

Memorize to Generalize: on the Necessity of Interpolation in High Dimensional Linear Regression

Chen Cheng



Stanford University

Department of Statistics



Chen Cheng
Stanford Stat



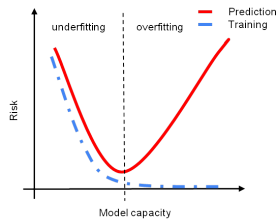
John Duchi
Stanford Stat & EE



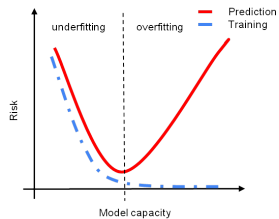
Rohith Kuditipudi
Stanford CS

Introduction: why study memorization?

Classical statistical wisdom

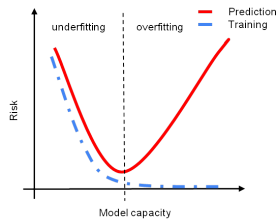


Classical statistical wisdom



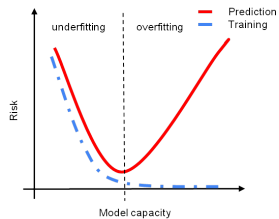
- bigger models tend to overfit

Classical statistical wisdom



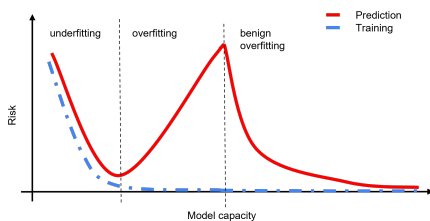
- bigger models tend to overfit
- need to limit model capacity

Classical statistical wisdom

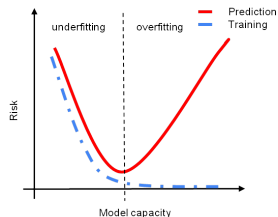


- bigger models tend to overfit
- need to limit model capacity

Modern empirical wisdom

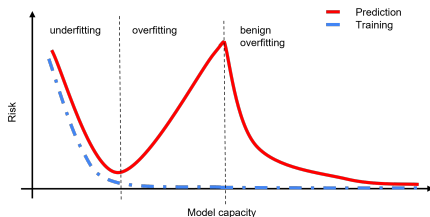


Classical statistical wisdom



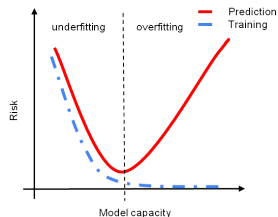
- bigger models tend to overfit
- need to limit model capacity

Modern empirical wisdom



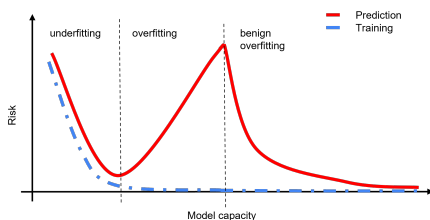
- overparameterized model

Classical statistical wisdom



- bigger models tend to overfit
- need to limit model capacity

Modern empirical wisdom

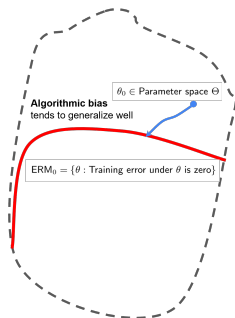


- overparameterized model
- interpolate data

When is it sufficient to overfit?

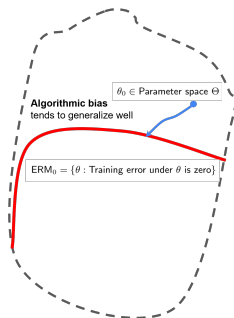
When is it sufficient to overfit?

Benign overfitting



When is it sufficient to overfit?

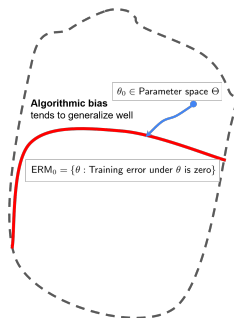
Benign overfitting



Literatures

When is it sufficient to overfit?

Benign overfitting

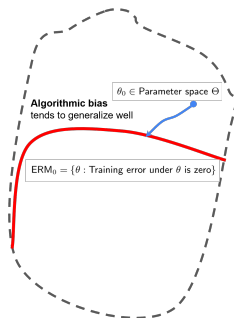


Literatures

- Surprises in high-dimensional ridgeless least squares interpolation
Hastie et al., 2018.

When is it sufficient to overfit?

Benign overfitting

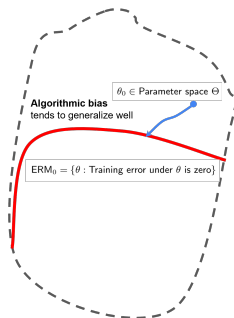


Literatures

- Surprises in high-dimensional ridgeless least squares interpolation
Hastie et al., 2018.
- The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime
Montanari et al., 2019.

When is it sufficient to overfit?

Benign overfitting

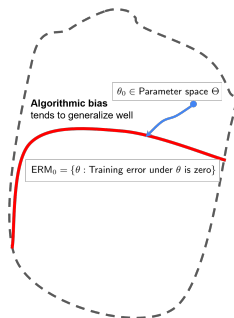


Literatures

- Surprises in high-dimensional ridgeless least squares interpolation
Hastie et al., 2018.
- The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime
Montanari et al., 2019.
- Just interpolate: Kernel “ridgeless” regression can generalize
Liang and Rakhlin, 2020.

When is it sufficient to overfit?

Benign overfitting

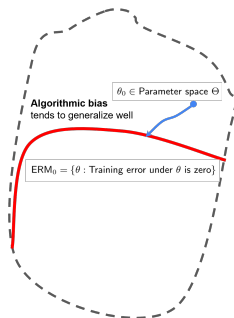


Literatures

- Surprises in high-dimensional ridgeless least squares interpolation
Hastie et al., 2018.
- The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime
Montanari et al., 2019.
- Just interpolate: Kernel “ridgeless” regression can generalize
Liang and Rakhlin, 2020.
- Two models of double descent for weak features
Belkin et al., 2020.

When is it sufficient to overfit?

Benign overfitting



Literatures

- Surprises in high-dimensional ridgeless least squares interpolation
Hastie et al., 2018.
- The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime
Montanari et al., 2019.
- Just interpolate: Kernel “ridgeless” regression can generalize
Liang and Rakhlin, 2020.
- Two models of double descent for weak features
Belkin et al., 2020.
- ...

