## **Double Descent / Implicit Regularization** + PAC-Bayes

CPSC 532S: Modern Statistical Learning Theory 30 March 2022 cs.ubc.ca/~dsuth/532S/22/



### Admin

- A3 late-late deadline is tonight
- A4 out, due next Friday late deadline Sunday, late-late deadline Tuesday
- Last regular class is Monday
  - Overview of online learning, privacy (~equivalent), connection to stability
- Project presentations in class next Wednesday
  - Details were posted on Piazza a bit ago
  - Will post schedule tonight
- Project reports due next Friday usual late policy (-5 Saturday, -10 Sunday)
- Final will be take-home, available over most of the finals period Exact dates TBA – contact me if this really matters to you for whatever reason







### Nakkiran et al. blog post's companion notebook











### Nakkiran et al. blog post's companion notebook





![](_page_6_Figure_0.jpeg)

 $L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d,$ of rank *n* 

 $L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top}(Xw - y)$ of rank *n* 

 $L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = 0$ of rank *n* 

![](_page_9_Picture_2.jpeg)

 $L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top}(Xw - y) \quad w_{0} = 0$ of rank *n* 

 $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1})$ 

![](_page_10_Picture_3.jpeg)

 $L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = 0$ of rank *n* 

 $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right) w_{k-1} + \frac{\eta}{n} X^{\mathsf{T}} y$ 

![](_page_11_Picture_3.jpeg)

 $L_{S}(w) = \frac{1}{2n} ||Xw - y||^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with}$ of rank *n* 

 $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right)$ 

$$n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = \tilde{\eta}$$

$$X = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\mathsf{T}} X)^{\ell} X^{\mathsf{T}} y$$

![](_page_12_Figure_4.jpeg)

 $L_{S}(w) = \frac{1}{2n} ||Xw - y||^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with}$ of rank *n* 

 $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right)$ 

 $X = U \Sigma V^{\mathsf{T}}$ 

$$n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = \tilde{\eta}$$

$$X = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\mathsf{T}} X)^{\ell} X^{\mathsf{T}} y$$

![](_page_13_Figure_5.jpeg)

 $L_{S}(w) = \frac{1}{2n} ||Xw - y||^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with}$ of rank *n* 

 $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right)$ 

 $X = U\Sigma V^{\mathsf{T}} = U\tilde{\Sigma}\tilde{V}^{\mathsf{T}}$ 

$$n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = \tilde{\eta}$$

$$X = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\mathsf{T}} X)^{\ell} X^{\mathsf{T}} y$$

![](_page_14_Figure_5.jpeg)

$$L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = \frac{1}{n} \|Xw - y\|^{2} \quad \text{of rank } n$$

$$w_{k} = w_{k-1} - \eta \nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right) w_{k-1} + \frac{\eta}{n} X^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\mathsf{T}} X)^{\ell} X^{\mathsf{T}} y$$

 $X = U\Sigma V^{\mathsf{T}} = U\tilde{\Sigma}\tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = \begin{bmatrix} \Sigma & 0 \end{bmatrix} \in \mathbb{R}^{n \times d}$ 

4

![](_page_15_Figure_5.jpeg)

$$L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = 0$$
  
of rank  $n$   
$$w_{k} = w_{k-1} - \eta \nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\top} X\right) w_{k-1} + \frac{\eta}{n} X^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\top} X)^{\ell} X^{\top} y$$
  
$$X = U \Sigma V^{\top} = U \tilde{\Sigma} \tilde{V}^{\top} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = [V \quad V_{2}] \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\top} \tilde{V} = I = \tilde{V}$$

![](_page_16_Figure_4.jpeg)

$$L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad \begin{array}{l} X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = 0 \\ \text{of rank } n & \tilde{\eta} \\ w_{k} = w_{k-1} - \eta \,\nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\top} X\right) w_{k-1} + \frac{\eta}{n} X^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\top} X)^{\ell} X^{\top} y \\ X = U \Sigma V^{\top} = U \tilde{\Sigma} \tilde{V}^{\top} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = \begin{bmatrix} V \quad V_{2} \end{bmatrix} \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\top} \tilde{V} = I = \tilde{V} \\ \end{array}$$

$$L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad \begin{array}{l} X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = 0 \\ \text{of rank } n & \tilde{\eta} \\ w_{k} = w_{k-1} - \eta \,\nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\top} X\right) w_{k-1} + \frac{\eta}{n} X^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\top} X)^{\ell} X^{\top} y \\ X = U \Sigma V^{\top} = U \tilde{\Sigma} \tilde{V}^{\top} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = [V \quad V_{2}] \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\top} \tilde{V} = I = \tilde{V} \\ \end{array}$$

# $w_k = \tilde{\eta} \sum_{k=1}^{k-1} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y$ $\ell = 0$

![](_page_17_Figure_4.jpeg)

$$L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad \begin{array}{l} X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = 0 \\ \text{of rank } n & \tilde{\eta} \\ w_{k} = w_{k-1} - \eta \,\nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right) w_{k-1} + \frac{\eta}{n} X^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\mathsf{T}} X)^{\ell} X^{\mathsf{T}} y \\ X = U \Sigma V^{\mathsf{T}} = U \tilde{\Sigma} \tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = [V \quad V_{2}] \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\mathsf{T}} \tilde{V} = I = \tilde{V} \\ w_{k} = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \quad \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left(I - \tilde{\eta} \left[ \sum_{\ell=0}^{2} 0 \right] \right)^{\ell} \quad \tilde{V}^{\mathsf{T}} \quad \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} \end{array}$$

$$\begin{split} L_{S}(w) &= \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = 0 \\ \text{of rank } n \quad & \tilde{\eta} \\ w_{k} &= w_{k-1} - \eta \, \nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\top} X\right) w_{k-1} + \frac{\eta}{n} X^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\top} X)^{\ell} X^{\top} y \\ X &= U \Sigma V^{\top} = U \tilde{\Sigma} \tilde{V}^{\top} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = [V \quad V_{2}] \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\top} \tilde{V} = I = \tilde{V} \\ w_{k} &= \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\top} \tilde{\Sigma} \tilde{V}^{\top})^{\ell} \quad \tilde{V} \tilde{\Sigma}^{\top} U^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left(I - \tilde{\eta} \begin{bmatrix} \Sigma^{2} & 0 \\ 0 & 0 \end{bmatrix} \right)^{\ell} \quad \tilde{V}^{\top} \quad \tilde{V} \tilde{\Sigma}^{\top} U^{\top} V \end{split}$$

