

Model: $F(\mathbf{w}) = \mathbf{h}_{\mathbf{w}}$ **Model Class:** $\mathcal{H} = \text{range}(F)$

$f(\mathbf{w}, x) = \mathbf{h}_{\mathbf{w}}(x)$ = prediction on x with params (“weights”) \mathbf{w}

Linear models: $f(\mathbf{w}, x) = \langle \beta_{\mathbf{w}}, x \rangle$

$F(\mathbf{w}) = \beta_{\mathbf{w}}$

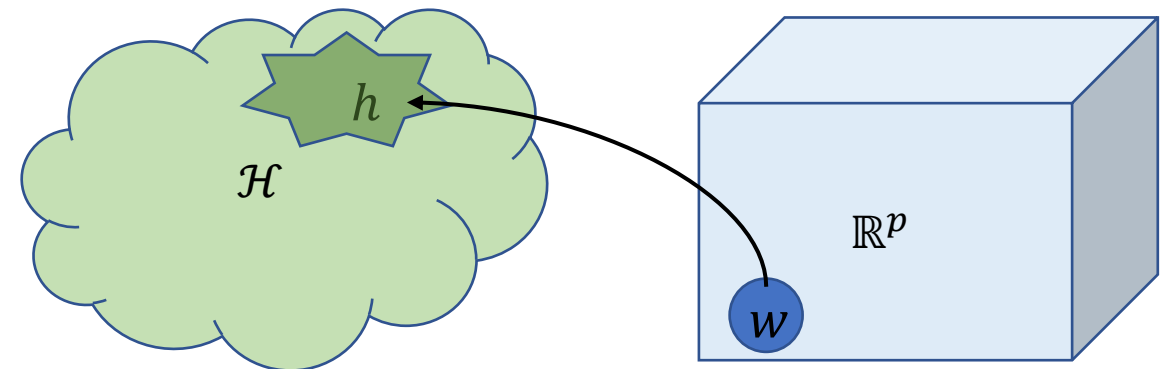
Loss: $L_S(\mathbf{w}) = \frac{1}{m} \sum_i \ell(f(\mathbf{w}, x_i), y_i)$ e.g. $\ell(\hat{y}, y) = (\hat{y} - y)^2$

GD on $L_S(\mathbf{w})$: $\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla_{\mathbf{w}} L_S(\mathbf{w})$ $F(\mathbf{w}_k) \rightarrow ???$

With $\eta \rightarrow 0$: $\dot{\mathbf{w}}(t) = -\nabla_{\mathbf{w}} L_S(\mathbf{w})$ $F(\mathbf{w}(t)) \rightarrow ???$

D-homogenous: $F(c\mathbf{w}) = c^D F(\mathbf{w})$, i.e. $f(c\mathbf{w}, x) = c^D f(\mathbf{w}, x)$

- **1-homogenous:** standard linear $F(\mathbf{w}) = \mathbf{w}$, $f(\mathbf{w}, x) = \langle \mathbf{w}, x \rangle$
- **2-homogenous:**
 - Matrix factorization $F(\mathbf{U}, \mathbf{V}) = \mathbf{UV}$
 - 2-Layer ReLU: $f(\mathbf{w}, x) = \sum_j w_{2,j} [\langle w_{1,j}, x \rangle]_+$
- **D-homogenous:**
 - D layer linear network
 - D layer linear conv net
 - D layer ReLU net



Model: $F(\mathbf{w}) = \mathbf{h}_{\mathbf{w}}$ **Model Class:** $\mathcal{H} = \text{range}(F)$

$f(\mathbf{w}, x) = \mathbf{h}_{\mathbf{w}}(x)$ = prediction on x with params (“weights”) \mathbf{w}

Linear models: $f(\mathbf{w}, x) = \langle \beta_{\mathbf{w}}, x \rangle$

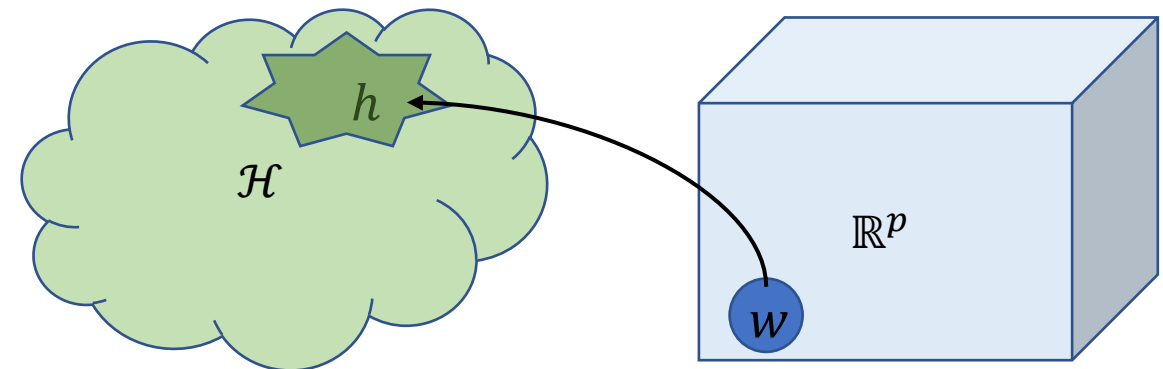
$F(\mathbf{w}) = \beta_{\mathbf{w}}$

Loss: $L_S(\mathbf{w}) = \frac{1}{m} \sum_i \ell(f(\mathbf{w}, x_i), y_i)$ e.g. $\ell(\hat{y}, y) = (\hat{y} - y)^2$

GD on $L_S(\mathbf{w})$: $\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla_{\mathbf{w}} L_S(\mathbf{w})$ $F(\mathbf{w}_k) = \mathbf{h}_{\mathbf{w}_k} \rightarrow ???$

With $\eta \rightarrow 0$: $\dot{\mathbf{w}}(t) = -\nabla_{\mathbf{w}} L_S(\mathbf{w})$ $F(\mathbf{w}(t)) = \mathbf{h}_{\mathbf{w}(t)} \rightarrow ???$

- How is the optimization geometry and dynamics on \mathbf{h} (or β), and the implicit bias effected by the parametrization?
- How is it related, or different, from explicitly $\|\mathbf{w}\|_2$ regularization?
- How is it effect by optimization choices?
- How it is related to the Kernel regime?



Model: $F(\mathbf{w}) = \mathbf{h}_{\mathbf{w}}$ **Model Class:** $\mathcal{H} = \text{range}(F)$

$f(\mathbf{w}, x) = \mathbf{h}_{\mathbf{w}}(x)$ = prediction on x with params (“weights”) \mathbf{w}

Linear models: $f(\mathbf{w}, x) = \langle \beta_{\mathbf{w}}, x \rangle$ $F(\mathbf{w}) = \beta_{\mathbf{w}}$

$$\text{Loss: } L_S(\mathbf{w}) = \frac{1}{m} \sum_i \ell(f(\mathbf{w}, x_i), y_i)$$

Kernel Regime: training behaves according to 1st order approximation about $w^{(0)}$,

$$f(\mathbf{w}, x) \approx f(w_0, x) + \langle \mathbf{w} - w_0, \phi_0(x) \rangle$$

where: $\phi_0(x) = \nabla_w f(w_0, x)$ corresponding to the **tangent kernel**

$$K_0(x, x') = \langle \nabla_w f(w_0, x), \nabla_w f(w_0, x') \rangle$$

(we will focus on “unbiased initialization”: $f(w^{(0)}, x) = 0$, i.e. $h_0 = 0$)

In this regime, $w(t) \rightarrow \text{argmin}_{L_S(w)=0} \|w - w_0\|_2$ and $\mathbf{h}_{w(t)} \rightarrow \min_{h(x_i)=y_i} \|\mathbf{h} - \mathbf{h}_0\|_{K_0}$

[Jacot et al 2018]: Width $\rightarrow \infty$ leads to Kernel Regime

[Chizat Bach 2018]: Scale $\rightarrow \infty$ leads to Kernel Regime

Kernel Regime and Scale of Init

- For D -homogenous model, $f(cw, x) = c^D f(w, x)$, consider gradient flow with:

$$\dot{w}_\alpha = -\nabla L_S(w) \quad \text{and} \quad w_\alpha(0) = \alpha w_0 \quad \text{with unbiased } f(w_0, x) = 0$$

We are interested in $w_\alpha(\infty) = \lim_{t \rightarrow \infty} w_\alpha(t)$

- For squared loss, under some conditions [Chizat and Bach 18]:

$$\lim_{\alpha \rightarrow \infty} \sup_t \left\| w_\alpha \left(\frac{1}{\alpha^{D-1}} t \right) - w_K(t) \right\| = 0$$

Gradient flow of linear least squares w.r.t
tangent kernel K_0 at initialization
 $\dot{w}_K = -\nabla_w \hat{L}(x \mapsto \langle w, \phi_{K_0}(x) \rangle)$

and so $f(w_\alpha(\infty), x) \xrightarrow{\alpha \rightarrow \infty} \hat{h}_K(x)$ where $\hat{h}_K = \arg \min \|h\|_{K_0} \text{ s.t. } h(x_i) = y_i$