When is it necessary to overfit?

When is it necessary to overfit?

Competing considerations

When is it necessary to overfit?

Competing considerations

- Privacy and security.

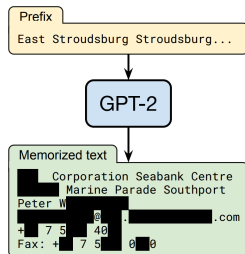


Figure from Carlini et al., 2021

When is it necessary to overfit?

Competing considerations

- Privacy and security.

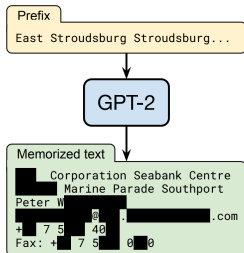


Figure from Carlini et al., 2021

- Can we generalize well without memorization?

When is it necessary to overfit?

Competing considerations

- Privacy and security.

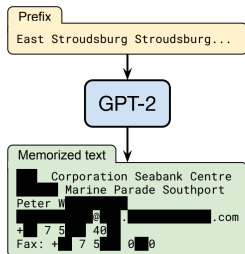


Figure from Carlini et al., 2021

- Can we generalize well without memorization?

Inspiring line of works

- Does Learning Require Memorization? A Short Tale about a Long Tail
Feldman, 2019.
- When Is Memorization of Irrelevant Training Data Necessary for High-Accuracy Learning?
Brown, Bun, Feldman, Smith, Talwar, 2021.

When is it necessary to overfit?

Competing considerations

- Privacy and security.

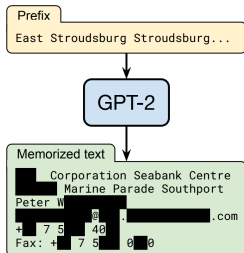


Figure from Carlini et al., 2021

- Can we generalize well without memorization?

Inspiring line of works

- Does Learning Require Memorization? A Short Tale about a Long Tail
Feldman, 2019.
- When Is Memorization of Irrelevant Training Data Necessary for High-Accuracy Learning?
Brown, Bun, Feldman, Smith, Talwar, 2021.

Takeaways

When is it necessary to overfit?

Competing considerations

- Privacy and security.

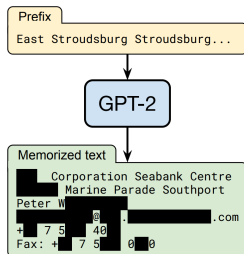


Figure from Carlini et al., 2021

- Can we generalize well without memorization?

Inspiring line of works

- Does Learning Require Memorization? A Short Tale about a Long Tail
Feldman, 2019.
- When Is Memorization of Irrelevant Training Data Necessary for High-Accuracy Learning?
Brown, Bun, Feldman, Smith, Talwar, 2021.

Takeaways

- Heavy-tailed distributions.

When is it necessary to overfit?

Competing considerations

- Privacy and security.

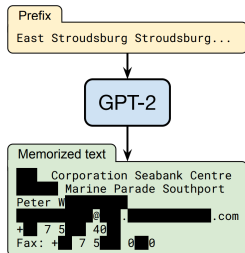


Figure from Carlini et al., 2021

- Can we generalize well without memorization?

Inspiring line of works

- Does Learning Require Memorization? A Short Tale about a Long Tail
Feldman, 2019.
- When Is Memorization of Irrelevant Training Data Necessary for High-Accuracy Learning?
Brown, Bun, Feldman, Smith, Talwar, 2021.

Takeaways

- Heavy-tailed distributions.
- Need to memorize each class.

When is it necessary to overfit?

Competing considerations

- Privacy and security.

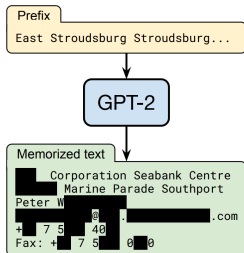


Figure from Carlini et al., 2021

- Can we generalize well without memorization?

Inspiring line of works

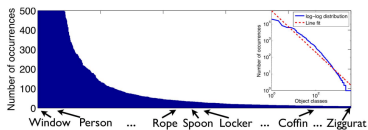
- Does Learning Require Memorization? A Short Tale about a Long Tail
Feldman, 2019.
- When Is Memorization of Irrelevant Training Data Necessary for High-Accuracy Learning?
Brown, Bun, Feldman, Smith, Talwar, 2021.

Takeaways

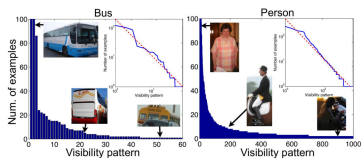
- Heavy-tailed distributions.
- Need to memorize each class.
- Combinatorial setup.

Real life data have heavy-tailed distributions

Real life data have heavy-tailed distributions

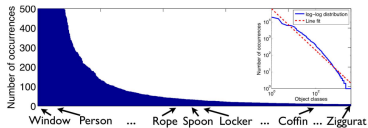


(a) The number of examples by object class in SUN dataset

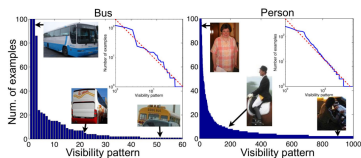


(b) Distributions of the visibility patterns for bus and person

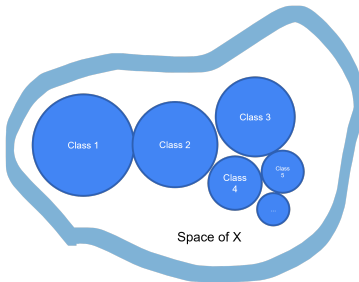
Real life data have heavy-tailed distributions



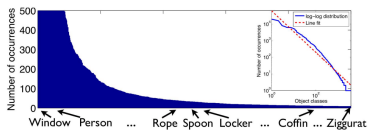
(a) The number of examples by object class in SUN dataset



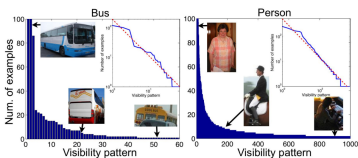
(b) Distributions of the visibility patterns for bus and person



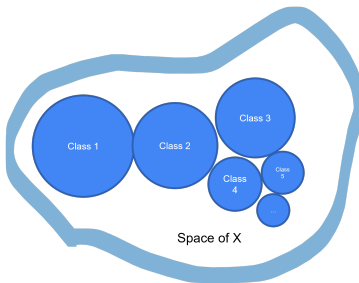
Real life data have heavy-tailed distributions



