

# Risk Comparisons in Linear Regression

## Implicit Regularization Dominates Explicit Regularization

Jingfeng Wu<sup>1</sup> Peter Bartlett<sup>\*13</sup> Jason Lee<sup>\*1</sup>  
Sham Kakade<sup>\*23</sup> Bin Yu<sup>\*1</sup>

<sup>1</sup>UC Berkeley <sup>2</sup>Harvard <sup>3</sup>Google DeepMind

### Linear regression

**task.**  $x \sim N(0, \Sigma)$ ,  $y = x^\top w^* + N(0, 1)$  for  $\|w^*\|_\Sigma \lesssim 1$  **problem determined by  $(\Sigma, w^*)$**

**risk.**  $R(w) = \mathbb{E}(y - x^\top w)^2 - \mathbb{E}(y - x^\top w^*)^2$   
 $= \|w - w^*\|_\Sigma^2$

**data.**  $n$  iid samples  $(x_1, y_1), \dots, (x_n, y_n)$   $X = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$   $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

### Algorithms

**ridge.**  $w_\lambda^{\text{ridge}} = \arg \min \frac{1}{n} \sum_{i=1}^n \|x_i^\top w - y_i\|^2 + \lambda \|w\|^2$   
 $= (X^\top X + n\lambda I)^{-1} X^\top Y$  **hyperparameter:  $\lambda \geq 0$**

**gradient descent.**

$w_0 = 0$   
for  $s = 1, \dots, t$ ,  
 $w_s = w_{s-1} - \frac{\eta}{n} X^\top (X w_{s-1} - Y)$   
 $w_t^{\text{gd}} = w_t$

**hyperparameter:  $t \geq 0$**

**stochastic gradient descent.**

$w_0 = 0, \eta_0 = \eta, N = n/\log n$   
for  $i = 1, \dots, n$ ,  
 $\eta_i = \begin{cases} 0.1\eta_{i-1} & \text{if } i \% N = 0 \\ \eta_{i-1} & \text{else} \end{cases}$   
 $w_i = w_{i-1} - \eta_i (x_i^\top w_{i-1} - y_i) x_i$   
 $w_n^{\text{sgd}} = w_n$  **hyperparameter:  $0 < \eta \lesssim 1/\text{tr}(\Sigma)$**

### Prior results

[Tsigler & Bartlett, 2023]

For all  $\lambda \geq 0$ , in expectation

$$\mathbb{E}R(w_\lambda^{\text{ridge}}) \gtrsim \tilde{\lambda}^2 \|w^*\|_{\Sigma_{0:k^*}^{-1}}^2 + \|w^*\|_{\Sigma_{k^*:\infty}}^2 + \min \left\{ \frac{D}{n}, 1 \right\}$$

**critical index**  $k^* = \min \left\{ k : \lambda + \frac{\sum_{i>k} \lambda_i}{n} \geq c\lambda_{k+1} \right\}$

**effective regularization**  $\tilde{\lambda} = \lambda + \frac{\sum_{i>k^*} \lambda_i}{n}$

**effective dimension**  $D = k^* + \frac{1}{\tilde{\lambda}^2} \sum_{i>k^*} \lambda_i^2$

[Wu\*, Zou\*, Braverman, Gu, Kakade, 2022]