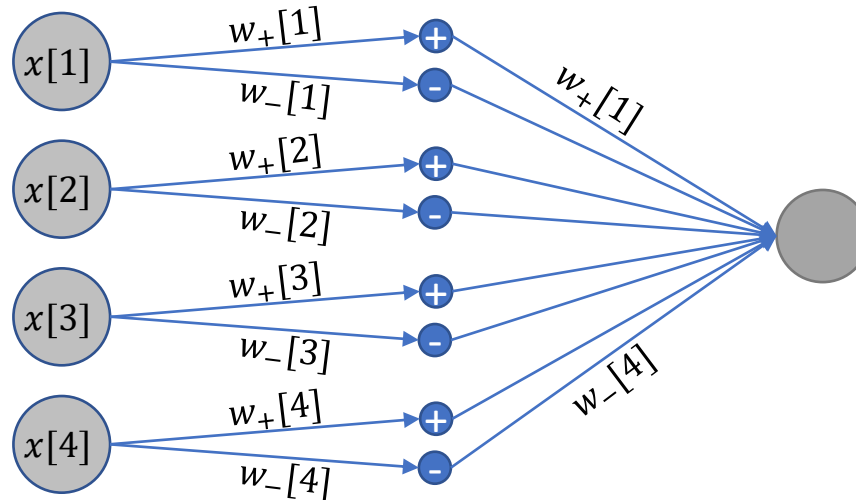
Consider a “linear diagonal net” (ie linear regression with squared parametrization):

$$f(\mathbf{w}, \mathbf{x}) = \sum_j (w_+[j]^2 - w_-[j]^2) x[j] = \langle \beta(\mathbf{w}), \mathbf{x} \rangle \quad \text{with } \beta(\mathbf{w}) = \mathbf{w}_+^2 - \mathbf{w}_-^2$$

And initialization $\mathbf{w}_\alpha(0) = \alpha \mathbf{1}$ (so that $\beta(\mathbf{w}_\alpha(0)) = 0$).

What’s the implicit bias of grad flow w.r.t square loss $L_s(\mathbf{w}) = \sum_i (f(\mathbf{w}, x_i) - y_i)^2$?

$$\beta_\alpha(\infty) = \lim_{t \rightarrow \infty} \beta(\mathbf{w}_\alpha(t))$$



$$f(\mathbf{w}, \mathbf{x}) = \mathbf{w}^\top \text{diag}(\mathbf{w}) \begin{bmatrix} +x \\ -x \end{bmatrix}$$

$$\beta(t) = w_+(t)^2 - w_-(t)^2$$

$$L = \|X\beta - y\|_2^2$$

$$\dot{w}_+(t) = -\nabla_{w_+} L(t) = -2X^\top r(t) \circ \frac{d\beta}{dw_+}$$

$$\beta(t) = w_+(t)^2 - w_-(t)^2$$

$$L = \|X\beta - y\|_2^2$$

$$\dot{w}_+(t) = -\nabla_{w_+} L(t) = -2X^\top r(t) \circ 2w_+(t) \quad w_+(t) = w_+(0) \circ \exp\left(-2X^\top \int_0^t r(\tau) d\tau\right)$$

$$\dot{w}_-(t) = -\nabla_{w_-} L(t) = +2X^\top r(t) \circ 2w_-(t) \quad w_-(t) = w_-(0) \circ \exp\left(+2X^\top \int_0^t r(\tau) d\tau\right)$$

$$\beta(t) = \alpha^2 \left(e^{-4X^\top \int_0^t r(\tau) d\tau} - e^{4X^\top \int_0^t r(\tau) d\tau} \right) \quad r(t) = X\beta(t) - y$$

$$s = 4 \int_0^\infty r(\tau) d\tau \in \mathbb{R}^m$$

$$\beta(\infty) = \alpha^2 \left(e^{-X^\top s} - e^{X^\top s} \right) = 2\alpha^2 \sinh X^\top s$$

$$X\beta(\infty) = y$$

$$\min Q(\beta) \quad s.t. \quad X\beta = y$$

$$\nabla Q(\beta^*) = X^\top \mathbf{v}$$

$$X\beta^* = y$$

$$\beta(\infty) = \alpha^2 \left(e^{-X^\top s} - e^{X^\top s} \right) = 2\alpha^2 \sinh X^\top s$$

$$X\beta(\infty) = y$$

$$\nabla Q(\beta) = \sinh^{-1} \frac{\beta}{2\alpha^2}$$

$$Q(\beta) = \sum_i \int \sinh^{-1} \frac{\beta[i]}{2\alpha^2} = \alpha^2 \sum_i \left(\frac{\beta[i]}{\alpha^2} \sinh^{-1} \frac{\beta[i]}{2\alpha^2} - \sqrt{4 + \left(\frac{\beta[i]}{\alpha^2} \right)^2} \right)$$

$$\min Q(\beta) \quad \text{s.t.} \quad X\beta = y$$

$$\nabla Q(\beta^*) = X^\top \mathbf{v}$$

$$X\beta^* = y$$

$$\sinh^{-1} \frac{\beta(\infty)}{2\alpha^2} = X^\top \mathbf{s}$$

$$X\beta(\infty) = y$$

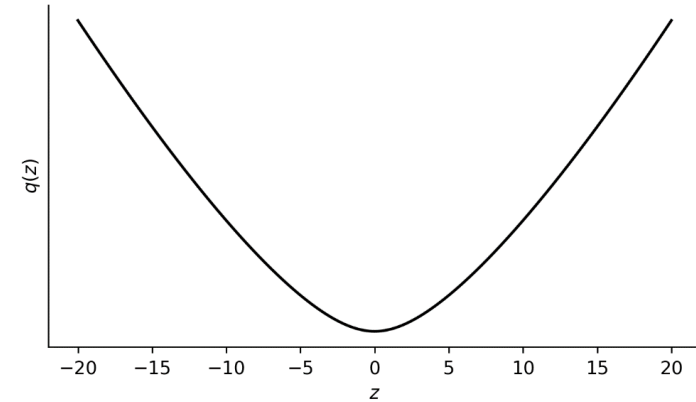
$$f(\mathbf{w}, \mathbf{x}) = \sum_j (\mathbf{w}_+[j]^2 - \mathbf{w}_-[j]^2) x[j] = \langle \boldsymbol{\beta}(\mathbf{w}), \mathbf{x} \rangle \quad \text{with } \boldsymbol{\beta}(\mathbf{w}) = \mathbf{w}_+^2 - \mathbf{w}_-^2$$

$$\boldsymbol{\beta}_\alpha(\infty) = \lim_{t \rightarrow \infty} \boldsymbol{\beta}(\mathbf{w}_\alpha(t)) \quad \text{where } \mathbf{w}_\alpha(0) = \alpha \mathbf{1} \text{ and } \dot{\mathbf{w}}_\alpha = -\nabla \sum_i (f(\mathbf{w}_\alpha, \mathbf{x}_i) - y_i)^2$$

$$\boldsymbol{\beta}_\alpha(\infty) = \arg \min_{\mathbf{x} \boldsymbol{\beta} = \mathbf{y}} Q_\alpha(\boldsymbol{\beta})$$

$$\text{where } Q_\alpha(\boldsymbol{\beta}) = \sum_j q\left(\frac{\beta[j]}{\alpha^2}\right) \text{ and}$$

$$q(b) = 2 - \sqrt{4 + b^2} + b \sinh^{-1}\left(\frac{b}{2}\right)$$



$$\boldsymbol{\beta}_\alpha(\infty) \xrightarrow{\alpha \rightarrow \infty} \hat{\boldsymbol{\beta}}_{L2} = \arg \min_{\mathbf{x} \boldsymbol{\beta} = \mathbf{y}} \|\boldsymbol{\beta}\|_2 \quad \text{“Kernel Regime” with NTK } K_0(\mathbf{x}, \mathbf{x}') = 4\langle \mathbf{x}, \mathbf{x}' \rangle$$