![](_page_18_Figure_4.jpeg)

![](_page_18_Picture_5.jpeg)

$$L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = 0$$
  
of rank  $n$   
$$w_{k} = w_{k-1} - \eta \nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right) w_{k-1} + \frac{\eta}{n} X^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\mathsf{T}} X)^{\ell} X^{\mathsf{T}} y$$
  
$$X = U \Sigma V^{\mathsf{T}} = U \tilde{\Sigma} \tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = [V \quad V_{2}] \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\mathsf{T}} \tilde{V} = I = \tilde{V}$$
  
$$w_{k} = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \quad \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left(I - \tilde{\eta} \begin{bmatrix} \Sigma^{2} & 0 \\ 0 & 0 \end{bmatrix} \right)^{\ell} \tilde{V}^{\mathsf{T}} \quad \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}}$$

$$L_{S}(w) = \frac{1}{2n} \|Xw - y\|^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = 0$$
  
of rank  $n$   
$$w_{k} = w_{k-1} - \eta \nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right) w_{k-1} + \frac{\eta}{n} X^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\mathsf{T}} X)^{\ell} X^{\mathsf{T}} y$$
  
$$X = U \Sigma V^{\mathsf{T}} = U \tilde{\Sigma} \tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = [V \quad V_{2}] \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\mathsf{T}} \tilde{V} = I = \tilde{V}$$
  
$$w_{k} = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \quad \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left(I - \tilde{\eta} \begin{bmatrix} \Sigma^{2} & 0 \\ 0 & 0 \end{bmatrix} \right)^{\ell} \tilde{V}^{\mathsf{T}} \quad \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}}$$

# $= \tilde{\eta} \sum_{\ell=0}^{k-1} \begin{bmatrix} V & V_2 \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^2)^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} U^{\mathsf{T}} y$

![](_page_19_Figure_5.jpeg)

![](_page_19_Picture_6.jpeg)

$$\begin{split} L_{S}(w) &= \frac{1}{2n} \|Xw - y\|^{2} \quad \begin{array}{l} X \in \mathbb{R}^{n \times d} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} \\ \text{of rank } n \quad & \tilde{\eta} \\ w_{k} &= w_{k-1} - \eta \nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\top} X\right) w_{k-1} + \frac{\eta}{n} X^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\top} X)^{\ell} X^{\top} y \\ X &= U \Sigma V^{\top} = U \tilde{\Sigma} \tilde{V}^{\top} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = [V \quad V_{2}] \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\top} \tilde{V} = I = \tilde{V} \\ w_{k} &= \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\top} \tilde{\Sigma} \tilde{V}^{\top})^{\ell} \quad \tilde{V} \tilde{\Sigma}^{\top} U^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left(I - \tilde{\eta} \left[ \sum_{0}^{2^{2}} 0 \\ 0 & 0 \end{array} \right] \right)^{\ell} \quad \tilde{V}^{\top} \quad \tilde{V} \tilde{\Sigma}^{\top} U^{\top} y \\ &= \tilde{\eta} \sum_{\ell=0}^{k-1} \left[ V \quad V_{2} \right] \begin{bmatrix} (I - \tilde{\eta} \Sigma^{2})^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} U^{\top} y \quad &= \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\top} y \end{split}$$

$$\begin{split} L_{S}(w) &= \frac{1}{2n} \|Xw - y\|^{2} \quad \stackrel{X \in \mathbb{R}^{n \times d}}{\text{of rank } n} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = 0 \\ w_{k} &= w_{k-1} - \eta \nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\top} X\right) w_{k-1} + \frac{\eta}{n} X^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\top} X)^{\ell} X^{\top} y \\ X &= U \Sigma V^{\top} = U \tilde{\Sigma} \tilde{V}^{\top} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = [V \quad V_{2}] \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\top} \tilde{V} = I = \tilde{V} \\ w_{k} &= \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\top} \tilde{\Sigma} \tilde{V}^{\top})^{\ell} \quad \tilde{V} \tilde{\Sigma}^{\top} U^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left(I - \tilde{\eta} \left[ \sum_{0}^{2} & 0 \\ 0 & 0 \end{bmatrix} \right)^{\ell} \quad \tilde{V}^{\top} \quad \tilde{V} \tilde{\Sigma}^{\top} U^{\top} y \\ &= \tilde{\eta} \sum_{\ell=0}^{k-1} [V \quad V_{2}] \begin{bmatrix} (I - \tilde{\eta} \Sigma^{2})^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} U^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\top} y \\ &= \tilde{\eta} \sum_{\ell=0}^{k-1} [V \quad V_{2}] \begin{bmatrix} (I - \tilde{\eta} \Sigma^{2})^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} U^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\top} y \\ &= \tilde{\eta} \sum_{\ell=0}^{k-1} [V \quad V_{2}] \begin{bmatrix} (I - \tilde{\eta} \Sigma^{2})^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} U^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\top} y \end{bmatrix}$$

![](_page_20_Figure_4.jpeg)