(a) The number of examples by object class in SUN dataset



(b) Distributions of the visibility patterns for bus and person



Have to memorize for each class

Carefully constructed combinatorial settings

Carefully constructed combinatorial settings

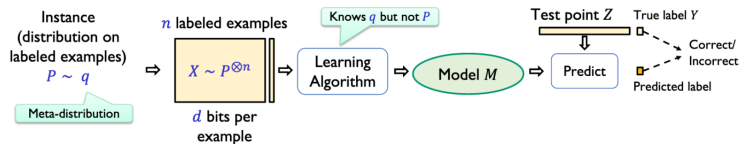


Figure from **Brown et al., 2021**

Carefully constructed combinatorial settings

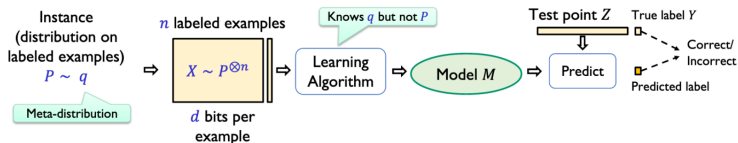
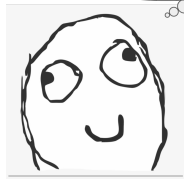


Figure from **Brown et al., 2021**

Can we do it for something simpler?
I don't know, like $y=X\theta+w$?



Average John

A general formulation

A general formulation

- Data pairs (x_i, y_i) from

$$y_i = f(x_i; \theta, w_i)$$

and a hypothesis class \mathcal{H} containing θ .

A general formulation

- Data pairs (x_i, y_i) from

$$y_i = f(x_i; \theta, w_i)$$

and a hypothesis class \mathcal{H}
containing θ .

- The cost of not-fitting

$$\begin{array}{ll} \underset{\hat{\theta} \in \mathcal{H}}{\text{minimize}} & \text{Pred}(\hat{\theta}) \\ \text{subject to} & \text{Train}(\hat{\theta}) \geq \epsilon^2 \end{array}$$

A general formulation

- Data pairs (x_i, y_i) from

$$y_i = f(x_i; \theta, w_i)$$

and a hypothesis class \mathcal{H} containing θ .

- The cost of not-fitting

$$\begin{array}{ll} \underset{\hat{\theta} \in \mathcal{H}}{\text{minimize}} & \text{Pred}(\hat{\theta}) \\ \text{subject to} & \text{Train}(\hat{\theta}) \geq \epsilon^2 \end{array}$$

A simpler model

A general formulation

- Data pairs (x_i, y_i) from

$$y_i = f(x_i; \theta, w_i)$$

and a hypothesis class \mathcal{H} containing θ .

- The cost of not-fitting

$$\begin{array}{ll} \underset{\hat{\theta} \in \mathcal{H}}{\text{minimize}} & \text{Pred}(\hat{\theta}) \\ \text{subject to} & \text{Train}(\hat{\theta}) \geq \epsilon^2 \end{array}$$

A simpler model

- Linear model for $X \in \mathbb{R}^{n \times d}$

$$y = X\theta + w$$

A general formulation

- Data pairs (x_i, y_i) from

$$y_i = f(x_i; \theta, w_i)$$

and a hypothesis class \mathcal{H} containing θ .

- The cost of not-fitting

$$\begin{array}{ll} \underset{\hat{\theta} \in \mathcal{H}}{\text{minimize}} & \text{Pred}(\hat{\theta}) \\ \text{subject to} & \text{Train}(\hat{\theta}) \geq \epsilon^2 \end{array}$$

A simpler model

- Linear model for $X \in \mathbb{R}^{n \times d}$

$$y = X\theta + w$$

- $d \geq n$ so we can interpolate

A general formulation

- Data pairs (x_i, y_i) from

$$y_i = f(x_i; \theta, w_i)$$

and a hypothesis class \mathcal{H} containing θ .

- The cost of not-fitting

$$\underset{\hat{\theta} \in \mathcal{H}}{\text{minimize}} \quad \text{Pred}(\hat{\theta})$$

$$\text{subject to} \quad \text{Train}(\hat{\theta}) \geq \epsilon^2$$

A simpler model

- Linear model for $X \in \mathbb{R}^{n \times d}$

$$y = X\theta + w$$

- $d \geq n$ so we can interpolate
- “memorization”: if we have to fit substantially below the inherent noise floor

Main results: necessity of memorization in linear regression

Let's start from the isotropic Gaussian case

Let's start from the isotropic Gaussian case

Problem setup

Consider the standard overparameterized linear model $y = X\theta + w$, with

Let's start from the isotropic Gaussian case

Problem setup

Consider the standard overparameterized linear model $y = X\theta + w$, with

- random isotropic i.i.d. design matrix $X = \mathbb{R}^{n \times d}$ ($d \geq n$)

Let's start from the isotropic Gaussian case

Problem setup

Consider the standard overparameterized linear model $y = X\theta + w$, with

- random isotropic i.i.d. design matrix $X = \mathbb{R}^{n \times d}$ ($d \geq n$)
- Bayesian setup, with unknown signal $\theta \sim \mathcal{N}(0, I_d/d)$

Let's start from the isotropic Gaussian case

Problem setup

Consider the standard overparameterized linear model $y = X\theta + w$, with

- random isotropic i.i.d. design matrix $X = \mathbb{R}^{n \times d}$ ($d \geq n$)
- Bayesian setup, with unknown signal $\theta \sim \mathcal{N}(0, I_d/d)$
- noise $w \sim \mathcal{N}(0, \sigma^2 I_n)$

Let's start from the isotropic Gaussian case

Problem setup

Consider the standard overparameterized linear model $y = X\theta + w$, with

- random isotropic i.i.d. design matrix $X = \mathbb{R}^{n \times d}$ ($d \geq n$)
- Bayesian setup, with unknown signal $\theta \sim \mathcal{N}(0, I_d/d)$
- noise $w \sim \mathcal{N}(0, \sigma^2 I_n)$

ℓ_2 error

$$\text{Train}_X(\hat{\theta}) = \frac{1}{n} \mathbb{E}_{w, \theta} \left[\left\| X\hat{\theta} - y \right\|_2^2 \right]$$

$$\text{Pred}_X(\hat{\theta}) = \mathbb{E}_{x, w, \theta} \left[\left\| x^\top \theta - x^\top \hat{\theta} \right\|_2^2 \right]$$

Cost of not-fitting for linear regression

Cost of not-fitting for linear regression

We want to solve

$$\min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}) \quad \text{s.t.} \quad \text{Train}_X(\hat{\theta}) \geq \epsilon^2,$$

where $\mathcal{H} = \{\hat{\theta}(X, y) \text{ square integrable}\}$.