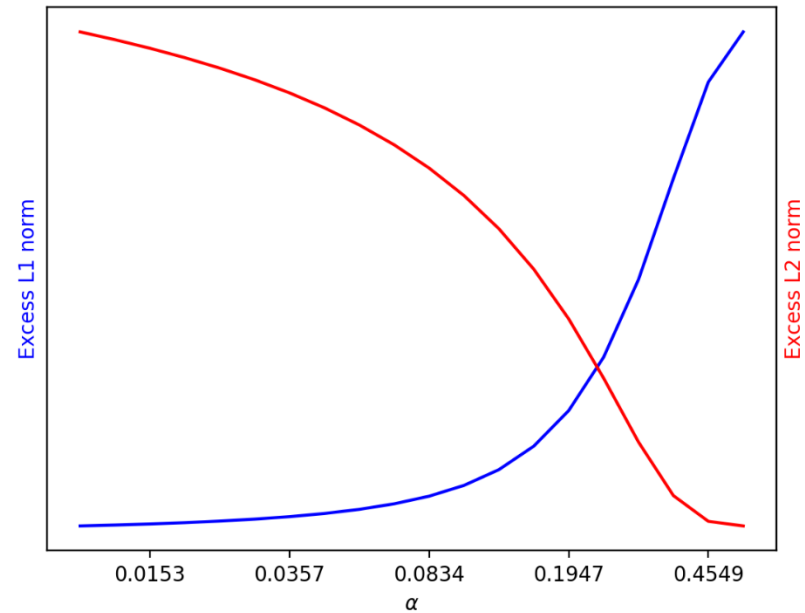
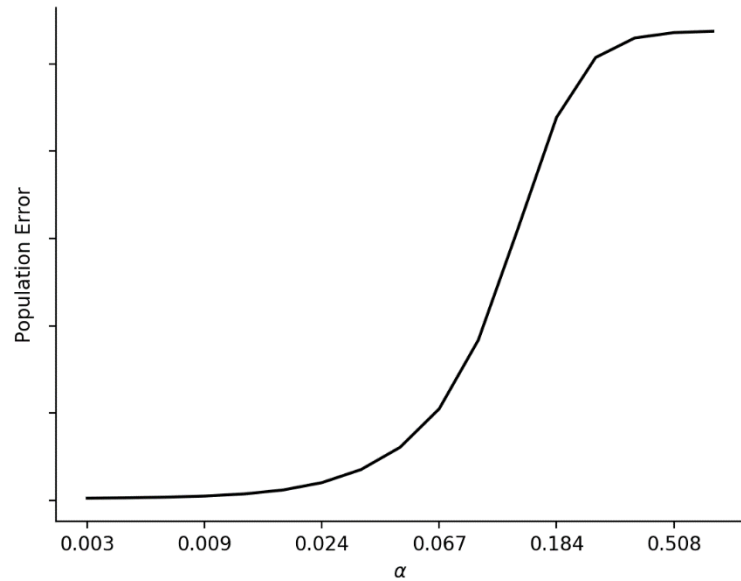
$$\alpha \geq \sqrt{2(1+\epsilon) \left(1 + \frac{2}{\epsilon}\right) \|\boldsymbol{\beta}_{L2}^*\|_2} \implies \|\hat{\boldsymbol{\beta}}_\alpha\|_2^2 \leq (1+\epsilon) \|\boldsymbol{\beta}_{L2}^*\|_2^2$$

$$\boldsymbol{\beta}_\alpha(\infty) \xrightarrow{\alpha \rightarrow 0} \hat{\boldsymbol{\beta}}_{L1} = \arg \min_{\mathbf{x} \boldsymbol{\beta} = \mathbf{y}} \|\boldsymbol{\beta}\|_1 \quad \text{“Rich Regime”}$$

$$\alpha \leq \min \left\{ (2(1+\epsilon) \|\boldsymbol{\beta}_{L1}^*\|_1)^{-\frac{2+\epsilon}{2\epsilon}}, \exp\left(-\frac{d}{\epsilon \|\boldsymbol{\beta}_{L1}^*\|_1}\right) \right\} \implies \|\hat{\boldsymbol{\beta}}_\alpha\|_1 \leq (1+\epsilon) \|\boldsymbol{\beta}_{L1}^*\|_1$$

Sparse Learning

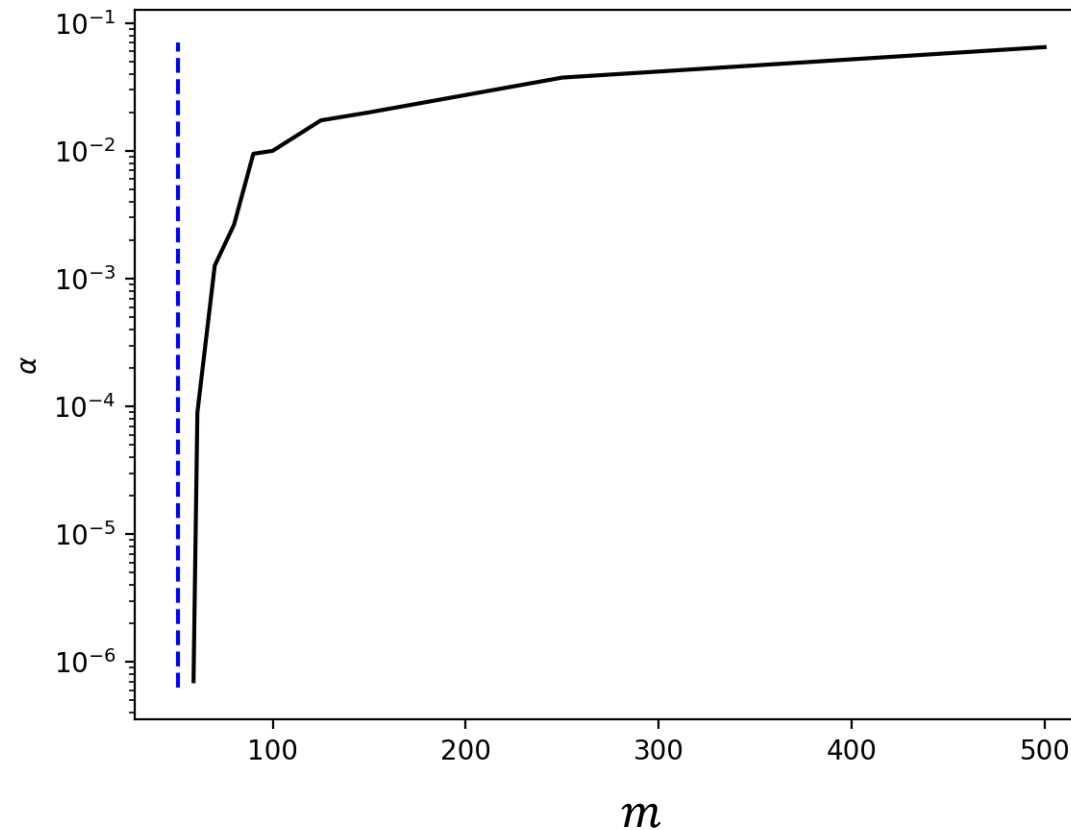
$$y_i = \langle \beta^*, x_i \rangle + N(0, 0.01)$$
$$d = 1000, \quad \|\beta^*\|_0 = 5, \quad m = 100$$



Sparse Learning

$$y_i = \langle \beta^*, x_i \rangle + N(0, 0.01)$$
$$d = 1000, \quad \|\beta^*\|_0 = k$$

How small does α need to be to get $L(\beta_\alpha(\infty)) < 0.025$



$$\beta = F(w) = w_+^2 - w_-^2 \quad \dot{w} = \nabla_w L(\beta) \quad w(0) = \alpha \mathbf{1}$$

$$\begin{aligned} \dot{\beta} &= \frac{d\beta}{dw} \dot{w} = -\nabla F(w(t))^\top \left(\nabla F(w(t)) \nabla L(\beta(t)) \right) = -\rho^{-1} \nabla L(\beta) \\ \rho &= \left(\nabla F(w(t))^\top \nabla F(w(t)) \right)^{-1} = \text{diag}(w_+^2 + w_-^2)^{-1} \\ \nabla F^\top &= [\text{diag}(w_+) \text{diag}(w_-)] \end{aligned}$$

Problem 1: $w(t)$ as a function of $\beta(t)$

$$\text{Claim: } \frac{d}{dt} (w_+(t)w_-(t)) = -2X^\top r(t) (w_+(t) \circ w_-(t) - w_-(t) \circ w_+(t)) = 0$$

$$\rightarrow w_+(t)w_-(t) = \alpha^2, \text{ combined with } \beta = w_+^2 - w_-^2 \rightarrow w_\pm^2(t) = \frac{\sqrt{\beta^2 + 4\alpha^4} \pm \beta}{2}$$

$$\rightarrow \rho^{-1} = \text{diag}(w_+^2 + w_-^2) = \text{diag}(\sqrt{\beta^2 + 4\alpha^4})$$

$$\text{Induced dynamics: } \dot{\beta}_\alpha = -\sqrt{\beta_\alpha^2 + 4\alpha^4} \odot \nabla L_s(\beta_\alpha)$$

Problem 2: Is $\rho = \text{diag}(\beta^2 + 4\alpha^4)^{-\frac{1}{2}}$ a Hessian map? Solve $\rho = \nabla^2 \Psi$

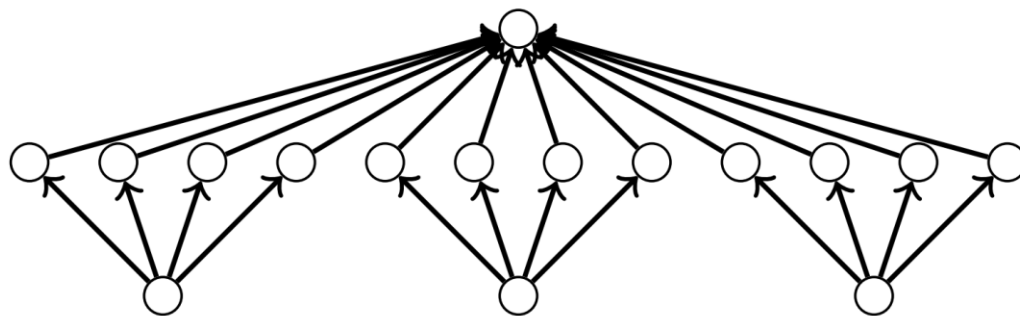
$$\Psi = \sum_i \int \int (\beta_i^2 + 4\alpha^4)^{-\frac{1}{2}} d\beta_i d\beta_i = \alpha^2 \sum_i \left(\text{const} - \sqrt{4 + \left(\frac{\beta}{\alpha^2}\right)^2} + \frac{\beta}{\alpha^2} \sinh^{-1} \left(\frac{\beta}{2\alpha^2} \right) \right)$$