![](_page_20_Picture_5.jpeg)

$$\begin{split} L_{S}(w) &= \frac{1}{2n} \|Xw - y\|^{2} \quad \stackrel{X \in \mathbb{R}^{n \times d}}{\text{of rank } n} \text{ with } n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = 0 \\ w_{k} &= w_{k-1} - \eta \, \nabla L_{S}(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right) w_{k-1} + \frac{\eta}{n} X^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\mathsf{T}} X)^{\ell} X^{\mathsf{T}} y \\ X &= U \Sigma V^{\mathsf{T}} = U \tilde{\Sigma} \tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{n \times d} \quad \tilde{V} = \begin{bmatrix} V \quad V_{2} \end{bmatrix} \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\mathsf{T}} \tilde{V} = I = \tilde{V} \\ w_{k} &= \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \quad \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left(I - \tilde{\eta} \begin{bmatrix} \Sigma^{2} & 0 \\ 0 & 0 \end{bmatrix} \right)^{\ell} \quad \tilde{V}^{\mathsf{T}} \quad \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} \\ &= \tilde{\eta} \sum_{\ell=0}^{k-1} \begin{bmatrix} V \quad V_{2} \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^{2})^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} U^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\mathsf{T}} y \\ &= \tilde{\eta} \sum_{\ell=0}^{k-1} \begin{bmatrix} V \quad V_{2} \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^{2})^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} U^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\mathsf{T}} y \\ &= \tilde{\eta} \sum_{\ell=0}^{k-1} \begin{bmatrix} V \quad V_{2} \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^{2})^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} U^{\mathsf{T}} y = \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\mathsf{T}} y \end{bmatrix}$$

 $= \tilde{\eta} V \left[ \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \Sigma^2)^{\ell} \right] \Sigma U^{\mathsf{T}} y$ 

![](_page_21_Figure_4.jpeg)

![](_page_21_Picture_5.jpeg)

 $L_{S}(w) = \frac{1}{2n} ||Xw - y||^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with}$ of rank *n*  $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right)$  $X = U\Sigma V^{\mathsf{T}} = U\tilde{\Sigma}\tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{\mathsf{T}}$  $w_k = \tilde{\eta} \sum_{k=1}^{N} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y$  $\ell = 0$  $= \tilde{\eta} \sum_{k=1}^{k-1} \begin{bmatrix} V & V_2 \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^2)^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix}$  $\ell = 0$  $= \tilde{\eta} V \left[ \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \Sigma^2)^{\ell} \right] \Sigma U^{\mathsf{T}} y$ Neumann series:  $\sum_{\ell=0}^{\infty} A^{\ell} = (I - A)^{-1}$  when  $||A||_{op} < 1$ 

![](_page_22_Figure_3.jpeg)

![](_page_22_Picture_4.jpeg)

 $L_{S}(w) = \frac{1}{2n} ||Xw - y||^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with}$ of rank *n*  $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right)$  $X = U\Sigma V^{\mathsf{T}} = U\tilde{\Sigma}\tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{\mathsf{T}}$  $w_k = \tilde{\eta} \sum_{k=1}^{N} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y$  $\ell = 0$  $= \tilde{\eta} \sum_{k=1}^{k-1} \begin{bmatrix} V & V_2 \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^2)^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix}$  $\ell = 0$  $= \tilde{\eta} V \left[ \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \Sigma^2)^{\ell} \right] \Sigma U^{\mathsf{T}} y$ Neumann series:  $\sum_{\ell=0}^{\infty} A^{\ell} = (I - A)^{-1}$  when  $||A||_{op} < 1$   $\tilde{\eta} < 2/\lambda_{max}(\Sigma^2)$ 

$$n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = \tilde{\eta}$$

$$X = \int W_{k-1} + \frac{\eta}{n} X^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\top} X)^{\ell} X^{\top} y$$

$$n \times d \quad \tilde{V} = \begin{bmatrix} V \quad V_{2} \end{bmatrix} \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\top} \tilde{V} = I = \tilde{V}$$

$$= \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left( I - \tilde{\eta} \begin{bmatrix} \Sigma^{2} & 0 \\ 0 & 0 \end{bmatrix} \right)^{\ell} \quad \tilde{V}^{\top} \quad \tilde{V} \tilde{\Sigma}^{\top} U^{\top}$$

$$\int U^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\top} y$$

![](_page_23_Figure_3.jpeg)

![](_page_23_Picture_4.jpeg)

 $L_{S}(w) = \frac{1}{2n} ||Xw - y||^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with}$ of rank *n*  $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right)$  $X = U\Sigma V^{\mathsf{T}} = U\tilde{\Sigma}\tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{\mathsf{T}}$  $w_k = \tilde{\eta} \sum_{k=1}^{N} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y$  $\ell = 0$  $= \tilde{\eta} \sum_{k=1}^{k-1} \begin{bmatrix} V & V_2 \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^2)^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix}$  $\ell = 0$  $= \tilde{\eta} V \left[ \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \Sigma^2)^{\ell} \right] \Sigma U^{\mathsf{T}} y \to \tilde{\eta} V$ Neumann series:  $\sum_{\ell=0}^{\infty} A^{\ell} = (I - A)^{-1}$  when  $||A||_{op} < 1$   $\tilde{\eta} < 2/\lambda_{max}(\Sigma^2)$ 

$$n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\mathsf{T}} (Xw - y) \quad w_{0} = \tilde{\eta} X$$

$$\left[I - (I - \tilde{\eta}\Sigma^2)\right]^{-1}\Sigma U^{\mathsf{T}} y$$

![](_page_24_Figure_4.jpeg)

![](_page_24_Picture_5.jpeg)

 $L_{S}(w) = \frac{1}{2n} ||Xw - y||^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with}$ of rank *n*  $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right)$  $X = U\Sigma V^{\mathsf{T}} = U\tilde{\Sigma}\tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{\mathsf{T}}$  $w_k = \tilde{\eta} \sum_{k=1}^{n} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y$  $\ell = 0$  $= \tilde{\eta} \sum_{k=1}^{k-1} \begin{bmatrix} V & V_2 \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^2)^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix}$  $\ell = 0$  $= \tilde{\eta} V \left[ \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \Sigma^2)^{\ell} \right] \Sigma U^{\mathsf{T}} y \to \tilde{\eta} V \left[ I - (I - \tilde{\eta} \Sigma^2) \right]^{-1} \Sigma U^{\mathsf{T}} y = V \Sigma^{-1} U^{\mathsf{T}} y$ Neumann series:  $\sum_{\ell=0}^{\infty} A^{\ell} = (I - A)^{-1}$  when  $||A||_{\text{op}} < 1$   $\tilde{\eta} < 2/\lambda_{\max}(\Sigma^2)$ 