Cost of not-fitting for linear regression

We want to solve

$$\min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}) \quad \text{s.t.} \quad \text{Train}_X(\hat{\theta}) \geq \epsilon^2,$$

where $\mathcal{H} = \{\hat{\theta}(X, y) \text{ square integrable}\}$. Let $\mathcal{H}(\epsilon) = \{\text{Train}_X(\hat{\theta}) \geq \epsilon^2\}$ and

$$\text{Cost}_X(\epsilon) := \min_{\hat{\theta} \in \mathcal{H}(\epsilon)} \text{Pred}_X(\hat{\theta}) - \min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}).$$

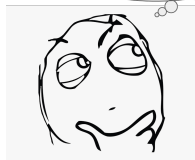
Cost of not-fitting for linear regression

We want to solve

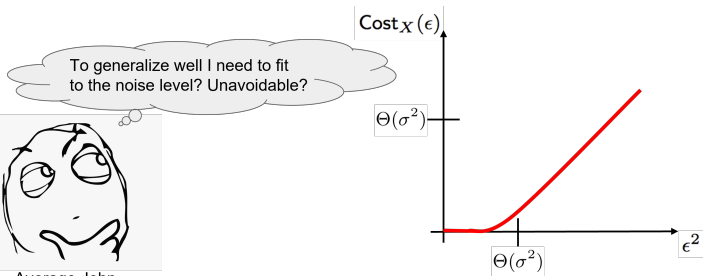
$$\min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}) \quad \text{s.t.} \quad \text{Train}_X(\hat{\theta}) \geq \epsilon^2,$$

where $\mathcal{H} = \{\hat{\theta}(X, y) \text{ square integrable}\}$. Let $\mathcal{H}(\epsilon) = \{\text{Train}_X(\hat{\theta}) \geq \epsilon^2\}$ and

$$\text{Cost}_X(\epsilon) := \min_{\hat{\theta} \in \mathcal{H}(\epsilon)} \text{Pred}_X(\hat{\theta}) - \min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}).$$



Average John



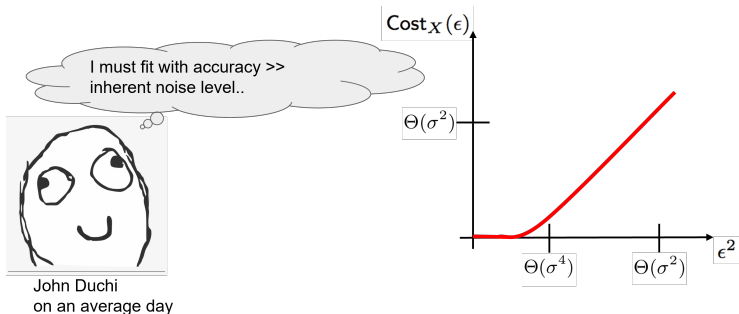
Cost of not-fitting for linear regression

We want to solve

$$\min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}) \quad \text{s.t.} \quad \text{Train}_X(\hat{\theta}) \geq \epsilon^2,$$

where $\mathcal{H} = \{\hat{\theta}(X, y) \text{ square integrable}\}$. Let $\mathcal{H}(\epsilon) = \{\text{Train}_X(\hat{\theta}) \geq \epsilon^2\}$ and

$$\text{Cost}_X(\epsilon) := \min_{\hat{\theta} \in \mathcal{H}(\epsilon)} \text{Pred}_X(\hat{\theta}) - \min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}).$$

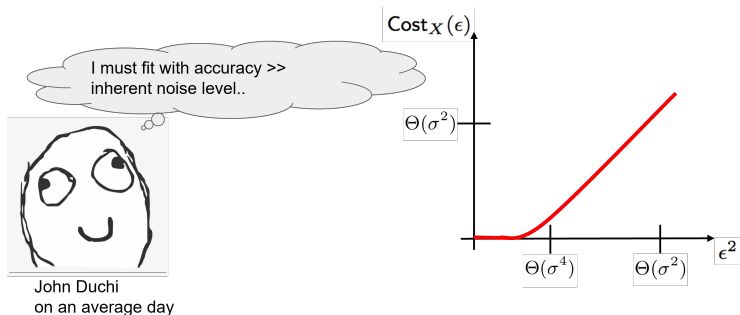


Cost of not-fitting for linear regression

Theorem 1 (Cheng, Duchi, Kuditipudi '22)

Under proportional asymptotics, namely $d/n \rightarrow \gamma$ as $n \rightarrow \infty$ for some $\gamma > 1$,

- (no-cost phase) $\lim_{n \rightarrow \infty} \text{Cost}_X(\epsilon) > 0$ iff $\epsilon^2 > \epsilon_\sigma^2 := \frac{\sigma^4}{\sigma^2 + 1 - 1/\gamma} + o(\sigma^4)$
- (linear-growth phase) $\lim_{n \rightarrow \infty} \text{Cost}_X(\epsilon) \geq C_\gamma \epsilon^2$ for $\epsilon^2 \geq c_\gamma \sigma^4$.

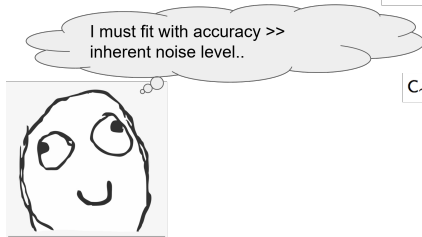


Cost of not-fitting for linear regression

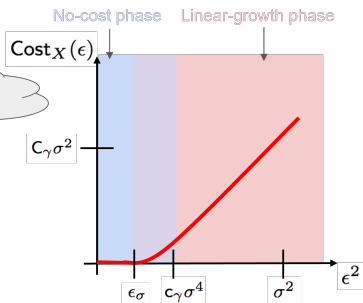
Theorem 1 (Cheng, Duchi, Kuditipudi '22)

Under proportional asymptotics $d/n \rightarrow \gamma > 1$,

- (no-cost phase) $\lim_{n \rightarrow \infty} \text{Cost}_X(\epsilon) > 0$ iff $\epsilon^2 > \epsilon_\sigma^2 := \frac{\sigma^4}{\sigma^2 + 1 - 1/\gamma} + o(\sigma^4)$
- (linear-growth phase) $\lim_{n \rightarrow \infty} \text{Cost}_X(\epsilon) \geq C_\gamma \epsilon^2$ for $\epsilon^2 \geq c_\gamma \sigma^4$.



John Duchi
on an average day



Proof sketch: strong duality and random matrix theory

The proof consists of three parts.