Width and Initialization Scale

$$f((U, V), x) = \sum_{i=1..d, j=1..k} u_{i,j} v_{i,j} x[i] = \langle UV^\top, \text{diag}(x) \rangle$$

$$U, V \in \mathbb{R}^{d \times k} \quad \beta_{U,V} = \text{diag}(UV^\top)$$

- Initialization: $u_{i,j}, v_{i,j} \sim \text{iid } N\left(0, \sigma^2 = \frac{\alpha^2}{\sqrt{k}}\right)$ s.t. $\text{Var}[\beta(0)[i]] = \alpha^2$



Width and Initialization Scale

$$f((U, V), x) = \sum_{i=1..d, j=1..k} u_{i,j} v_{i,j} x[i] = \langle UV^\top, \text{diag}(x) \rangle$$

$$U, V \in \mathbb{R}^{d \times k} \quad \beta_{U,V} = \text{diag}(UV^\top)$$

- Initialization: $u_{i,j}, v_{i,j} \sim \text{iid } N\left(0, \sigma^2 = \frac{\alpha^2}{\sqrt{k}}\right)$ s.t. $\text{Var}[\beta(0)[i]] = \alpha^2$
- Symmetrized problem: $\tilde{f}(W, x) = \langle WW^\top, \tilde{X} \rangle$

$$W = \begin{bmatrix} U \\ V \end{bmatrix} \text{ so } WW^\top = \begin{bmatrix} UU^\top & UV^\top \\ VU^\top & VV^\top \end{bmatrix}$$

$$\tilde{X} := \frac{1}{2} \begin{bmatrix} 0 & \text{diag}(x) \\ \text{diag}(x) & 0 \end{bmatrix}$$

- Relevant scale: $WW^\top \approx \sqrt{k} \alpha^2 I$
 - $\beta(\infty) \xrightarrow{k \rightarrow \infty, \alpha \rightarrow 0} \arg \min_{x \beta=y} Q_\mu(\beta)$ $\mu = \lim \alpha^4 \sqrt{k} = \lim \sigma \sqrt{k}$
 - $\alpha = o(1/\sqrt[4]{k})$, i.e. $\sigma = o(1/\sqrt{k}) \rightarrow \ell_1$
 - $\alpha = \omega(1/\sqrt[4]{k})$, i.e. $\sigma = \omega(1/\sqrt{k}) \rightarrow \ell_2$
 - $\sqrt{k} \alpha^2 \rightarrow 0$ leads to the kernel regime, even if $|\beta(0)| \approx \alpha^2 \rightarrow 0$

Width and Initialization Scale

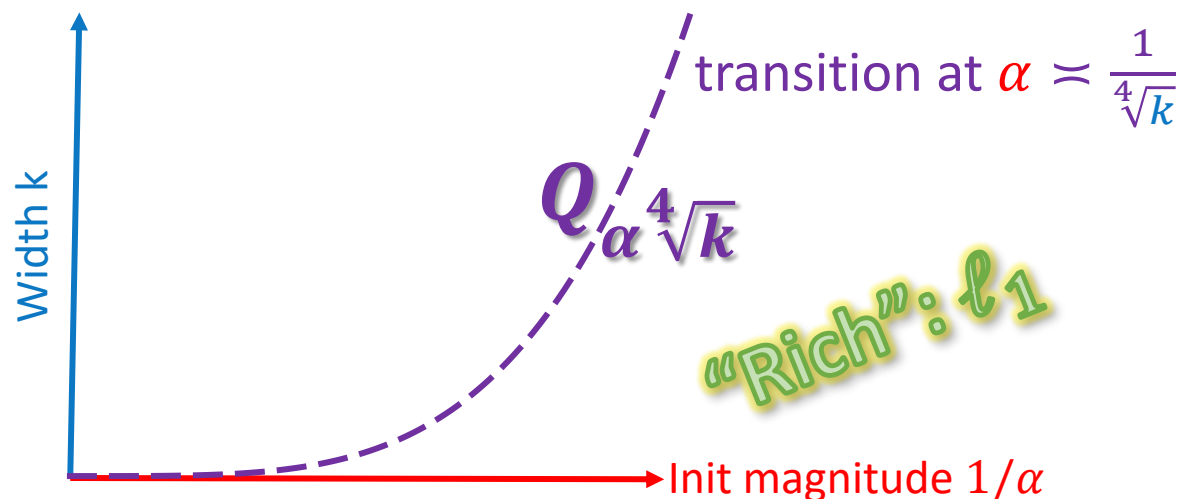
$$f((U, V), x) = \sum_{i=1..d, j=1..k} u_{i,j} v_{i,j} x[i] = \langle UV^\top, \text{diag}(x) \rangle$$

$$U, V \in \mathbb{R}^{d \times k} \quad \beta_{U,V} = \text{diag}(UV^\top)$$

- Initialization: $u_{i,j}, v_{i,j} \sim \text{iid } N\left(0, \sigma^2 = \frac{\alpha^2}{\sqrt{k}}\right)$ s.t. $\text{Var}[\beta(0)[i]] = \alpha^2$
- Symmetrized problem: $\tilde{f}(W, x) = \langle WW^\top, \tilde{X} \rangle$

$$W = \begin{bmatrix} U \\ V \end{bmatrix} \text{ so } WW^\top = \begin{bmatrix} UU^\top & UV^\top \\ VU^\top & VV^\top \end{bmatrix}$$

$$\tilde{X} := \frac{1}{2} \begin{bmatrix} 0 & \text{diag}(x) \\ \text{diag}(x) & 0 \end{bmatrix}$$

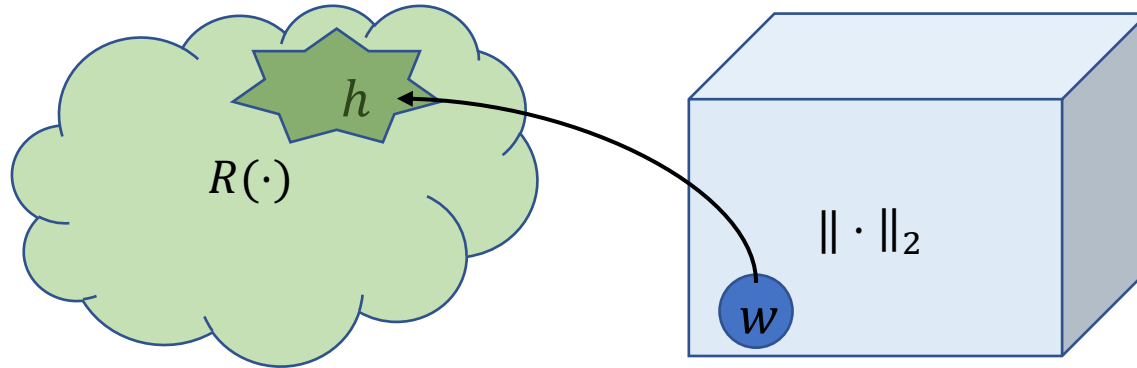


Is implicit bias of GD just ℓ_2 in param space + mapping to func space?

Is initializing to $w(0) = \alpha \mathbf{1}$ the same as regularizing distance to $\alpha \mathbf{1}$?

$$\beta_{\alpha}^R = F \left(\arg \min_{L_S(w)=0} \|w - \alpha \mathbf{1}\|_2^2 \right) = \arg \min_{X\beta=y} R_{\alpha}(\beta)$$

Where $R_{\alpha}(\beta) = \min_{F(w)=\beta} \|w - \alpha \mathbf{1}\|_2^2$



Is implicit bias of GD just ℓ_2 in param space + mapping to func space?

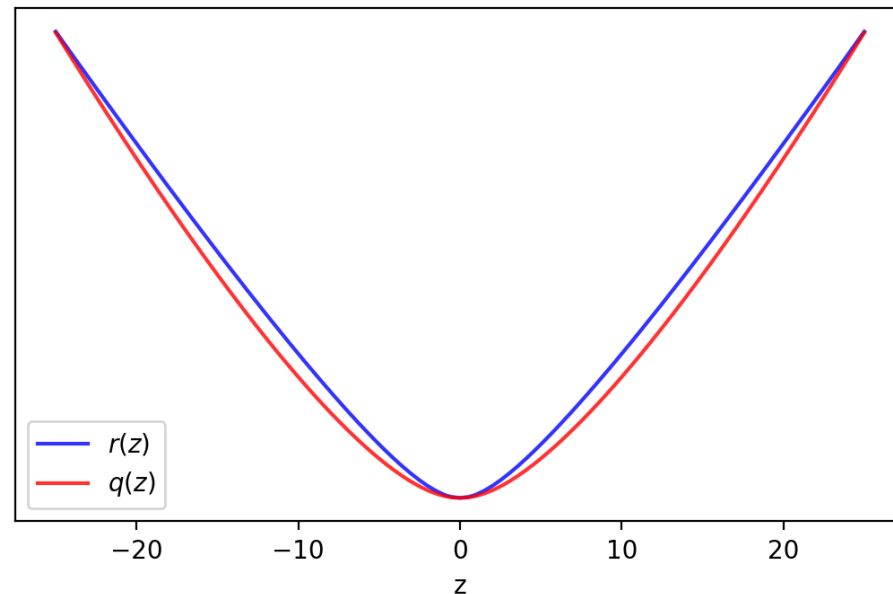
Is initializing to $w(0) = \alpha \mathbf{1}$ the same as regularizing distance to $\alpha \mathbf{1}$?