$$n < d, \quad \nabla L_{S}(w) = \frac{1}{n} X^{\top} (Xw - y) \quad w_{0} = \tilde{\eta}$$

$$K = \int W_{k-1} + \frac{\eta}{n} X^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} (I - \tilde{\eta} X^{\top} X)^{\ell} X^{\top} y$$

$$n \times d \quad \tilde{V} = \begin{bmatrix} V \quad V_{2} \end{bmatrix} \in \mathbb{R}^{d \times d}; \quad \tilde{V}^{\top} \tilde{V} = I = \tilde{V}$$

$$= \tilde{\eta} \sum_{\ell=0}^{k-1} \tilde{V} \left( I - \tilde{\eta} \begin{bmatrix} \Sigma^{2} & 0 \\ 0 & 0 \end{bmatrix} \right)^{\ell} \quad \tilde{V}^{\top} \quad \tilde{V} \tilde{\Sigma}^{\top} U^{\top}$$

$$\int U^{\top} y = \tilde{\eta} \sum_{\ell=0}^{k-1} V (I - \tilde{\eta} \Sigma^{2})^{\ell} \Sigma U^{\top} y$$

![](_page_25_Figure_4.jpeg)

![](_page_25_Picture_5.jpeg)

 $L_{S}(w) = \frac{1}{2n} ||Xw - y||^{2} \quad X \in \mathbb{R}^{n \times d} \text{ with}$ of rank *n*  $w_k = w_{k-1} - \eta \nabla L_S(w_{k-1}) = \left(I - \frac{\eta}{n} X^{\mathsf{T}} X\right)$  $X = U\Sigma V^{\mathsf{T}} = U\tilde{\Sigma}\tilde{V}^{\mathsf{T}} \quad \tilde{\Sigma} = [\Sigma \quad 0] \in \mathbb{R}^{\mathsf{T}}$  $w_k = \tilde{\eta} \sum_{k=1}^{n} (I - \tilde{\eta} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} \tilde{\Sigma} \tilde{V}^{\mathsf{T}})^{\ell} \tilde{V} \tilde{\Sigma}^{\mathsf{T}} U^{\mathsf{T}} y$  $\ell = 0$  $= \tilde{\eta} \sum_{k=1}^{k-1} \begin{bmatrix} V & V_2 \end{bmatrix} \begin{bmatrix} (I - \tilde{\eta} \Sigma^2)^{\ell} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma \\ 0 \end{bmatrix}$  $\ell = 0$  $= \tilde{\eta} V \left[ \sum_{\ell=0}^{k-1} (I - \tilde{\eta} \Sigma^2)^{\ell} \right] \Sigma U^{\mathsf{T}} y \to \tilde{\eta} V \left[ I - (I - \tilde{\eta} \Sigma^2) \right]^{-1} \Sigma U^{\mathsf{T}} y = V \Sigma^{-1} U^{\mathsf{T}} y = X^{\dagger} y$ Neumann series:  $\sum_{\ell=0}^{\infty} A^{\ell} = (I - A)^{-1}$  when  $||A||_{\text{op}} < 1$   $\tilde{\eta} < 2/\lambda_{\max}(\Sigma^2)$ 

![](_page_26_Figure_4.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

![](_page_27_Picture_5.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

![](_page_28_Picture_5.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

- So, the 1,000-degree polynomial picture is what GD would give
- If we track  $w_0 \neq 0$  in same analysis, get  $w_{\infty} = V_2 V_2^{\mathsf{T}} w_0 + X^{\dagger} y$  (proof)

![](_page_29_Picture_6.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

- Does this same idea hold for other losses / models? Not necessarily.
- So, the 1,000-degree polynomial picture is what GD would give
- If we track  $w_0 \neq 0$  in same analysis, get  $w_{\infty} = V_2 V_2^{\mathsf{T}} w_0 + X^{\dagger} y$  (proof)

![](_page_30_Picture_7.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

- Does this same idea hold for other losses / models? Not necessarily.
- So, the 1,000-degree polynomial picture is what GD would give
  - Logistic regression:

![](_page_31_Picture_8.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

- Does this same idea hold for other losses / models? Not necessarily.
- So, the 1,000-degree polynomial picture is what GD would give
  - Logistic regression:
    - Separable: norm diverges in direction of max-margin separator (Soudry et al.)

![](_page_32_Picture_9.jpeg)

![](_page_32_Picture_10.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

- So, the 1,000-degree polynomial picture is what GD would give Does this same idea hold for other losses / models? Not necessarily.
- - Logistic regression:
    - Separable: norm diverges in direction of max-margin separator (Soudry et al.) Non-separable: biased towards max-margin, but complicated (<u>Ji/Telgarsky</u>)

![](_page_33_Picture_10.jpeg)

![](_page_33_Figure_11.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

- So, the 1,000-degree polynomial picture is what GD would give • Does this same idea hold for other losses / models? Not necessarily.
- - Logistic regression:
    - Separable: norm diverges in direction of max-margin separator (Soudry et al.) • Non-separable: biased towards max-margin, but complicated (Ji/Telgarsky)

    - Also see <u>Telgarsky notes section 10</u>

![](_page_34_Picture_11.jpeg)

![](_page_34_Figure_12.jpeg)

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

- So, the 1,000-degree polynomial picture is what GD would give • Does this same idea hold for other losses / models? Not necessarily.
- - Logistic regression:
    - Separable: norm diverges in direction of max-margin separator (Soudry et al.) • Non-separable: biased towards max-margin, but complicated (Ji/Telgarsky) • Also see <u>Telgarsky notes section 10</u>
  - Matrix factorization models: <u>conjectured</u> min nuclear norm, slightly controversial

![](_page_35_Picture_12.jpeg)

![](_page_35_Figure_13.jpeg)
#### Implicit regularization of gradient descent

• We just showed that gradient descent for OLS with X of rank n, starting from zero with  $\eta < 2n / \sigma_{\max}(X)^2$ , converges to the minimum-norm interpolator  $X^{\dagger}y$ 