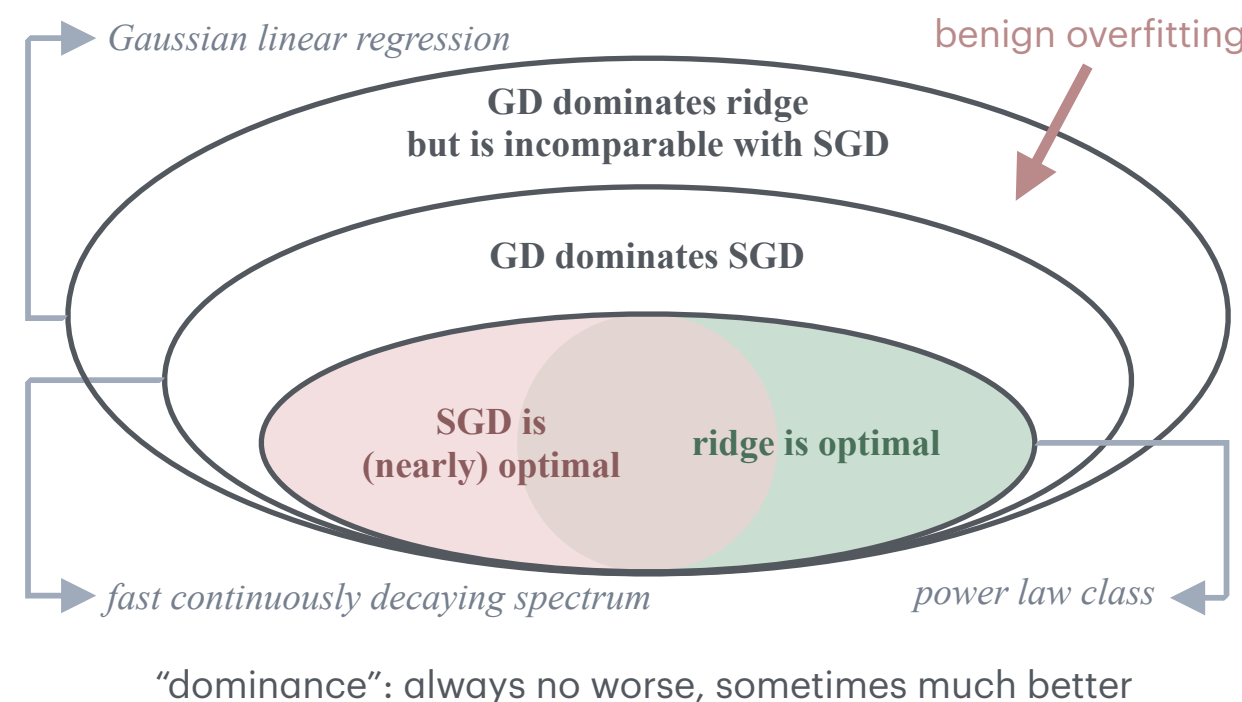
For all  $0 < \eta \lesssim 1/\text{tr}(\Sigma)$ , in expectation

$$\mathbb{E}R(w_\eta^{\text{sgd}}) \approx \left\| \prod_{i=1}^n (I - \eta_i \Sigma) w^* \right\|_\Sigma^2 + \frac{D}{N}$$

**effective steps**  $N = n/\log n$

**critical index**  $k^* := \min \left\{ \frac{1}{\eta N} \geq c\lambda_{k+1} \right\}$

**effective dimension**  $D = k^* + \eta^2 N^2 \sum_{i>k^*} \lambda_i^2$



### GD dominates ridge

**Theorem**

For all  $0 < \eta \lesssim 1/\text{tr}(\Sigma)$  and  $t \geq 0$ , w.h.p.

$$R(w_t^{\text{gd}}) \lesssim \tilde{\lambda}^2 \|w^*\|_{\Sigma_{0:k^*}^{-1}}^2 + \|w^*\|_{\Sigma_{k^*:\infty}}^2 + \frac{D}{n}$$

**critical index**  $k^* = \min \left\{ k : \frac{1}{\eta t} + \frac{\sum_{i>k} \lambda_i}{n} \geq c\lambda_{k+1} \right\}$

**effective regularization**  $\tilde{\lambda} = \frac{1}{\eta t} + \frac{\sum_{i>k^*} \lambda_i}{n}$

**effective dimension**  $D = k^* + \frac{1}{\tilde{\lambda}^2} \sum_{i>k^*} \lambda_i^2$

**Corollary**

For every Gaussian linear regression,  $n \geq 1$ , and  $\lambda \geq 0$ , there is  $t$  such that: w.h.p.

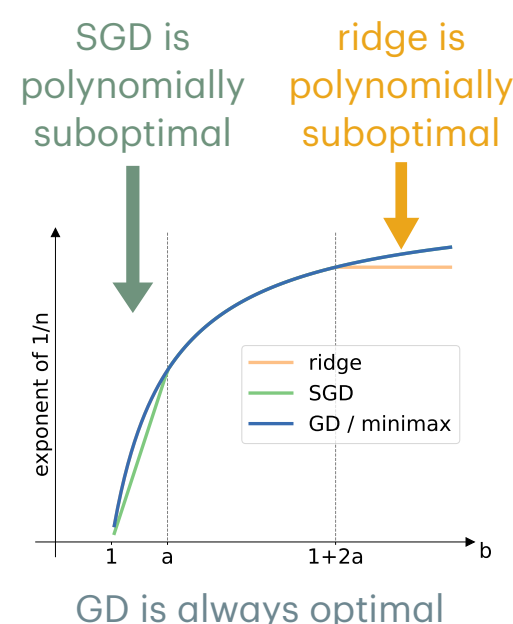
$$R(w_t^{\text{gd}}) \lesssim \mathbb{E}R(w_\lambda^{\text{ridge}})$$

