

Asymptotics of the Sketched Pseudoinverse

December 2022

Daniel LeJeune

Rice University

Pratik Patil

Carnegie Mellon University

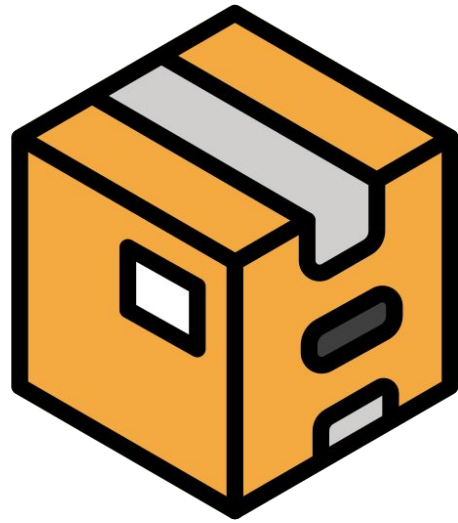
Department of Electrical and Computer Engineering Department of Statistics and Machine Learning

Joint work with: Hamid Javadi, Richard Baraniuk, Ryan Tibshirani

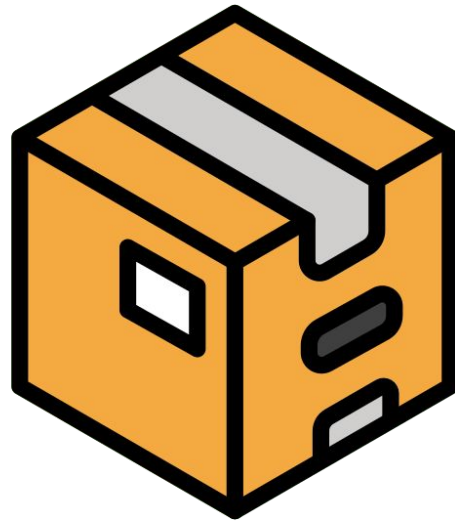
Too Much Data



Too Much Data



Too Much Data



Age=27, Height=5'11", ABO=A, ALT=36, Glucose=78, Creatinine=0.99, Sodium=132, Carbon Dioxide=21...



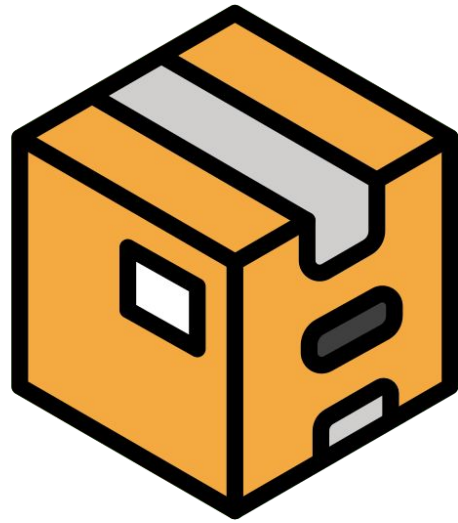
Age=56, Height=5'3", ABO=AB, ALT=40, Glucose=98, Creatinine=0.63, Sodium=182, Carbon Dioxide=25...



Age=41, Height=6'4", ABO=O, ALT=36, Glucose=84, Creatinine=0.79, Sodium=156, Carbon Dioxide=22...



Too Much Data



Age=27, Height=5'11", ABO=A, ALT=36, Glucose=78, Creatinine=0.99, Sodium=132, Carbon Dioxide=21...



Age=56, Height=5'3", ABO=AB, ALT=40, Glucose=98, Creatinine=0.63, Sodium=182, Carbon Dioxide=25...

