = \arg\min \left\{ \| \theta \|^2_2 : y = \Phi(X)\theta \right\}$
 $\theta_0 = 0$

$$\hat{\theta} = (\Phi(X)\Phi(X)^T)^+ \Phi(X)^T y$$

Implicit regularization: which $\hat{\theta} \in \text{ERM}_0$

What about generalization? Next lecture!



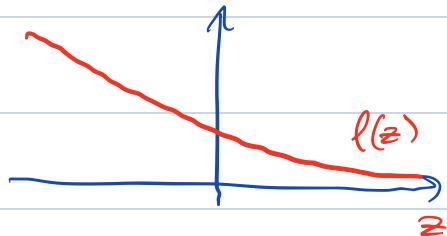
②

CLASSIFICATION SETTING

Setting: $\{(y_i, \alpha_i)\}_{i \leq m}$ $y_i \in \{+1, -1\}$

$$\hat{f}(x; \theta) = \text{sign}(\langle \theta, x \rangle)$$

$$\hat{R}_m(\theta) = \sum_{i=1}^m l(y_i \langle \theta, x_i \rangle)$$



Steepest descent (SD):

* $\|\cdot\|$ norm (not necessarily $\|\cdot\|_2$ -norm)

* SD algo:

- initialization θ^0
- step size η_t
- update:

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

$$v^t = \frac{\theta^{t+1} - \theta^t}{\eta_t} \text{ this is equivalent to}$$

(12)

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

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

$$\text{Dual norm: } \|x\|_* = \sup_{\|v\| \leq 1} \langle x, v \rangle$$

$$u_t := -\nabla \hat{R}_m(\theta^t)$$

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

$$= \min_{\delta} \left\{ \frac{1}{2} \delta^2 - \delta \|u_t\|_* \right\} = -\frac{1}{2} \|u_t\|_*^2$$

$$\epsilon = \|u_t\|_*$$

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

SD update:

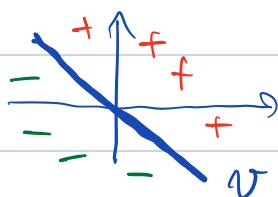
$$\left\{ \begin{array}{l} \theta^{t+1} = \theta^t - \eta_t v^t \\ v^t = \underset{\|v\| \leq \|\nabla \hat{R}_m(\theta^t)\|_*}{\text{argmax}} \langle v, \nabla \hat{R}_m(\theta^t) \rangle \end{array} \right.$$

X why minus

X Standard G(1) if $\|\cdot\|_2$

Setting:

- Data is separable: $\exists v \in \mathbb{R}^d, y_i \langle x_i, v \rangle > 0 \forall i$



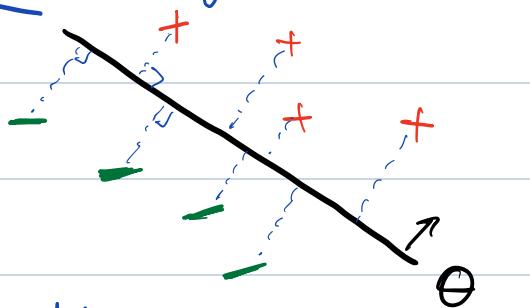
$$\bullet l(z) = e^{-z}$$

THM [Gmosek et al., '18] (B)
 Assume separable data + $\|\alpha_i\|_\infty \leq B \quad \forall i \leq m$
 $\gamma_t \leq \left[\frac{1}{B^2 R(\theta^t)} \wedge \gamma_{\text{max}} \right]$
 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 $t \rightarrow \infty$.

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



$$\tilde{\alpha}_i = y_i \alpha_i$$

* Margin: $\gamma = \max_{\|\theta\|=1} \min_{i \leq m} \langle \theta, \tilde{\alpha}_i \rangle$

$$= \max_{\|\theta\|=1} \min_{i \leq m} \langle e_i, \tilde{\chi} \theta \rangle$$

$$= \max_{\|\theta\| \leq 1} \min_{n \in \Delta_{m-1}} \langle n, \tilde{\chi} \theta \rangle$$

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

14

$$\text{Slater's condition} \\ = \min_{\boldsymbol{\theta} \in \Delta_{m-1}} \max_{\|\boldsymbol{\theta}\| \leq 1} \langle \boldsymbol{\theta}, \tilde{\mathbf{X}}^T \mathbf{n} \rangle$$

$$= \min_{\boldsymbol{\theta} \in \Delta_{m-1}} \|\tilde{\mathbf{X}}^T \mathbf{n}\|_*$$

$$\text{Here: } \hat{R}_m(\boldsymbol{\theta}) = \sum_{i=1}^m e^{-y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle} = \sum_{i=1}^m n_i(\boldsymbol{\theta})$$

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

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

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

Proof: * WLOG $y_i = +1$

* Denote for simplicity $R(\boldsymbol{\theta}) := \hat{R}_m(\boldsymbol{\theta}) = \sum_{i=1}^m e^{-\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle}$

* C° time gradient flow :

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

Lemme:

$$(1) \int_0^\infty \|\nabla R(\theta_t)\|_*^2 dt < \infty$$

$$(2) R(\theta_t) \rightarrow 0$$

$$(3) \int_0^\infty \|\nabla R(\theta_t)\|_* dt = \infty$$

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: $\exists v : \langle v, x_i \rangle > 0 \quad \forall i$

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

$$\nabla R(\theta^t) \rightarrow 0 \Rightarrow R(\theta_t) \rightarrow 0$$

(i.e. $\min_i \langle x_i, \theta_t \rangle \rightarrow \infty$)

$$(3) \quad \|\theta^t\| \leq \|\theta_0\| + \int_0^t \|v_s\| ds \\ = \|\theta_0\| + \int_0^t \|\nabla R(\theta_s)\|_* ds$$

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

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

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

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

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

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

17

$$\Rightarrow \min_{i \in [n]} \frac{\langle \alpha_i, \theta^T \rangle}{\|\theta^T\|} \geq \frac{-\log R(\theta^*) + \int_0^T \frac{\|\nabla R(\theta^t)\|_*^2}{R(\theta^t)} dt}{\|\theta^*\| + \int_0^T \|\nabla R(\theta^t)\|_* dt} \xrightarrow{T \rightarrow \infty} \infty$$

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

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

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