**Proof.** If  $D > n$ , set  $t = 0$ ; otherwise, set  $t = 1/(\eta\lambda)$ .

### Power law class

$\lambda_i \approx i^{-a}$   $\lambda_i (u_i^\top w^*)^2 \approx i^{-b}$  for  $a, b > 1$

	1<b<a	a<b<1+2a	b>1+2a
ridge	$O(n^{-\frac{b-1}{b}})$		$\Omega(n^{-\frac{2a}{1+2a}})$
SGD	$\tilde{\Omega}(n^{-\frac{b-1}{a}})$	$\tilde{O}(n^{-\frac{b-1}{b}})$	
GD		$O(n^{-\frac{b-1}{b}})$	
minimax		$\Omega(n^{-\frac{b-1}{b}})$	



### GD is incomparable with SGD

**Theorem**

For all  $0 < \eta \lesssim 1/\text{tr}(\Sigma)$  and  $t \geq 0$

$$\mathbb{E}R(w_t^{\text{gd}}) \gtrsim \left( \frac{\sum_{i>\ell^*} \lambda_i}{n} \right)^2 \|w^*\|_{\Sigma_{0:\ell^*}^{-1}}^2 + \|w^*\|_{\Sigma_{\ell^*:\infty}}^2 + \min \left\{ \frac{D}{n}, 1 \right\}$$

**effective dimension**  $D = k^* + \frac{1}{\tilde{\lambda}^2} \sum_{i>k^*} \lambda_i^2$  as before...

**benign overfitting index**  $\ell^* = \min \left\{ k : \frac{\sum_{i>k} \lambda_i}{n} \geq c\lambda_{k+1} \right\}$

**Corollary**

$n \geq 1$ . For a sequence of  $d$ -dim problems

$$d \geq n^2 \quad w^* = \begin{bmatrix} n^{0.45} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} n^{-0.9} & & & \\ & 1/d & & \\ & & \ddots & \\ & & & 1/d \end{bmatrix}$$

we have  $\|w^*\|_\Sigma^2 \leq 1$ , moreover

• for all  $0 < \eta \lesssim 1$  and  $t \geq 0$ ,  $\mathbb{E}R(w_t^{\text{gd}}) = \Omega(n^{-0.2})$

• for  $\eta \approx 1$ ,  $\mathbb{E}R(w_\eta^{\text{sgd}}) = O(\log(n)/n)$

### GD dominates SGD in a significant subset

**Theorem**

For all  $0 < \eta \lesssim 1/\text{tr}(\Sigma)$  and  $0 \leq t \lesssim n$ , w.h.p.

$$R(w_t^{\text{gd}}) \lesssim \left\| (I - \eta \Sigma)^{t/2} w^* \right\|_\Sigma^2 + \frac{D}{n} + \left( \frac{D_1}{n} \right)^2$$

**critical index**  $k^* := \min \left\{ \frac{1}{\eta t} \geq c\lambda_{k+1} \right\}$

**effective dimension**  $D = k^* + \eta^2 t^2 \sum_{i>k^*} \lambda_i^2$

**order-1 effective dim**  $D_1 = k^* + \eta t \sum_{i>k^*} \lambda_i$

**Assumption**

Spectrum decays fast and continuously:

$$\text{for all } \tau > 1, \quad \tau \sum_{\lambda_i < 1/\tau} \lambda_i \lesssim \#\{\lambda_i \geq 1/\tau\}$$

**Corollary**

For every Gaussian linear regression satisfying the above,  $n \geq 1$ , and  $0 \leq \eta \lesssim 1$ , there is  $t$  such that

$$\mathbb{E}R(w_t^{\text{gd}}) \lesssim \mathbb{E}R(w_\eta^{\text{sgd}})$$

**Proof.** Assumption implies  $D_1 \lesssim k^* \leq D$ .