$$\beta_{\alpha}^R = F \left(\arg \min_{L_S(w)=0} \|w - \alpha \mathbf{1}\|_2^2 \right) = \arg \min_{X\beta=y} R_{\alpha}(\beta)$$

Where $R_{\alpha}(\beta) = \min_{F(w)=\beta} \|w - \alpha \mathbf{1}\|_2^2$

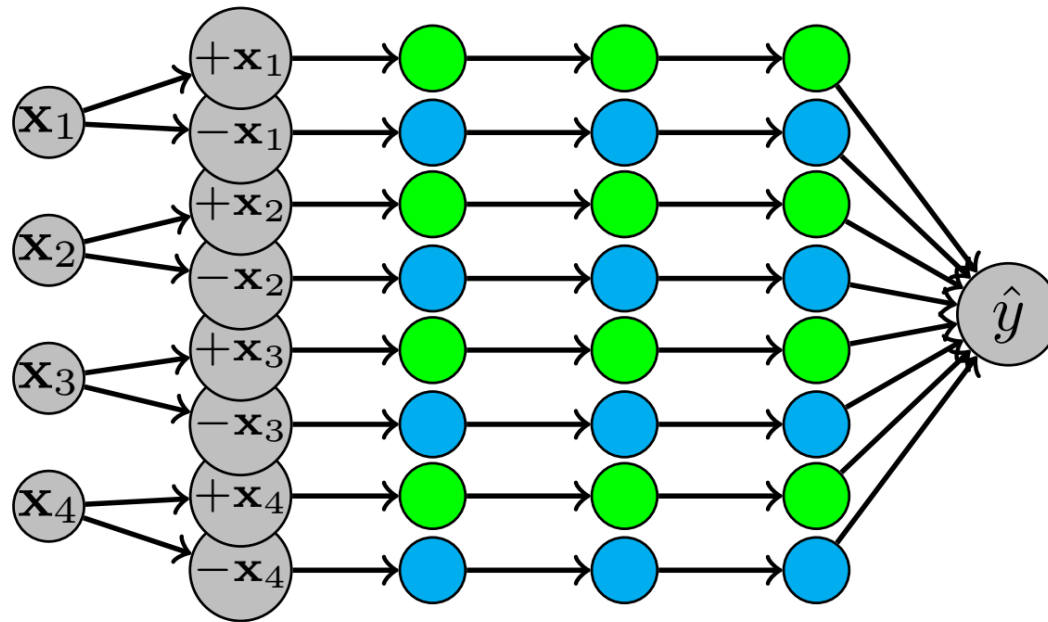
$R_{\alpha}(\beta) = \sum_j r \left(\frac{\beta[j]}{\alpha^2} \right)$ where $r(b)$ is solution of quartic equation:

$$r^4 - 6r^3 + (12 - 2b^2)r^2 - (8 + 10b^2)r + b^2 + b^4 = 0$$



Deep Diagonal Linear Net

$$\beta(t) = w_+(t)^D - w_-(t)^D$$



Deep Diagonal Linear Net

$$\beta(t) = w_+(t)^D - w_-(t)^D$$

$$r(t) = X\beta(t) - y$$

$$\beta(t) = \alpha^D \left(\left(1 + \alpha^{D-2} D(D-2) X^\top \int_0^t r(\tau) d\tau \right)^{\frac{-1}{D-2}} - \left(1 - \alpha^{D-2} D(D-2) X^\top \int_0^t r(\tau) d\tau \right)^{\frac{-1}{D-2}} \right)$$

KKT for $\min Q(\beta)$ s. t. $X\beta = y$:

$$\nabla Q(\beta^*) = X^\top \nu$$

$$X\beta^* = y$$

$$L = \|X\beta - y\|_2^2$$

$$\frac{dw}{dt} = -\nabla L$$

$$\frac{d\beta}{dt} = \frac{d\beta}{dw} \cdot \frac{dw}{dt}$$

$$s = \int_0^\infty r(\tau) d\tau \in \mathbb{R}^m$$

$$\beta(\infty) = \alpha^D h_D(X^\top s)$$

$$X\beta(\infty) = y$$

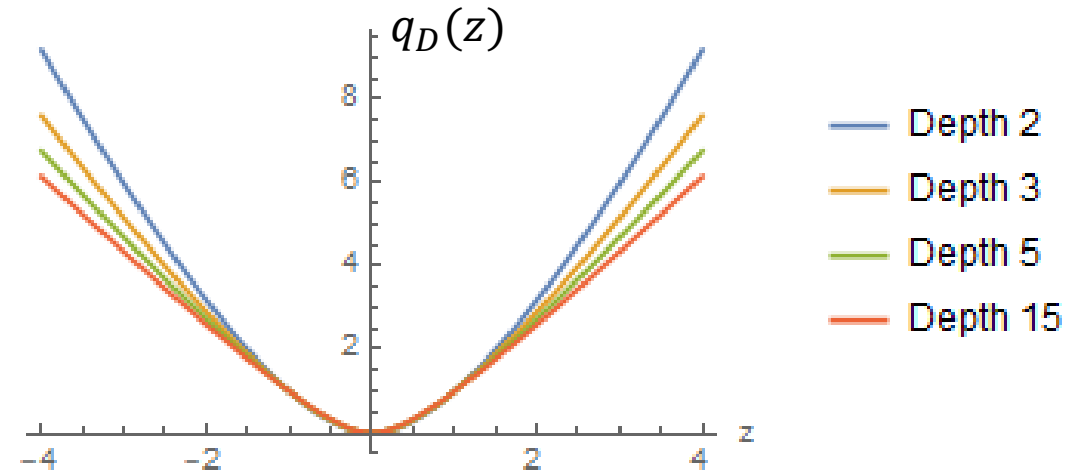
$$h_D(z) = \alpha^D \left((1 + \alpha^{D-2} D(D-2) z)^{\frac{-1}{D-2}} - (1 - \alpha^{D-2} D(D-2) z)^{\frac{-1}{D-2}} \right)$$

$$q_D = \int h_D^{-1}$$

$$Q_D(\beta) = \sum_i q_D \left(\frac{\beta[i]}{\alpha^D} \right)$$

Deep Diagonal Linear Net

$$\beta(t) = w_+(t)^D - w_-(t)^D \quad \beta(\infty) = \arg \min Q_D \left(\beta / \alpha^D \right) \quad s.t. \quad X\beta = y$$



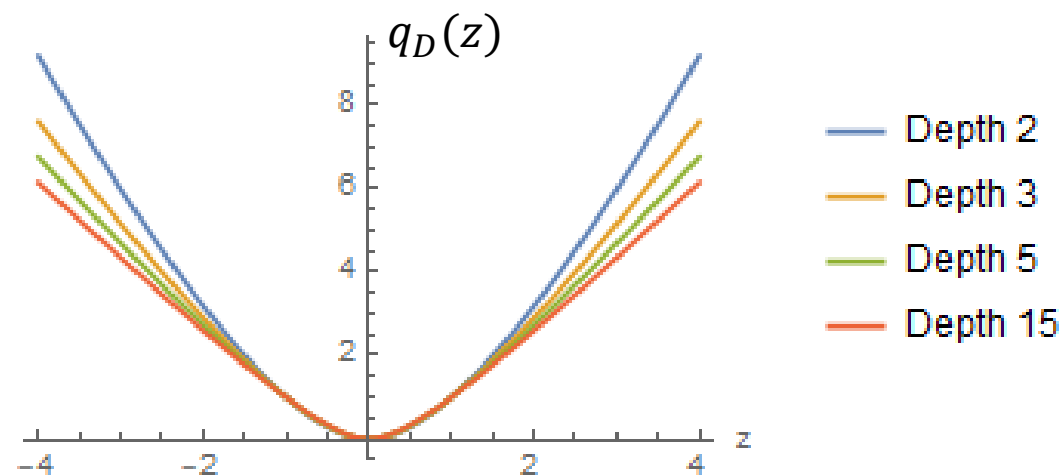
$$h_D(z) = \alpha^D \left((1 + \alpha^{D-2} D(D-2)z)^{\frac{-1}{D-2}} - (1 - \alpha^{D-2} D(D-2)z)^{\frac{-1}{D-2}} \right)$$

$$q_D = \int h_D^{-1}$$

$$Q_{D,\alpha}(\beta) = \sum_i q_D \left(\frac{\beta[i]}{\alpha^D} \right)$$

Deep Diagonal Linear Net

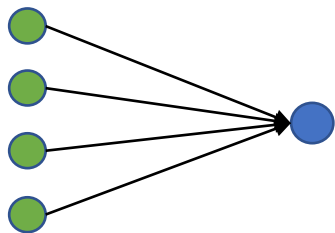
$$\beta(t) = w_+(t)^D - w_-(t)^D \quad \beta(\infty) = \arg \min Q_D \left(\beta / \alpha^D \right) \text{ s.t. } X\beta = y$$