The proof consists of three parts.

- **Strong duality for linear estimators.** Starting from linear hypothesis class $\hat{\theta} = Ay$, we solve the nonconvex minimization problem

$$\min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}) \quad \text{s.t.} \quad \text{Train}_X(\hat{\theta}) \geq \epsilon^2,$$

using strong duality.

The proof consists of three parts.

- **Strong duality for linear estimators.** Starting from linear hypothesis class $\hat{\theta} = Ay$, we solve the nonconvex minimization problem

$$\min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}) \quad \text{s.t.} \quad \text{Train}_X(\hat{\theta}) \geq \epsilon^2,$$

using strong duality. Bayes optimal estimator is linear without constraint—without constraint? Still linear!

The proof consists of three parts.

- **Strong duality for linear estimators.** Starting from linear hypothesis class $\hat{\theta} = Ay$, we solve the nonconvex minimization problem

$$\min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}) \quad \text{s.t.} \quad \text{Train}_X(\hat{\theta}) \geq \epsilon^2,$$

using strong duality. Bayes optimal estimator is linear without constraint—without constraint? Still linear!

- **Derive asymptotics using RMT.** With the exact form of the (approximate) minimizer, we derive asymptotic limits of threshold value ϵ_σ , cost of not-fitting $\text{Cost}_X(\epsilon)$ by random matrix theory.

The proof consists of three parts.

- **Strong duality for linear estimators.** Starting from linear hypothesis class $\hat{\theta} = Ay$, we solve the nonconvex minimization problem

$$\min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}) \quad \text{s.t.} \quad \text{Train}_X(\hat{\theta}) \geq \epsilon^2,$$

using strong duality. Bayes optimal estimator is linear without constraint—without constraint? Still linear!

- **Derive asymptotics using RMT.** With the exact form of the (approximate) minimizer, we derive asymptotic limits of threshold value ϵ_σ , cost of not-fitting $\text{Cost}_X(\epsilon)$ by random matrix theory.
- **Upgrade by functional strong duality.** Finally, we upgrade to any square integrable estimator $\hat{\theta}(X, y)$ by showing a functional strong duality result.

Strong duality for linear hypothesis class

For linear estimator $\hat{\theta} = Ay$, let $\mathcal{P}(A) := \text{Pred}_X(\hat{\theta})$ and $\mathcal{T}(A) := \text{Train}_X(\hat{\theta})$.

For linear estimator $\hat{\theta} = Ay$, let $\mathcal{P}(A) := \text{Pred}_X(\hat{\theta})$ and $\mathcal{T}(A) := \text{Train}_X(\hat{\theta})$.

Reduction to QCQP

$$\begin{aligned} \underset{A \in \mathbb{R}^{d \times n}}{\text{minimize}} \quad & \mathcal{P}(A) = \frac{1}{d} \|AX - I\|_F^2 + \sigma^2 \|A\|_F^2 \\ \text{subject to} \quad & \mathcal{T}(A) = \frac{1}{nd} \|XAX - X\|_F^2 + \frac{\sigma^2}{n} \|XA - I\|_F^2 \geq \epsilon^2. \end{aligned}$$

Strong duality for linear hypothesis class

For linear estimator $\hat{\theta} = Ay$, let $\mathcal{P}(A) := \text{Pred}_X(\hat{\theta})$ and $\mathcal{T}(A) := \text{Train}_X(\hat{\theta})$.

Reduction to QCQP

$$\begin{aligned} \underset{A \in \mathbb{R}^{d \times n}}{\text{minimize}} \quad & \mathcal{P}(A) = \frac{1}{d} \|AX - I\|_F^2 + \sigma^2 \|A\|_F^2 \\ \text{subject to} \quad & \mathcal{T}(A) = \frac{1}{nd} \|XAX - X\|_F^2 + \frac{\sigma^2}{n} \|XA - I\|_F^2 \geq \epsilon^2. \end{aligned}$$

Strong duality

- The problem—while nonconvex—has **quadratic objective and a single quadratic constraint**. Strong duality holds!

Strong duality for linear hypothesis class

For linear estimator $\hat{\theta} = Ay$, let $\mathcal{P}(A) := \text{Pred}_X(\hat{\theta})$ and $\mathcal{T}(A) := \text{Train}_X(\hat{\theta})$.

Reduction to QCQP

$$\begin{aligned} \underset{A \in \mathbb{R}^{d \times n}}{\text{minimize}} \quad & \mathcal{P}(A) = \frac{1}{d} \|AX - I\|_F^2 + \sigma^2 \|A\|_F^2 \\ \text{subject to} \quad & \mathcal{T}(A) = \frac{1}{nd} \|XAX - X\|_F^2 + \frac{\sigma^2}{n} \|XA - I\|_F^2 \geq \epsilon^2. \end{aligned}$$

Strong duality

- The problem—while nonconvex—has **quadratic objective and a single quadratic constraint**. Strong duality holds!
- Optimality condition with $\rho_n := \rho_n(\epsilon)$

$$A(\rho_n) = \left(I - \rho_n \sigma^2 \left(I - \frac{\rho_n}{d} X^\top X \right)^{-1} \right) (X^\top X + d\sigma^2 I)^{-1} X^\top$$

Strong duality for linear hypothesis class

For linear estimator $\hat{\theta} = Ay$, let $\mathcal{P}(A) := \text{Pred}_X(\hat{\theta})$ and $\mathcal{T}(A) := \text{Train}_X(\hat{\theta})$.

Reduction to QCQP

$$\begin{aligned} \underset{A \in \mathbb{R}^{d \times n}}{\text{minimize}} \quad & \mathcal{P}(A) = \frac{1}{d} \|AX - I\|_F^2 + \sigma^2 \|A\|_F^2 \\ \text{subject to} \quad & \mathcal{T}(A) = \frac{1}{nd} \|XAX - X\|_F^2 + \frac{\sigma^2}{n} \|XA - I\|_F^2 \geq \epsilon^2. \end{aligned}$$

Strong duality

- The problem—while nonconvex—has **quadratic objective and a single quadratic constraint**. Strong duality holds!
- Optimality condition with $\rho_n := \rho_n(\epsilon)$

$$A(\rho_n) = \left(I - \rho_n \sigma^2 \left(I - \frac{\rho_n}{d} X^\top X \right)^{-1} \right) (X^\top X + d\sigma^2 I)^{-1} X^\top$$

Ridge estimator when $\rho = 0$, optimal with ϵ_σ^2 training error.