• If we track  $w_0 \neq 0$  in same analysis, get  $w_{\infty} = V_2 V_2^{\top} w_0 + X^{\dagger} y$  (proof)

- So, the 1,000-degree polynomial picture is what GD would give • Does this same idea hold for other losses / models? Not necessarily.
- - Logistic regression:
    - Separable: norm diverges in direction of max-margin separator (Soudry et al.) • Non-separable: biased towards max-margin, but complicated (Ji/Telgarsky) Also see <u>Telgarsky notes section 10</u>
  - Matrix factorization models: <u>conjectured</u> min nuclear norm, slightly controversial
  - Deep learning: ???







model predictors  $h_{n,N}$  learned on a subset of MNIST ( $n = 10^4$ , 10 classes). The interpolation threshold is achieved at  $N = 10^4$ .

Fig. 2. Double-descent risk curve for the RFF model on MNIST. Shown are test risks (log scale), coefficient  $\ell_2$  norms (log scale), and training risks of the RFF 6





**Classical regime** (left of peak): unique ERM

model predictors  $h_{n,N}$  learned on a subset of MNIST ( $n = 10^4$ , 10 classes). The interpolation threshold is achieved at  $N = 10^4$ .

Fig. 2. Double-descent risk curve for the RFF model on MNIST. Shown are test risks (log scale), coefficient  $\ell_2$  norms (log scale), and training risks of the RFF 6





**Classical regime** (left of peak): unique ERM

model predictors  $h_{n,N}$  learned on a subset of MNIST ( $n = 10^4$ , 10 classes). The interpolation threshold is achieved at  $N = 10^4$ .

Fig. 2. Double-descent risk curve for the RFF model on MNIST. Shown are test risks (log scale), coefficient  $\ell_2$  norms (log scale), and training risks of the RFF 6





**Classical regime** (left of peak): unique ERM

model predictors  $h_{n,N}$  learned on a subset of MNIST ( $n = 10^4$ , 10 classes). The interpolation threshold is achieved at  $N = 10^4$ .

Fig. 2. Double-descent risk curve for the RFF model on MNIST. Shown are test risks (log scale), coefficient  $\ell_2$  norms (log scale), and training risks of the RFF 6





Fig. 3. interpolation threshold (black dashed line) is observed at  $n \cdot K$ .



Double-descent risk curve for a fully connected neural network Fig. 4. Double-descent risk curve for random forests on MNIST. The doubleon MNIST. Shown are training and test risks of a network with a single descent risk curve is observed for random forests with increasing model layer of H hidden units, learned on a subset of MNIST ( $n = 4 \cdot 10^3$ , d = 784, complexity trained on a subset of MNIST ( $n = 10^4$ , 10 classes). Its complex-K = 10 classes). The number of parameters is  $(d + 1) \cdot H + (H + 1) \cdot K$ . The ity is controlled by the number of trees N<sub>tree</sub> and the maximum number of leaves allowed for each tree  $N_{leaf}^{max}$ .















# More data hurts!

75 100 125 150 175 200 Embedding Dimension (Transformer Model Size)

Nakkiran et al. ICLR-20



#### Test Error





procedure  $\mathcal{T}$ , with respect to distribution  $\mathcal{D}$  and parameter  $\epsilon > 0$ , is defined as:

where  $\operatorname{Error}_{S}(M)$  is the mean error of model M on train samples S.

Our main hypothesis can be informally stated as follows:

predicting labels based on n samples from D then:

that increases its effective complexity will decrease the test error.

that increases its effective complexity will decrease the test error.

effective complexity might decrease or increase the test error.

**Definition 1 (Effective Model Complexity)** *The* Effective Model Complexity (*EMC*) of a training

- $\mathrm{EMC}_{\mathcal{D},\epsilon}(\mathcal{T}) := \max \left\{ n \mid \mathbb{E}_{S \sim \mathcal{D}^n}[\mathrm{Error}_S(\mathcal{T}(S))] \le \epsilon \right\}$
- Hypothesis 1 (Generalized Double Descent hypothesis, informal) For any natural data distribution D, neural-network-based training procedure T, and small  $\epsilon > 0$ , if we consider the task of
- **Under-paremeterized regime.** If  $\text{EMC}_{\mathcal{D},\epsilon}(\mathcal{T})$  is sufficiently smaller than n, any perturbation of  $\mathcal{T}$
- **Over-parameterized regime.** If  $\text{EMC}_{\mathcal{D},\epsilon}(\mathcal{T})$  is sufficiently larger than n, any perturbation of  $\mathcal{T}$
- Critically parameterized regime. If  $\text{EMC}_{\mathcal{D},\epsilon}(\mathcal{T}) \approx n$ , then a perturbation of  $\mathcal{T}$  that increases its







(pause)

- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis

- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$

- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$

- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$
  - Make predictions/decision based on posterior mean/median, MAP, single draw, ...



- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$
- Make predictions/decision based on posterior mean/median, MAP, single draw, ... This is optimal if you believe in your prior + likelihood!



- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$
- Make predictions/decision based on posterior mean/median, MAP, single draw, ... This is optimal if you believe in your prior + likelihood!
- Frequentists say: "but how good is it actually???"



- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$
- Make predictions/decision based on posterior mean/median, MAP, single draw, ... This is optimal if you believe in your prior + likelihood! Frequentists say: "but how good is it actually???"
- - What if your model class / prior / ... are wrong?



- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$
- Make predictions/decision based on posterior mean/median, MAP, single draw, ... This is optimal if you believe in your prior + likelihood!
  - Frequentists say: "but how good is it actually???"
  - What if your model class / prior / ... are wrong?
- Tempered likelihood (<u>Zhang 06</u>) / SafeBayes (<u>Grünwald 12</u>):



- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$
- Make predictions/decision based on posterior mean/median, MAP, single draw, ... This is optimal if you believe in your prior + likelihood!
  - Frequentists say: "but how good is it actually???"
  - What if your model class / prior / ... are wrong?
- Tempered likelihood (<u>Zhang 06</u>) / SafeBayes (<u>Grünwald 12</u>):
  - If your model is misspecified, can be provably better to use  $\mathscr{L}^{\lambda}$  for  $\lambda < 1$



- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$
- Make predictions/decision based on posterior mean/median, MAP, single draw, ... This is optimal if you believe in your prior + likelihood!
  - Frequentists say: "but how good is it actually???"
  - What if your model class / prior / ... are wrong?
- Tempered likelihood (<u>Zhang 06</u>) / SafeBayes (<u>Grünwald 12</u>):
  - If your model is misspecified, can be provably better to use  $\mathscr{L}^{\lambda}$  for  $\lambda < 1$
  - No longer quite Bayesian inference, but turns a prior into a posterior



- Bayesians say:
  - Start with a prior distribution  $\pi(h)$  on choice of hypothesis
  - Observe data S with likelihood  $\mathscr{L}(S \mid h)$
  - End up with posterior distribution  $\rho(h \mid S) \propto \mathscr{L}(S \mid h) \pi(h)$
- Make predictions/decision based on posterior mean/median, MAP, single draw, ... This is optimal if you believe in your prior + likelihood!
  - Frequentists say: "but how good is it actually???"
  - What if your model class / prior / ... are wrong?
- Tempered likelihood (<u>Zhang 06</u>) / SafeBayes (<u>Grünwald 12</u>):
  - If your model is misspecified, can be provably better to use  $\mathscr{L}^{\lambda}$  for  $\lambda < 1$
  - No longer quite Bayesian inference, but turns a prior into a posterior
- PAC-Bayes: analyzes any prior-posterior pair (potentially even totally unrelated)



• We start with some prior  $\pi$  (independent of the data S) on hypotheses

- Our learning algorithm sees S and gives us a posterior  $\rho$

• We start with some prior  $\pi$  (independent of the data S) on hypotheses

- We start with some prior  $\pi$  (independent of the data S) on hypotheses
- Our learning algorithm sees S and gives us a posterior  $\rho$
- We'll analyze  $L_{\mathscr{D}}(\rho) = \mathbb{E}_{h\sim\rho}[L_{\mathscr{D}}(h)]$  based on  $L_{S}(\rho) = \mathbb{E}_{h\sim\rho}[L_{S}(h)]$

- Our learning algorithm sees S and gives us a posterior  $\rho$
- We'll analyze  $L_{\mathcal{D}}(\rho) = \mathbb{E}_{h\sim\rho}[L_{\mathcal{D}}(h)]$  based on  $L_{S}(\rho) = \mathbb{E}_{h\sim\rho}[L_{S}(h)]$
- McAllester-style bound (SSBD theorem 31.1):

• We start with some prior  $\pi$  (independent of the data S) on hypotheses

- We start with some prior  $\pi$  (independent of the data S) on hypotheses
- Our learning algorithm sees S and gives us a posterior  $\rho$
- We'll analyze  $L_{\mathscr{D}}(\rho) = \mathbb{E}_{h\sim\rho}[L_{\mathscr{D}}(h)]$  based on  $L_{S}(\rho) = \mathbb{E}_{h\sim\rho}[L_{S}(h)]$ • McAllester-style bound (SSBD theorem 31.1):
- - If  $\ell(h, z) \in [0, 1]$ , with probability at least  $1 \delta$  over  $S \sim \mathscr{D}^n$ ,  $\leq \sqrt{\frac{\mathrm{KL}(\rho \| \pi) + \log \frac{n}{\delta}}{2(n-1)}}$   $\frac{\rho(h)}{\pi(h)} \text{ (the usual KL divergence)}$

$$L_{\mathcal{D}}(\rho) - L_{S}(\rho) \leq$$

where 
$$\operatorname{KL}(\rho \| \pi) = \mathbb{E}_{h \sim \rho} \log \frac{\rho}{\pi}$$

- We start with some prior  $\pi$  (independent of the data S) on hypotheses
- Our learning algorithm sees S and gives us a posterior  $\rho$
- We'll analyze  $L_{\mathscr{D}}(\rho) = \mathbb{E}_{h\sim\rho}[L_{\mathscr{D}}(h)]$  based on  $L_{S}(\rho) = \mathbb{E}_{h\sim\rho}[L_{S}(h)]$ • McAllester-style bound (SSBD theorem 31.1):
- - If  $\ell(h, z) \in [0, 1]$ , with probability at least  $1 \delta$  over  $S \sim \mathscr{D}^n$ , 
    $$\begin{split} L_{\mathscr{D}}(\rho) - L_{S}(\rho) &\leq \sqrt{\frac{\mathrm{KL}(\rho \| \pi) + \log \frac{n}{\delta}}{2(n-1)}} \\ \text{where } \mathrm{KL}(\rho \| \pi) &= \mathbb{E}_{h \sim \rho} \log \frac{\rho(h)}{\pi(h)} \text{ (the usual KL divergence)} \end{split}$$

$$L_{\mathcal{D}}(\rho) - L_{S}(\rho) \leq$$

 $\mathcal{I}(\mathcal{I})$ Proved in SSBD chapter 31 (not bad at all)

 $L_{\mathcal{D}}(\rho) - L_{S}(\rho) \leq$ 

• What's the best learning algorithm, according to this bound?

ng algorithm?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$

 $L_{\mathcal{D}}(\rho) - L_{\mathcal{S}}(\rho) \leq$ 

- What's the best learning algorithm, according to this bound?
  - Turns out to be the Gibbs posterior:  $\rho(h) \propto \exp(-\lambda L_S(h)) \pi(h)$

ng algorithm?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$

according to this bound? or:  $\rho(h) \propto \exp(-\lambda L_S(h)) \pi(h)$ 

#### $L_{\mathcal{D}}(\rho) - L_{\mathcal{S}}(\rho) \leq$

- What's the best learning algorithm, according to this bound?
  - Turns out to be the Gibbs posterior:  $\rho(h) \propto \exp(-\lambda L_{S}(h)) \pi(h)$

ng algorithm?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$



#### $L_{O}(\rho) - L_{S}(\rho) \leq$

- What's the best learning algorithm, according to this bound?
  - Turns out to be the Gibbs posterior:  $\rho(h) \propto \exp(-\lambda L_S(h)) \pi(h)$