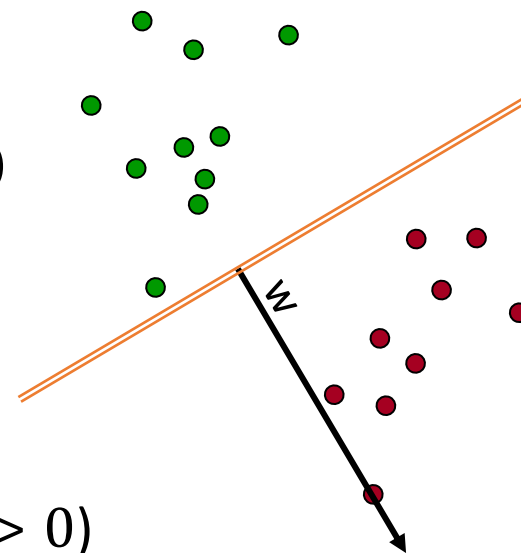
For all depth $D \geq 2$, $\beta(\infty) \xrightarrow{\alpha \rightarrow 0} \arg \min_{X\beta=y} \|\beta\|_1$

- Contrast with explicit reg: For $R_\alpha(\beta) = \min_{\beta=w_+^D - w_-^D} \|w - \alpha \mathbf{1}\|_2^2$, $R_\alpha(\beta) \xrightarrow{\alpha \rightarrow 0} \|\beta\|_{2/D}$
also observed by [Arora Cohen Hu Luo 2019]
- Also with **logistic loss**, $\beta(\infty) \rightarrow \propto \text{SOSP of } \|\beta\|_{2/D}$ [Gunasekar Lee Soudry Srebro 2018]
[Lyu Li 2019]

Implicit Bias in Logistic Regression

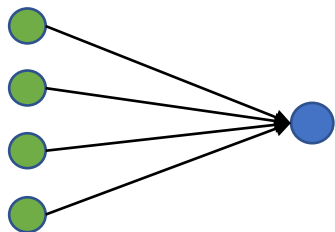


$$\arg \min_{w \in \mathbb{R}^n} \mathcal{L}(w) = \sum_{i=1}^m \ell(y_i \langle w, x_i \rangle)$$
$$\ell(z) = \log(1 + e^{-z})$$



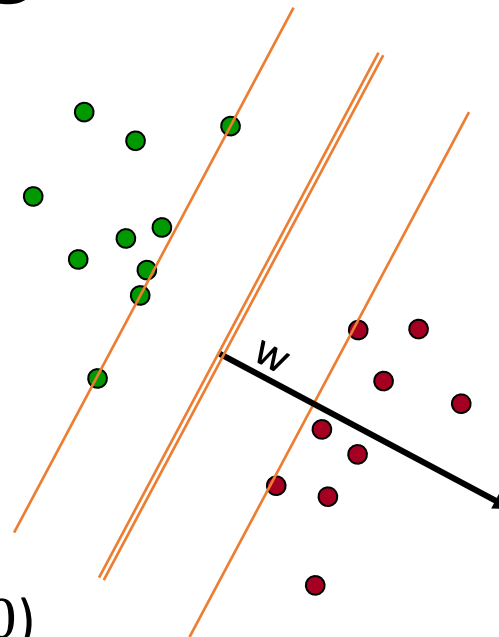
- Data $\{(x_i, y_i)\}_{i=1}^m$ linearly separable ($\exists_w \forall_i y_i \langle w, x_i \rangle > 0$)
- Where does gradient descent converge?
$$w(t) = w(t) - \eta \nabla \mathcal{L}(w(t))$$
 - $\inf \mathcal{L}(w) = 0$, but minima unattainable
 - GD diverges to infinity: $w(t) \rightarrow \infty$, $\mathcal{L}(w(t)) \rightarrow 0$
- **In what direction?** What does $\frac{w(t)}{\|w(t)\|}$ converge to?

Implicit Bias in Logistic Regression



$$\arg \min_{w \in \mathbb{R}^n} \mathcal{L}(w) = \sum_{i=1}^m \ell(y_i \langle w, x_i \rangle)$$

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



- Data $\{(x_i, y_i)\}_{i=1}^m$ linearly separable ($\exists_w \forall_i y_i \langle w, x_i \rangle > 0$)

- Where does gradient descent converge?

$$w(t) = w(t) - \eta \nabla \mathcal{L}(w(t))$$

- $\inf \mathcal{L}(w) = 0$, but minima unattainable
- GD diverges to infinity: $w(t) \rightarrow \infty$, $\mathcal{L}(w(t)) \rightarrow 0$
- **In what direction?** What does $\frac{w(t)}{\|w(t)\|}$ converge to?
- **Theorem:** $\frac{w(t)}{\|w(t)\|_2} \rightarrow \frac{\hat{w}}{\|\hat{w}\|_2}$ $\hat{w} = \arg \min \|w\|_2 \text{ s.t. } \forall_i y_i \langle w, x_i \rangle \geq 1$

Implicit Bias in Logistic Regression

$$\arg \min_{w \in \mathbb{R}^n} \mathcal{L}(w) = \sum_{i=1}^m \ell(y_i \langle w, x_i \rangle)$$
$$\ell(z) = \log(1 + e^{-z})$$

Theorem: $w(t) = \hat{w} \log t + \rho(t)$, with $\rho(t)$ bounded*

$$\hat{w} = \arg \min \|w\|_2 \text{ s.t. } \forall_i y_i \langle w, x_i \rangle \geq 1$$

- Holds for any initial point $w(0)$ and stepsize $\eta \leq 2$
- Holds for any monotonically decreasing strictly positive smooth loss s.t. $-\ell'(z)$ has a tight exponential tail. Asymptotically, all behave as:

$$\ell(z) = e^{-z}$$

*For data in general position. With degenerate data, $\rho(t) = O(\log \log t)$

Proof sketch: ($y_i = 1$ w.l.o.g.)

Write $w(t) = g(t)w_\infty + \rho(t)$ with $g(t) \rightarrow \infty$ and $\rho(t) = o(g(t))$.

Since we converge to zero error, $\forall_i \langle w_\infty, x_i \rangle > 0$

Since the loss derivative has an exponential tail:

$$-\nabla \mathcal{L}(w) \approx \sum_i e^{-\langle w(t), x_i \rangle} x_i^\top = \sum_i e^{-g(t)\langle w_\infty, x_i \rangle - \langle \rho(t), x_i \rangle} x_i^\top$$

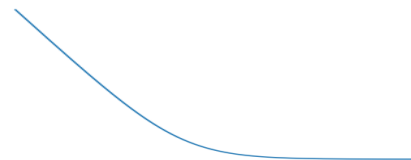
As $g(t) \rightarrow \infty$, only points with minimal $\langle w_\infty, x_i \rangle$ (points on the margin, “support vectors”) will dominate gradient

→ $\nabla \mathcal{L}(w)$ spanned by support vectors

→ $w(t)$ spanned by support vectors

Define $\hat{w} = \frac{w_\infty}{\min_i \langle w_\infty, x_i \rangle}$. We have:

$$\hat{w} = \sum \alpha_i w_i \quad \forall_i (\alpha_i \geq 0 \text{ and } \langle \hat{w}, x_i \rangle = 1) \text{ OR } (\alpha_i = 0 \text{ and } \langle \hat{w}, x_i \rangle > 1)$$

$$\ell_{\text{logistic}}(h(w), y) = \log(1 + e^{-yh(w)}) \approx e^{-yh(w)} = \ell_{\text{exp}}(h(w), y)$$


Consider gradient descent w.r.t. logistic loss $L_S(w) = \sum_i \ell(f(w, x_i); y_i)$
 (or other exp-tail loss) on a D-homogenous model $f(w, x)$:

Theorem [Nacson Gunasekar Lee S Soudry 2019][Lyu Li 2019]:

If $L_S(w) \rightarrow 0$, and small enough stepsize (ensuring convergence in direction):

$$w_\infty \propto \text{first order stationary point of} \quad (*)$$

$$\arg \min \|w\|_2 \text{ s.t. } \forall_i y_i f(w, x_i) \geq 1$$

Suggests implicit bias defined by $R_F(h) = \arg \min_{F(w)=h} \|w\|_2$ and

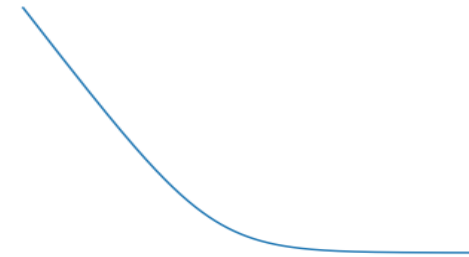
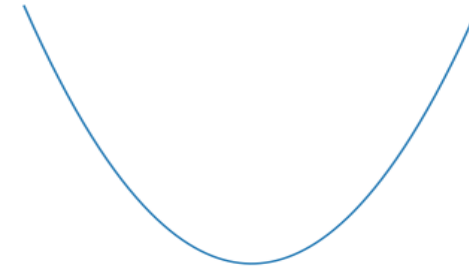
$$h_\infty = F(w_\infty) \propto \text{first order stationary point of} \quad (**)$$

$$\arg \min R_F(h) \text{ s.t. } y_i f(x_i) \geq 1$$

But need to be careful: f.o.s.p of $(*)$ does **not** imply f.o.s.p of $(**)$

Different Asymptotics

- For least squares (or any other loss with attainable minimum):
 - w_∞ depends on initial point w_0 and stepsize η
 - To get clean characterization, need to take $\eta \rightarrow 0$
 - If 0 is a saddle point, need to take $w_0 \rightarrow 0$
- For monotone decreasing loss (eg logistic)
 - w_∞ does NOT depend on initial w_0 and stepsize η
 - Don't need $\eta \rightarrow 0$ and $w_0 \rightarrow 0$
 - What happens at the beginning doesn't effect w_∞