Let X have singular values $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. The empirical spectral distribution of $\frac{1}{d}XX^\top$ is μ_n with its c.d.f. $H_n(s) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\lambda_i^2/d \leq s}$.

Let X have singular values $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. The empirical spectral distribution of $\frac{1}{d}XX^\top$ is μ_n with its c.d.f. $H_n(s) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\lambda_i^2/d \leq s}$.

Marchenko-Pastur law

$$\mu_n \Rightarrow \mu, H_n(s) \rightarrow H(s).$$

Let X have singular values $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. The empirical spectral distribution of $\frac{1}{d}XX^\top$ is μ_n with its c.d.f. $H_n(s) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\lambda_i^2/d \leq s}$.

Marchenko-Pastur law

$\mu_n \Rightarrow \mu, H_n(s) \rightarrow H(s)$.

$$dH(s) = \frac{\gamma}{2\pi} \frac{\sqrt{(\lambda_+ - s)(s - \lambda_-)}}{s} \mathbb{1}_{s \in [\lambda_-, \lambda_+]} ds,$$

with $\lambda_{\pm} := (1 \pm 1/\sqrt{\gamma})^2$.

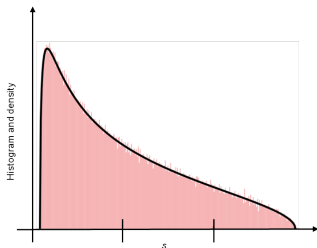
Let X have singular values $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. The empirical spectral distribution of $\frac{1}{d}XX^\top$ is μ_n with its c.d.f. $H_n(s) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\lambda_i^2/d \leq s}$.

Marchenko-Pastur law

$\mu_n \Rightarrow \mu, H_n(s) \rightarrow H(s)$.

$$dH(s) = \frac{\gamma}{2\pi} \frac{\sqrt{(\lambda_+ - s)(s - \lambda_-)}}{s} \mathbb{1}_{s \in [\lambda_-, \lambda_+]} ds,$$

with $\lambda_{\pm} := (1 \pm 1/\sqrt{\gamma})^2$.



Prediction and training errors

$$\mathcal{P}(A(\rho)) - \mathcal{P}(A(0)) = \frac{\rho^2 \sigma^4}{d} \text{Tr} \left(\left(I - \frac{\rho}{d} X^\top X \right)^{-2} \frac{X^\top X}{d} \left(\frac{X^\top X}{d} + \sigma^2 I \right)^{-1} \right)$$
$$\mathcal{T}(A(\rho)) = \frac{\sigma^4}{n} \text{Tr} \left(\left(I - \frac{\rho}{d} X^\top X \right)^{-2} \left(\frac{X^\top X}{d} + \sigma^2 I \right)^{-1} \right)$$

Prediction and training errors

$$\mathcal{P}(A(\rho)) - \mathcal{P}(A(0)) = \frac{\rho^2 \sigma^4}{d} \text{Tr} \left(\left(I - \frac{\rho}{d} X^\top X \right)^{-2} \frac{X^\top X}{d} \left(\frac{X^\top X}{d} + \sigma^2 I \right)^{-1} \right)$$

$$\mathcal{T}(A(\rho)) = \frac{\sigma^4}{n} \text{Tr} \left(\left(I - \frac{\rho}{d} X^\top X \right)^{-2} \left(\frac{X^\top X}{d} + \sigma^2 I \right)^{-1} \right)$$

Prediction and training errors in ESD

$$\mathcal{P}(A(\rho)) - \mathcal{P}(A(0)) = \frac{\rho^2 n}{d} \int \frac{\sigma^4 s}{(1 - \rho s)^2 (s + \sigma^2)} dH_n(s)$$

$$\mathcal{T}(A(\rho)) = \int \frac{\sigma^4}{(1 - \rho s)^2 (s + \sigma^2)} dH_n(s)$$

Prediction and training errors

$$\mathcal{P}(A(\rho)) - \mathcal{P}(A(0)) = \frac{\rho^2 \sigma^4}{d} \text{Tr} \left(\left(I - \frac{\rho}{d} X^\top X \right)^{-2} \frac{X^\top X}{d} \left(\frac{X^\top X}{d} + \sigma^2 I \right)^{-1} \right)$$

$$\mathcal{T}(A(\rho)) = \frac{\sigma^4}{n} \text{Tr} \left(\left(I - \frac{\rho}{d} X^\top X \right)^{-2} \left(\frac{X^\top X}{d} + \sigma^2 I \right)^{-1} \right)$$

Prediction and training errors in ESD

$$\mathcal{P}(A(\rho)) - \mathcal{P}(A(0)) = \frac{\rho^2 n}{d} \int \frac{\sigma^4 s}{(1 - \rho s)^2 (s + \sigma^2)} dH_n(s)$$

$$\mathcal{T}(A(\rho)) = \int \frac{\sigma^4}{(1 - \rho s)^2 (s + \sigma^2)} dH_n(s)$$

Limit of Lagrange multiplier

Since $\mathcal{T}(A(\rho_n)) = \epsilon^2$, would expect $\rho_n \rightarrow \rho_\epsilon$

$$\int \frac{\sigma^4}{(1 - \rho_\epsilon s)^2 (s + \sigma^2)} dH(s) = \epsilon^2$$

Limit of threshold

Taking $\rho_\epsilon = 0$ gives

$$\epsilon_\sigma^2 = \mathcal{T}(A(0)) \rightarrow \int \frac{\sigma^4}{s + \sigma^2} dH(s) = \frac{\sigma^4}{\sigma^2 + 1 - 1/\gamma} + o(\sigma^4)$$

Limit of threshold

Taking $\rho_\epsilon = 0$ gives

$$\epsilon_\sigma^2 = \mathcal{T}(A(0)) \rightarrow \int \frac{\sigma^4}{s + \sigma^2} dH(s) = \frac{\sigma^4}{\sigma^2 + 1 - 1/\gamma} + o(\sigma^4)$$

Limit of cost of not-fitting

$$\begin{aligned} \text{Cost}_X(\epsilon) &= \mathcal{P}(A(\rho_n)) - \mathcal{P}(A(0)) = \mathcal{P}(A(\rho_\epsilon)) - \mathcal{P}(A(0)) + o_n(1) \\ &\rightarrow \frac{\rho_\epsilon^2}{\gamma} \int \frac{\sigma^4 s}{(1 - \rho_\epsilon s)^2 (s + \sigma^2)} dH(s) \end{aligned}$$

Limit of threshold

Taking $\rho_\epsilon = 0$ gives

$$\epsilon_\sigma^2 = \mathcal{T}(A(0)) \rightarrow \int \frac{\sigma^4}{s + \sigma^2} dH(s) = \frac{\sigma^4}{\sigma^2 + 1 - 1/\gamma} + o(\sigma^4)$$