  - Typical choice (see 340): e.g. squared loss  $\leftrightarrow$  Gaussian likelihood

ng algorithm?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$



#### $L_{O}(\rho) - L_{S}(\rho) \leq$

- What's the best learning algorithm, according to this bound?
  - Turns out to be the Gibbs posterior:  $\rho(h) \propto \exp(-\lambda L_S(h)) \pi(h)$

  - Typical choice (see 340): e.g. squared loss  $\leftrightarrow$  Gaussian likelihood
- But the bound applies to any prior-posterior pair (with  $\pi$  independent of S)

**ng algorithm?**  
$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$



#### $L_{O}(\rho) - L_{S}(\rho) \leq$

- What's the best learning algorithm, according to this bound?
  - Turns out to be the Gibbs posterior:  $\rho(h) \propto \exp(-\lambda L_S(h)) \pi(h)$

  - Typical choice (see 340): e.g. squared loss  $\leftrightarrow$  Gaussian likelihood
- But the bound applies to any prior-posterior pair (with  $\pi$  independent of S)
  - For instance: could learn a  $\hat{h}$  with (S)GD and then add noise to it

ng algorithm?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$



#### $L_{O}(\rho) - L_{S}(\rho) \leq$

- What's the best learning algorithm, according to this bound?
  - Turns out to be the Gibbs posterior:  $\rho(h) \propto \exp(-\lambda L_S(h)) \pi(h)$

  - Typical choice (see 340): e.g. squared loss  $\leftrightarrow$  Gaussian likelihood
- - For instance: could learn a  $\hat{h}$  with (S)GD and then add noise to it
  - If h is in a flat minimum, then h + noise will still be good

**ng algorithm?**  
$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$

• Same as tempered likelihood / SafeBayes if  $\mathscr{L}(S \mid h) = -\log L_S(h) + \text{const}$ 

• But the bound applies to any prior-posterior pair (with  $\pi$  independent of S)


# What learni

#### $L_{O}(\rho) - L_{S}(\rho) \leq$

- What's the best learning algorithm, according to this bound?
  - Turns out to be the Gibbs posterior:  $\rho(h) \propto \exp(-\lambda L_S(h)) \pi(h)$

  - Typical choice (see 340): e.g. squared loss  $\leftrightarrow$  Gaussian likelihood
- But the bound applies to any prior-posterior pair (with  $\pi$  independent of S)
  - For instance: could learn a  $\hat{h}$  with (S)GD and then add noise to it
  - If  $\hat{h}$  is in a flat minimum, then h + noise will still be good
  - But note that if  $\rho \to \text{point mass}$  and  $\pi$  continuous,  $\text{KL}(\rho \| \pi) \to \infty$

ng algorithm?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$

• Same as tempered likelihood / SafeBayes if  $\mathscr{L}(S \mid h) = -\log L_S(h) + \text{const}$ 



- What's the best prior?
  - Bound on generalization gap is better if  $\rho$  is "closer" to  $\pi$

t prior?  

$$\sqrt{KL(\rho \| \pi) + \log \frac{n}{\delta}}$$

$$\frac{2(n-1)}{\delta}$$

- What's the best prior?
  - Bound on generalization gap is better if  $\rho$  is "closer" to  $\pi$ • S didn't make us "change our mind" too much – similar to MDL

t prior?  

$$\sqrt{KL(\rho \| \pi) + \log \frac{n}{\delta}}$$

$$\frac{1}{2(n-1)}$$

### $L_{\mathcal{D}}(\rho) - L_{S}(\rho) \leq$

- What's the best prior?
  - Bound on generalization gap is better if  $\rho$  is "closer" to  $\pi$ 
    - S didn't make us "change our mind" too much similar to MDL
  - But we also want a good ho, i.e. average training loss  $L_{\rm S}(
    ho)$  should be small

t prior?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$

- What's the best prior?
  - Bound on generalization gap is better if  $\rho$  is "closer" to  $\pi$ 
    - S didn't make us "change our mind" too much similar to MDL
  - But we also want a good  $\rho$ , i.e. average training loss  $L_{S}(\rho)$  should be small
- Notice  $\pi$  only shows up in the bound nothing to do with the learning algorithm

t prior?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$



- What's the best prior?
  - Bound on generalization gap is better if  $\rho$  is "closer" to  $\pi$ 
    - S didn't make us "change our mind" too much similar to MDL
  - But we also want a good  $\rho$ , i.e. average training loss  $L_{S}(\rho)$  should be small
- Notice  $\pi$  only shows up in the bound nothing to do with the learning algorithm • We could potentially pick a prior that actually depends on  $\mathscr{D}$

t prior?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$



- What's the best prior?
  - Bound on generalization gap is better if  $\rho$  is "closer" to  $\pi$ 
    - S didn't make us "change our mind" too much similar to MDL
  - But we also want a good  $\rho$ , i.e. average training loss  $L_{S}(\rho)$  should be small
- Notice  $\pi$  only shows up in the bound nothing to do with the learning algorithm • We could potentially pick a prior that actually depends on  $\mathscr{D}$ • ...as long as we can still bound  $KL(\rho \| \pi)$

t prior?  

$$\sqrt{\frac{KL(\rho \parallel \pi) + \log \frac{n}{\delta}}{2(n-1)}}$$



## **Other forms of PAC-Bayes bounds**

• Linear bound: 
$$L_{\mathscr{D}}(\rho) \leq \frac{1}{\beta}L_{S}(\rho)$$
 -

- - Can be much tighter (unfortunately) if  $KL(\rho \| \pi)/n$  is big
- Also variants based on <u>general f-divergences</u>, <u>Wasserstein</u>, …