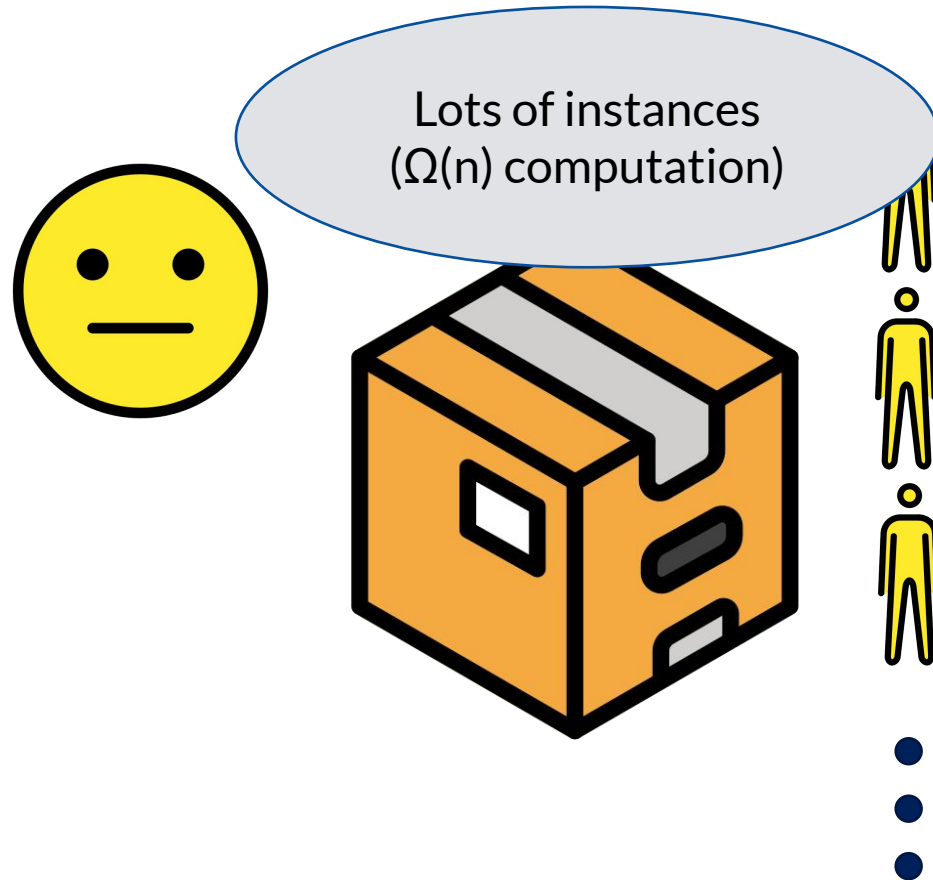


Age=41, Height=6'4", ABO=O, ALT=36, Glucose=84, Creatinine=0.79, Sodium=156, Carbon Dioxide=22...



Q: Odds of cancer in next 5 years?