Limit of cost of not-fitting

$$\begin{aligned} \text{Cost}_X(\epsilon) &= \mathcal{P}(A(\rho_n)) - \mathcal{P}(A(0)) = \mathcal{P}(A(\rho_\epsilon)) - \mathcal{P}(A(0)) + o_n(1) \\ &\rightarrow \frac{\rho_\epsilon^2}{\gamma} \int \frac{\sigma^4 s}{(1 - \rho_\epsilon s)^2 (s + \sigma^2)} dH(s) \end{aligned}$$

Theorem 1 (Cheng, Duchi, Kuditipudi '22)

Under proportional asymptotics $d/n \rightarrow \gamma > 1$,

- (no-cost phase) $\lim_{n \rightarrow \infty} \text{Cost}_X(\epsilon) > 0$ iff $\epsilon^2 > \epsilon_\sigma^2 := \frac{\sigma^4}{\sigma^2 + 1 - 1/\gamma} + o(\sigma^4)$
- (linear-growth phase) $\lim_{n \rightarrow \infty} \text{Cost}_X(\epsilon) \geq C_\gamma \epsilon^2$ for $\epsilon^2 \geq c_\gamma \sigma^4$.

Upgrade to general hypothesis class

It only remains to show the same conclusion holds for $\hat{\theta}(X, y)$ square integrable given Gaussianity.

Upgrade to general hypothesis class

It only remains to show the same conclusion holds for $\hat{\theta}(X, y)$ square integrable given Gaussianity.

$$\begin{aligned} & \underset{\hat{\theta}(X, y) \in L^2}{\text{minimize}} && \int \left\| \hat{\theta} - \left(X^\top X + d\sigma^2 I \right)^{-1} X^\top y \right\|_2^2 d\mu \\ & \text{subject to} && \int \left\| X\hat{\theta} - y \right\|_2^2 d\mu \geq \epsilon^2 \end{aligned}$$

where $\mu \stackrel{d}{=} \mathcal{N}(0, \frac{1}{d} X X^\top + \sigma^2 I)$.

Upgrade to general hypothesis class

It only remains to show the same conclusion holds for $\hat{\theta}(X, y)$ square integrable given Gaussianity.

$$\begin{aligned} & \underset{\hat{\theta}(X, y) \in L^2}{\text{minimize}} && \int \left\| \hat{\theta} - \left(X^\top X + d\sigma^2 I \right)^{-1} X^\top y \right\|_2^2 d\mu \\ & \text{subject to} && \int \left\| X\hat{\theta} - y \right\|_2^2 d\mu \geq \epsilon^2 \end{aligned}$$

where $\mu \stackrel{d}{=} \mathcal{N}(0, \frac{1}{d} X X^\top + \sigma^2 I)$. Strong duality in Hilbert space?

Upgrade to general hypothesis class

It only remains to show the same conclusion holds for $\hat{\theta}(X, y)$ square integrable given Gaussianity.

$$\begin{aligned} & \underset{\hat{\theta}(X, y) \in L^2}{\text{minimize}} && \int \left\| \hat{\theta} - \left(X^\top X + d\sigma^2 I \right)^{-1} X^\top y \right\|_2^2 d\mu \\ & \text{subject to} && \int \left\| X\hat{\theta} - y \right\|_2^2 d\mu \geq \epsilon^2 \end{aligned}$$

where $\mu \stackrel{d}{=} \mathcal{N}(0, \frac{1}{d} X X^\top + \sigma^2 I)$. Strong duality in Hilbert space?

$$\hat{\theta} - \left(X^\top X + d\sigma^2 I \right)^{-1} X^\top y - \rho X^\top (X\hat{\theta} - y) / d = 0$$

Upgrade to general hypothesis class

It only remains to show the same conclusion holds for $\hat{\theta}(X, y)$ square integrable given Gaussianity.

$$\begin{aligned} & \underset{\hat{\theta}(X, y) \in L^2}{\text{minimize}} && \int \left\| \hat{\theta} - \left(X^\top X + d\sigma^2 I \right)^{-1} X^\top y \right\|_2^2 d\mu \\ & \text{subject to} && \int \left\| X\hat{\theta} - y \right\|_2^2 d\mu \geq \epsilon^2 \end{aligned}$$

where $\mu \stackrel{d}{=} N(0, \frac{1}{d} X X^\top + \sigma^2 I)$. Strong duality in Hilbert space?

$$\hat{\theta} - \left(X^\top X + d\sigma^2 I \right)^{-1} X^\top y - \rho X^\top (X\hat{\theta} - y)/d = 0$$

We exactly have

$$\hat{\theta} = A(X)y$$

Functional strong duality

Functional strong duality

$$\begin{aligned} & \underset{\hat{\theta}(X, y_i) \in \mathbb{R}^d, 1 \leq i \leq m}{\text{minimize}} && \int \left\| \hat{\theta} - (X^\top X + d\sigma^2 I)^{-1} X^\top y \right\|_2^2 d\mu_m \\ & \text{subject to} && \int \left\| X\hat{\theta} - y \right\|_2^2 d\mu_m \geq \epsilon^2 \end{aligned}$$

where μ_m are empirical distributions for i.i.d. samples of $y \mid X$.

Functional strong duality

$$\begin{aligned} & \underset{\hat{\theta}(X, y_i) \in \mathbb{R}^d, 1 \leq i \leq m}{\text{minimize}} && \int \left\| \hat{\theta} - (X^\top X + d\sigma^2 I)^{-1} X^\top y \right\|_2^2 d\mu_m \\ & \text{subject to} && \int \left\| X\hat{\theta} - y \right\|_2^2 d\mu_m \geq \epsilon^2 \end{aligned}$$

where μ_m are empirical distributions for i.i.d. samples of $y \mid X$. Strong duality applies to finite dimensional problems! Take $m \rightarrow \infty$ and conclude by SLLN.

Cost of not-interpolating

Cost of not-fitting

$$\text{Cost}_X(\epsilon) := \min_{\hat{\theta} \in \mathcal{H}(\epsilon)} \text{Pred}_X(\hat{\theta}) - \min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}).$$

Cost of not-fitting

$$\text{Cost}_X(\epsilon) := \min_{\hat{\theta} \in \mathcal{H}(\epsilon)} \text{Pred}_X(\hat{\theta}) - \min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}).$$

Cost of not-interpolating

$$\overline{\text{Cost}}_X(\epsilon) := \min_{\hat{\theta} \in \mathcal{H}(\epsilon)} \text{Pred}_X(\hat{\theta}) - \min_{\hat{\theta} \in \mathcal{H}_0} \text{Pred}_X(\hat{\theta}).$$

The optimal interpolant is the OLS estimator $\hat{\theta}_{\text{ols}} = X^\top (XX^\top)^{-1}y$.