Squared Loss vs Logistic/Exp Loss

$$\ell_{\text{logistic}}(h(w), y) = \log(1 + e^{-yh(w)}) \approx e^{-yh(w)} = \ell_{\text{exp}}(h(w), y)$$

When $\ell \rightarrow 0$, ie $yh(w) \rightarrow \infty$

- For squared loss, under some conditions [Chizat and Bach 18]:

$$\lim_{\alpha \rightarrow \infty} \sup_t \left\| w_\alpha \left(\frac{1}{\alpha^{D-1}} t \right) - w_K(t) \right\| = 0$$

$$\Rightarrow h_\alpha(\infty) \rightarrow \hat{h}_K = \arg \min \|h\|_K \text{ s.t. } h(x_i) = y_i$$

- For logistic:

$$\forall_t \lim_{\alpha \rightarrow \infty} \sup_{\tau < t} \left\| w_\alpha \left(\frac{1}{\alpha^{D-1}} \tau \right) - w_K(\tau) \right\| = 0$$

Contrast with [Nacson Gunasekar Lee S Soudry 2019][Lyu Li 2019]:

$$\forall_\alpha \lim_{t \rightarrow \infty} \frac{w_\alpha(t)}{\|w_\alpha(t)\|} \propto \text{f.o.s.p of } \arg \min \|w\|_2 \text{ s.t. } y_i h(x_i) \geq 1$$

For our model $\beta = w_+^2 - w_-^2$ with logistic loss:

$\beta(\epsilon) = \beta(t) \text{ s.t. } L_S(\beta) = \epsilon$
 Uniquely defined since $L_S(\beta(t))$ monotonically decreases from 1 to 0.

$$\forall \alpha \lim_{\epsilon \rightarrow 0} \frac{\beta(\epsilon)}{\|\beta(\epsilon)\|} \propto \arg \min \|\beta\|_1 \text{ s.t. } y_i x_i^\top \beta \geq 1$$

Rich



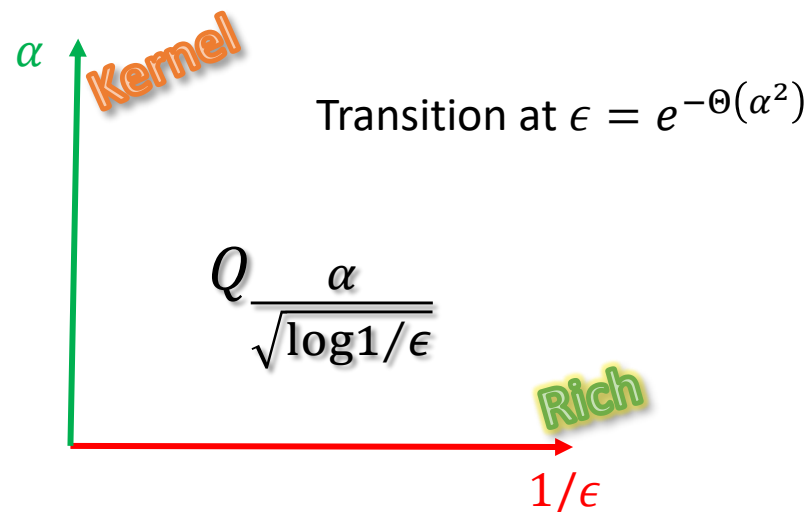
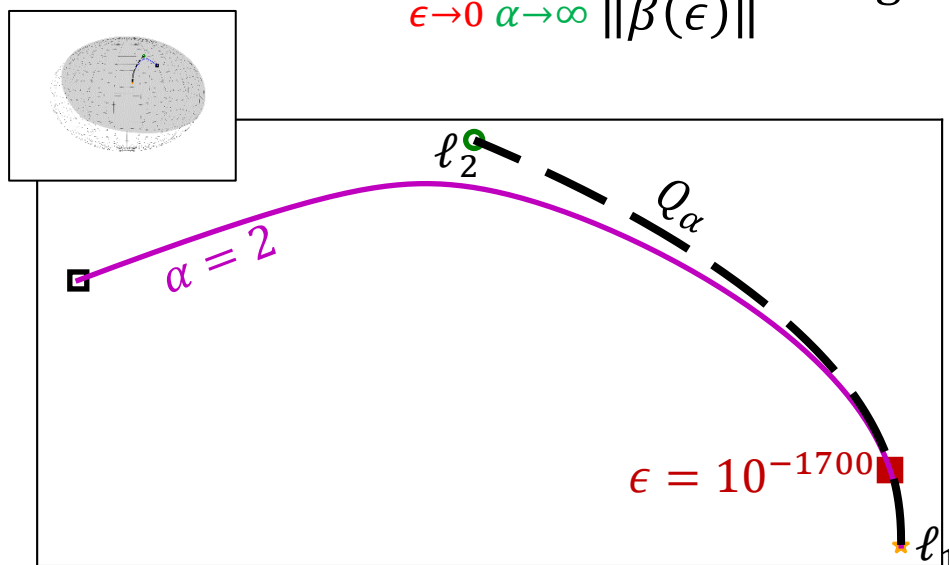
$$\lim_{\alpha \rightarrow \infty} \lim_{\epsilon \rightarrow 0} \frac{\beta(\epsilon)}{\|\beta(\epsilon)\|} \propto \arg \min \|\beta\|_1 \text{ s.t. } y_i x_i^\top \beta \geq 1$$

Rich

Contrast with:

$$\lim_{\epsilon \rightarrow 0} \lim_{\alpha \rightarrow \infty} \frac{\beta(\epsilon)}{\|\beta(\epsilon)\|} \propto \arg \min \|\beta\|_2 \text{ s.t. } y_i x_i^\top \beta \geq 1$$

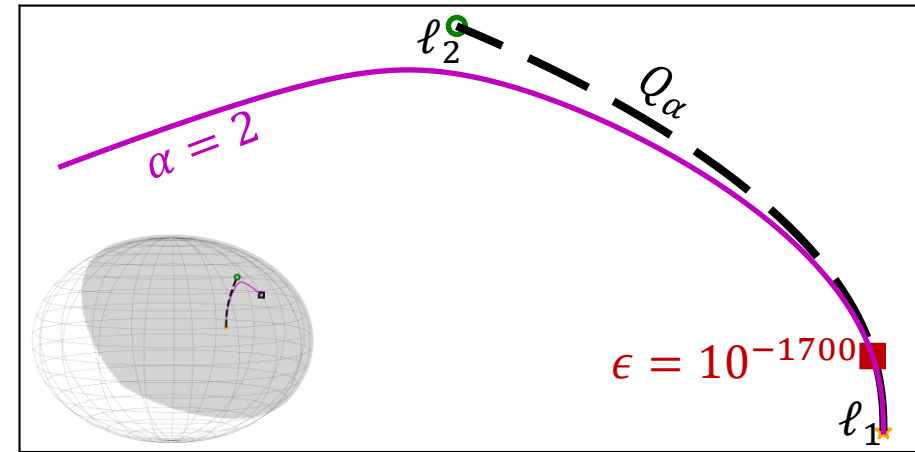
Kernel



Logistic Loss vs Squared Loss

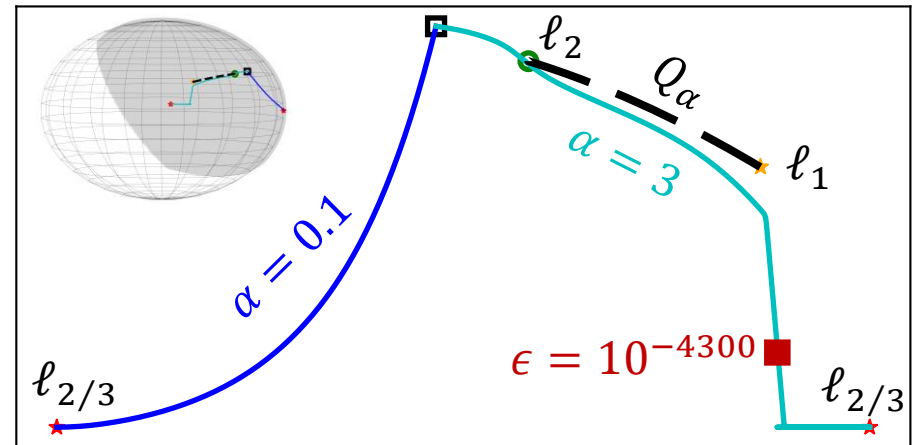
Depth two:

- Square loss: $\beta(\infty) \propto \arg \min_{X\beta=y} Q_\alpha(\beta)$
- Logistic loss: $\forall_\alpha \beta(\infty) \propto \arg \min_{X\beta=y} \|\beta\|_1$