Too Much Data



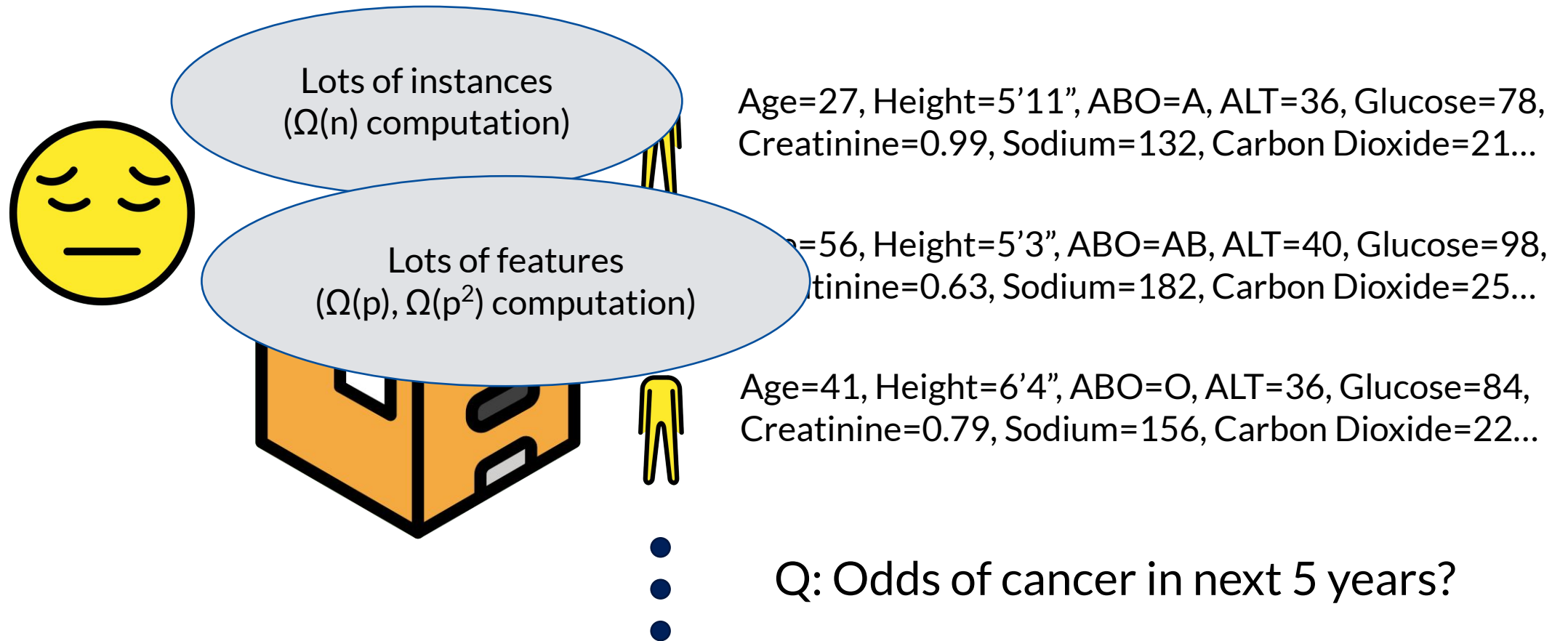
Age=27, Height=5'11", ABO=A, ALT=36, Glucose=78, Creatinine=0.99, Sodium=132, Carbon Dioxide=21...

Age=56, Height=5'3", ABO=AB, ALT=40, Glucose=98, Creatinine=0.63, Sodium=182, Carbon Dioxide=25...

Age=41, Height=6'4", ABO=O, ALT=36, Glucose=84, Creatinine=0.79, Sodium=156, Carbon Dioxide=22...

Q: Odds of cancer in next 5 years?

Too Much Data



Too Much Data



Lots of instances
($\Omega(n)$ computation)

Age=27, Height=5'11", ABO=A, ALT=36, Glucose=78,
Creatinine=0.99, Sodium=132, Carbon Dioxide=21...

Lots of features
($\Omega(p)$, $\Omega(p^2)$ computation)

Age=56, Height=5'3", ABO=AB, ALT=40, Glucose=98,
Creatinine=0.63, Sodium=182, Carbon Dioxide=25...

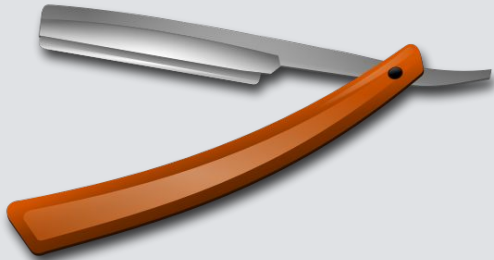
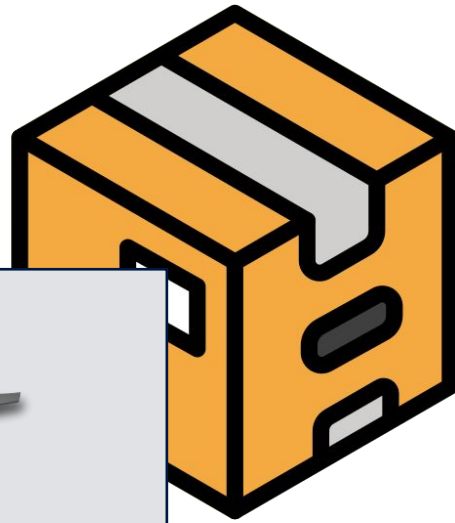
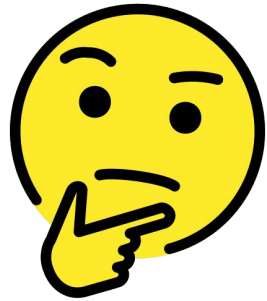
Curse of dimensionality
(Too many solutions)