Cost of not-fitting

$$\text{Cost}_X(\epsilon) := \min_{\hat{\theta} \in \mathcal{H}(\epsilon)} \text{Pred}_X(\hat{\theta}) - \min_{\hat{\theta} \in \mathcal{H}} \text{Pred}_X(\hat{\theta}).$$

Cost of not-interpolating

$$\overline{\text{Cost}}_X(\epsilon) := \min_{\hat{\theta} \in \mathcal{H}(\epsilon)} \text{Pred}_X(\hat{\theta}) - \min_{\hat{\theta} \in \mathcal{H}_0} \text{Pred}_X(\hat{\theta}).$$

The optimal interpolant is the OLS estimator $\hat{\theta}_{\text{ols}} = X^\top (XX^\top)^{-1}y$.

Theorem 2 (Cheng, Duchi, Kuditipudi '22)

Under proportional asymptotics $d/n \rightarrow \gamma > 1$,

- (no-cost phase) $\lim_{n \rightarrow \infty} \overline{\text{Cost}}_X(\epsilon) > 0$ iff $\epsilon^2 > \epsilon_{\sigma, \text{ols}}^2$.
- (linear-growth phase) $\lim_{n \rightarrow \infty} \overline{\text{Cost}}_X(\epsilon) \geq \bar{C}_\gamma \epsilon^2$ for $\epsilon^2 \geq \bar{c}_\gamma \sigma^4$.
- (threshold value) $\epsilon_\sigma < \epsilon_{\sigma, \text{ols}} \leq \kappa_\gamma \epsilon_\sigma$.

Relax assumptions

General covariance

General covariance

- The empirical spectral distribution of Σ converges.

General covariance

- The empirical spectral distribution of Σ converges.
- The condition number of Σ is bounded.

General covariance

- The empirical spectral distribution of Σ converges.
- The condition number of Σ is bounded.

General prior and noise distributions

General covariance

- The empirical spectral distribution of Σ converges.
- The condition number of Σ is bounded.

General prior and noise distributions

General covariance

- The empirical spectral distribution of Σ converges.
- The condition number of Σ is bounded.

General prior and noise distributions

- Gaussianity ensures model complexity. A counterexample when memorization does not happen is $\theta = e_j/\sqrt{d}$ with equal probability.

General covariance

- The empirical spectral distribution of Σ converges.
- The condition number of Σ is bounded.

General prior and noise distributions

- Gaussianity ensures model complexity. A counterexample when memorization does not happen is $\theta = e_j/\sqrt{d}$ with equal probability.
- For $\theta \sim (0, I_d/d)$ and $w \sim (0, \sigma^2 I_n)$, we restrict to linear estimators

$$\mathcal{H} = \left\{ \hat{\theta}(X, y) : \hat{\theta} = A(X)y \right\}.$$

General covariance

- The empirical spectral distribution of Σ converges.
- The condition number of Σ is bounded.

General prior and noise distributions

- Gaussianity ensures model complexity. A counterexample when memorization does not happen is $\theta = e_j/\sqrt{d}$ with equal probability.
- For $\theta \sim (0, I_d/d)$ and $w \sim (0, \sigma^2 I_n)$, we restrict to linear estimators

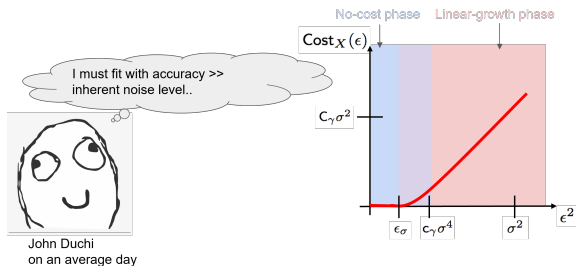
$$\mathcal{H} = \left\{ \hat{\theta}(X, y) : \hat{\theta} = A(X)y \right\}.$$

Theorem 3 (Cheng, Duchi, Kuditipudi '22)

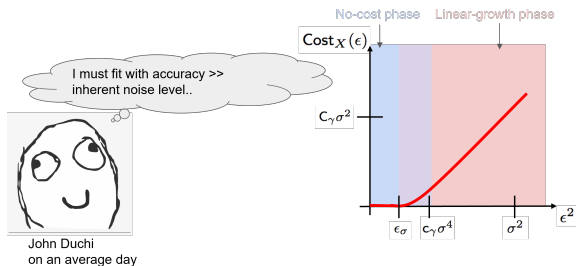
(Informal) *Under above conditions, we have to train till below $O(\sigma^4)$ error to generalize well.*

Concluding remarks

Necessity of memorization in linear regression

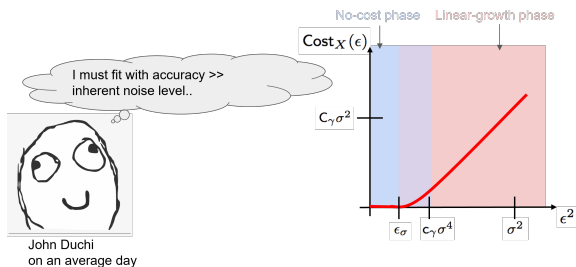


Necessity of memorization in linear regression



Similar results for other problems?

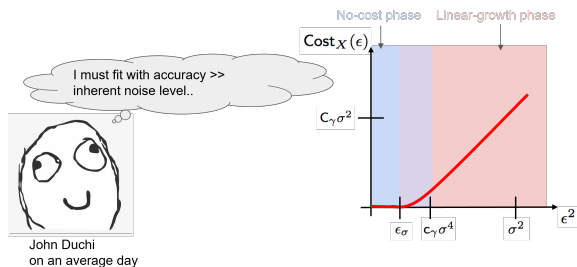
Necessity of memorization in linear regression



Similar results for other problems?

Implications for data cleaning and security. Can we have both?

Necessity of memorization in linear regression



Similar results for other problems?

Implications for data cleaning and security. Can we have both?

Motivation to construct datasets with multiple labels

- Theory of dataset with multiple labels. *Hilal Asi, Chen Cheng, John Duchi.*
- Surrogate consistency with data aggregation. *Chen Cheng, John Duchi.*

For more details: [arXiv:2202.09889](https://arxiv.org/abs/2202.09889)