Deeper Diagonal Nets:

- Squared loss, $\beta(\infty) \xrightarrow{\alpha \rightarrow 0} \arg \min_{X\beta=y} \|\beta\|_1$
- Logistic loss, $\beta(\infty) \propto \text{SOSP of } \|\beta\|_{2/D}$



Depth=3

[Moroshko Gunasekar Woodworth Lee S Soudry 2020 “Implicit Bias in Deep Linear Classification: Initialization Scale vs Training Accuracy”]

Other Control Choices

- Early Stopping (and not so early stopping)
- Shape/relative scale [Azulay Moroshko Nacson Woodworth S Globerson Soudry 2021]
- Stepsize [Nacson Ravichandran S Soudry 2022]
- **Stochasticity**
 - Batchsize [Pesme Pillaud-Vivien Flammarion 2021]
 - Label noise [HaoChen, Wei, Lee, Ma 2020][Blanc, Gupta, Valiant, Valiant 2020]
- ...

$$f(\mathbf{w}, x) = \sum_j (\mathbf{w}_+[j]^2 - \mathbf{w}_-[j]^2)x[j] = \langle \beta(\mathbf{w}), x \rangle$$

$$\text{with } \beta(\mathbf{w}) = \mathbf{w}_+^2 - \mathbf{w}_-^2$$

Simplest architecture displaying complex implicit bias phenomena

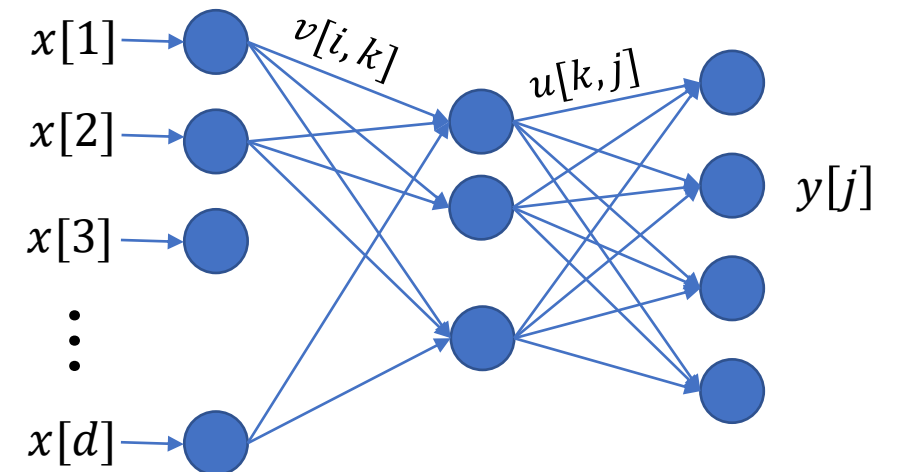
$$\min_{\beta \in \mathbb{R}^{d_1 \times d_2}} \hat{L}(\beta) = \|\mathcal{X}(\beta) - y\|_2^2 \quad \mathcal{X}(\beta)_i = \langle X_i, \beta \rangle \quad X_1, \dots, X_m \in \mathbb{R}^{d_1 \times d_2}, y \in \mathbb{R}^m$$

$$\beta = F(U, V) = UV^T, \quad U, V \in \mathbb{R}^{n \times n}$$

GD on U, V : $\dot{U}(t) = -\nabla_U \hat{L}(UV^T), \dot{V}(t) = -\nabla_V \hat{L}(UV^T)$

$$\Rightarrow \quad \dot{\beta} = -(\nabla F^T \nabla F) \nabla \hat{L}(\beta) = -(UU^T \nabla \hat{L}(\beta) + \nabla \hat{L}(\beta) VV^T)$$

- Matrix completion (X_i is indicator matrix)
- Matrix reconstruction from linear measurements
- Multi-task learning ($X_i = e_{\text{task of example } i} \cdot \phi(\text{example } i)^T$)



$$\min_{\beta \in \mathbb{R}^{d_1 \times d_2}} \hat{L}(\beta) = \|\mathcal{X}(\beta) - y\|_2^2 \quad \mathcal{X}(\beta)_i = \langle X_i, \beta \rangle \quad X_1, \dots, X_m \in \mathbb{R}^{d_1 \times d_2}, y \in \mathbb{R}^m$$

$$\beta = F(U, V) = UV^T, \quad U, V \in \mathbb{R}^{n \times n}$$

GD on U, V : $\dot{U}(t) = -\nabla_U \hat{L}(UV^T), \dot{V}(t) = -\nabla_V \hat{L}(UV^T)$

$$\Rightarrow \quad \dot{\beta} = -(\nabla F^T \nabla F) \nabla \hat{L}(\beta) = -(UU^T \nabla \hat{L}(\beta) + \nabla \hat{L}(\beta) VV^T)$$

$$W = \begin{bmatrix} U \\ V \end{bmatrix}, \quad \tilde{\beta} = WW^T = \begin{bmatrix} UU^T & UV^T \\ VU^T & VV^T \end{bmatrix} = \begin{bmatrix} UU^T & \beta \\ \beta^T & VV^T \end{bmatrix}$$

$$\min_{\tilde{\beta} \succcurlyeq 0} \hat{L}(\tilde{\beta}) = \|\tilde{\mathcal{X}}(\tilde{\beta}) - y\|_2^2 \quad \tilde{\mathcal{X}}(\tilde{\beta})_i = \langle \tilde{X}_i, \tilde{\beta} \rangle \quad \tilde{X}_i = \begin{bmatrix} 0 & X_i \\ X_i^T & 0 \end{bmatrix} \in \mathbb{S}_d, y \in \mathbb{R}^m$$

$$\tilde{\beta} = \tilde{F}(W) = WW^T, W \in \mathbb{R}^{d \times d}$$

$$\dot{\tilde{\beta}} = -(\nabla \tilde{F}^T \nabla \tilde{F}) \nabla \hat{L}(\tilde{\beta}) = -(WW^T \nabla \hat{L}(\tilde{\beta}) + \nabla \hat{L}(\tilde{\beta}) WW^T) = -(\tilde{\beta} \nabla \hat{L}(\tilde{\beta}) + \nabla \hat{L}(\tilde{\beta}) \tilde{\beta})$$

$$\min_{\beta \succcurlyeq 0} \hat{L}(\beta) = \|\mathcal{X}(\beta) - y\|_2^2 \quad \mathcal{X}(\beta)_i = \langle X_i, \beta \rangle \quad X_i \in \mathbb{S}_d, y \in \mathbb{R}^m$$

$$\dot{\beta} = -(\beta \nabla \hat{L}(\beta) + \nabla \hat{L}(\beta) \beta) = (-\beta \mathcal{X}^*(r(t)) - \mathcal{X}^*(r(t)) \beta)$$

$r(t) = \mathcal{X}(\beta) - y$

If X_i, β_0 commute:

$$\beta(t) = e^{\mathcal{X}^*(s_t)} \beta_0 e^{\mathcal{X}^*(s_t)} \quad s_t = -\int r_t dt \in \mathbb{R}^m$$

$$\in \{e^{\mathcal{A}^*(s)} \beta_0 e^{\mathcal{A}^*(s)} \mid s \in \mathbb{R}^m\}$$

- Independent of “steering” r_t
- Can use other loss, weights, or sample X_i ; But finite steps, as well as (infinitesimal) momentum, will fall off \mathcal{M} !
- Restricting to joint diagonalization $\beta = U \tilde{\beta} U^\top$, $\rho(\beta) = (A \mapsto \beta A + A \beta)^{-1} = \frac{1}{2} \beta^{-1}$ is a Hessian map:

$$\Psi(\beta) = \sum_i \tilde{\beta}[i] \log \frac{\tilde{\beta}[i]}{e} \rightarrow D_\Psi(\beta \parallel \alpha I) = \sum_i \tilde{\beta}[i] \left(\log \frac{1}{e\alpha} + \tilde{\beta}[i] \right)$$

If X_i don't commute, solution given by “time ordered exponential”:

$$\beta(t) = \left(\lim_{\epsilon \rightarrow 0} \prod_{\tau=t/\epsilon}^0 e^{-\epsilon \mathcal{X}^*(r_\tau)} \right) \beta_0 \left(\lim_{\epsilon \rightarrow 0} \prod_{\tau=0}^{t/\epsilon} e^{-\epsilon \mathcal{X}^*(r_\tau)} \right)$$

- With arbitrary (crazy) steering, can move in any direction and get to any psd matrix (even with $m = 2$ random measurement matrices)
- $\rho(\beta) = (A \mapsto \beta A + A \beta)^{-1}$ is not a Hessian map

The “complexity measure” approach

Identify $c(h)$ s.t.

- Optimization algorithm biases towards low $c(h)$
- $\mathcal{H}_{c(reality)} = \{h | c(h) \leq c(reality)\}$ has low capacity
- Reality is well explained by low $c(h)$

Can optimization bias can be described as **$\arg \min c(h) \text{ s.t. } L_S(h) = 0$** ??

- Not always [Dauber Feder Koren Livni 2020]
- Approximately? Enough to explain generalization??

Ultimate Question: What is the true Inductive Bias? What makes reality *efficiently* learnable by fitting a (huge) neural net with a specific algorithm?