 $+\frac{\mathrm{KL}(\rho \| \pi) + \log \frac{1}{\delta}}{2\beta(1-\beta)n} \text{ for any } \beta \in (0,1)$ 

• Catoni bound: for  $\alpha > 1$ ,  $\Phi_{\gamma}^{-1}(x) = (1 - \exp(-\gamma x))/(1 - \exp(-\gamma))$ ,  $L_{\mathcal{D}}(\rho) \le \inf_{\lambda > 1} \Phi_{\lambda/n}^{-1} \left( L_{S}(\rho) + \frac{\alpha}{\lambda} \left[ \operatorname{KL}(\rho \| \pi) - \log \varepsilon + 2\log \frac{\log(\alpha^{2}\lambda)}{\log \alpha} \right] \right)$ 



#### **NON-VACUOUS GENERALIZATION BOUNDS AT THE IM-**AGENET SCALE: A PAC-BAYESIAN COMPRESSION APPROACH

Wenda Zhou Columbia University New York, NY

Victor Veitch Columbia University New York, NY wz2335@columbia.edu victorveitch@gmail.com

Ryan P. Adams Princeton University Princeton, NJ rpa@princeton.edu

Morgane Austern **Columbia University** New York, NY ma3293@columbia.edu

**Peter Orbanz Columbia University** New York, NY porbanz@stat.columbia.edu

#### • Pre-pick a coding scheme to represent networks (e.g. compress the weights) • Train a network with SGD, sparsify it/etc to $\hat{h}$ , then add a little noise to weights

#### **NON-VACUOUS GENERALIZATION BOUNDS AT THE IM-**AGENET SCALE: A PAC-BAYESIAN COMPRESSION APPROACH

Wenda Zhou Columbia University New York, NY

Victor Veitch Columbia University New York, NY wz2335@columbia.edu victorveitch@gmail.com

Ryan P. Adams Princeton University Princeton, NJ rpa@princeton.edu

**Table 1:** Summary of bounds obtained from compression

Dataset	Orig. size	Comp. size	Robust. Adj.	Eff. Size	Error Bound	
					Top 1	Top 5
MNIST	$168.4{ m KiB}$	$8.1{ m KiB}$	$1.88{ m KiB}$	$6.23{ m KiB}$	< 46%	NA
ImageNet	$5.93{ m MiB}$	$452{ m KiB}$	$102{ m KiB}$	$350{ m KiB}$	< 96.5%	< 89%

Morgane Austern **Columbia University** New York, NY ma3293@columbia.edu

**Peter Orbanz Columbia University** New York, NY porbanz@stat.columbia.edu

#### • Pre-pick a coding scheme to represent networks (e.g. compress the weights) • Train a network with SGD, sparsify it/etc to $\hat{h}$ , then add a little noise to weights

• In practice, we don't actually use randomized predictors (usually)

- In practice, we don't actually use randomized predictors (usually)

• Possible to "derandomize" to a high-probability bound on  $L_{\Im}(h) - L_{S}(h)$ :

- In practice, we don't actually use randomized predictors (usually)

• Possible to "derandomize" to a high-probability bound on  $L_{O}(h) - L_{S}(h)$ :

• Show convergence of  $L_{\mathscr{D}}(h)$  to  $\mathbb{E}_{h\sim\rho}L_{\mathscr{D}}(h)$ ,  $L_{S}(h)$  to  $\mathbb{E}_{h\sim\rho}L_{S}(h)$ , under  $\rho$ 

- In practice, we don't actually use randomized predictors (usually)
- - Show convergence of  $L_{\mathcal{O}}(h)$  to  $\mathbb{E}$

• Possible to "derandomize" to a high-probability bound on  $L_{O}(h) - L_{S}(h)$ :

$$\mathbb{E}_{h\sim\rho}L_{\mathscr{D}}(h), L_{S}(h)$$
 to  $\mathbb{E}_{h\sim\rho}L_{S}(h)$ , under  $\rho$ 

• Or, use a margin-type loss to show 0-1 error doesn't change under  $\rho$ 

- In practice, we don't actually use randomized predictors (usually)
- Possible to "derandomize" to a high-probability bound on  $L_{\Im}(h) L_{S}(h)$ :
  - Show convergence of  $L_{\infty}(h)$  to  $\mathbb{E}$
- Or, use a margin-type loss to show 0-1 error doesn't change under  $\rho$ But...these then become "two-sided" bounds

$$\mathbb{E}_{h\sim\rho}L_{\mathscr{D}}(h), L_{S}(h)$$
 to  $\mathbb{E}_{h\sim\rho}L_{S}(h)$ , under  $\rho$ 

- In practice, we don't actually use randomized predictors (usually)
- Possible to "derandomize" to a high-probability bound on  $L_{\Im}(h) L_{S}(h)$ :
  - Show convergence of  $L_{\mathcal{O}}(h)$  to  $\mathbb{E}$
- Or, use a margin-type loss to show 0-1 error doesn't change under  $\rho$  But...these then become "two-sided" bounds
  - Subject to the <u>Nagarajan/Kolter</u> failure mode (their Appendix J)

**Uniform convergence may be unable to explain** generalization in deep learning

20

Vaishnavh Nagarajan Department of Computer Science Carnegie Mellon University Pittsburgh, PA vaishnavh@cs.cmu.edu

$$\mathbb{E}_{h\sim\rho}L_{\mathscr{D}}(h), L_{S}(h)$$
 to  $\mathbb{E}_{h\sim\rho}L_{S}(h)$ , under  $\rho$ 

J. Zico Kolter Department of Computer Science Carnegie Mellon University & Bosch Center for Artificial Intelligence Pittsburgh, PA zkolter@cs.cmu.edu

- Double descent
  - Classical behaviour, then descend again after interpolation peak
  - Highly dependent on learning algorithm's implicit regularization
- PAC-Bayes
  - <u>A Primer on PAC-Bayesian Learning (Guedi 2019)</u>

  - Still, significant practical issues, not always super explanatory
  - Lots of ongoing fancy variations

### Recap

Tightest bounds we know for deep learning (afaik)...but still not that tight