Age=41, Height=6'4", ABO=O, ALT=36, Glucose=84,
Creatinine=0.79, Sodium=156, Carbon Dioxide=22...

Q: Odds of cancer in next 5 years?



Too Much Data



Occam's razor:
The true function should not be
too complicated.



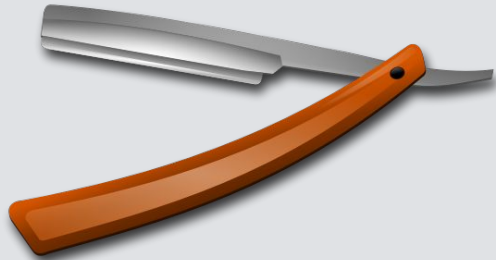
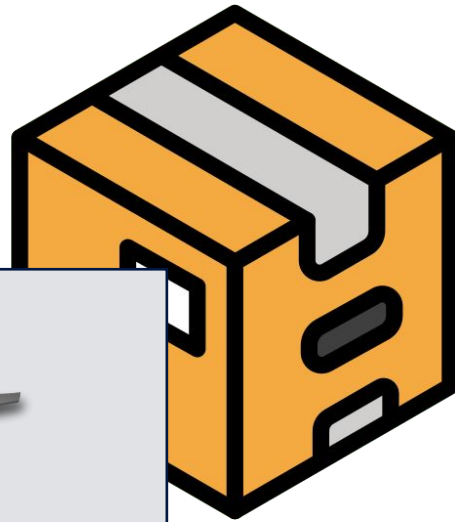
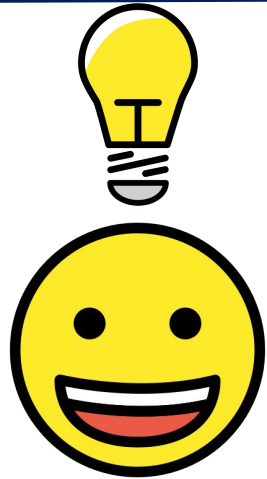
Age=27, Height=5'11", ABO=A, ALT=36, Glucose=78,
Creatinine=0.99, Sodium=132, Carbon Dioxide=21...

Age=56, Height=5'3", ABO=AB, ALT=40, Glucose=98,
Creatinine=0.63, Sodium=182, Carbon Dioxide=25...

Age=41, Height=6'4", ABO=O, ALT=36, Glucose=84,
Creatinine=0.79, Sodium=156, Carbon Dioxide=22...

Q: Odds of cancer in next 5 years?

Too Much Data



Occam's razor:
The true function should not be
too complicated.



Age=21, Height=5'10", ABO=O, ALT=36, Glucose=78, Creatinine=0.63, Sodium=156, Carbon Dioxide=21...

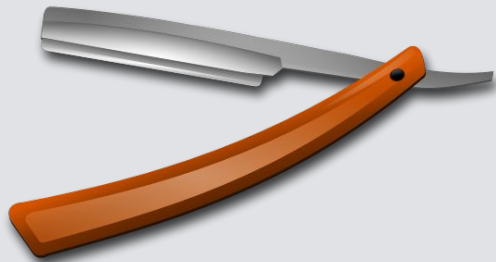
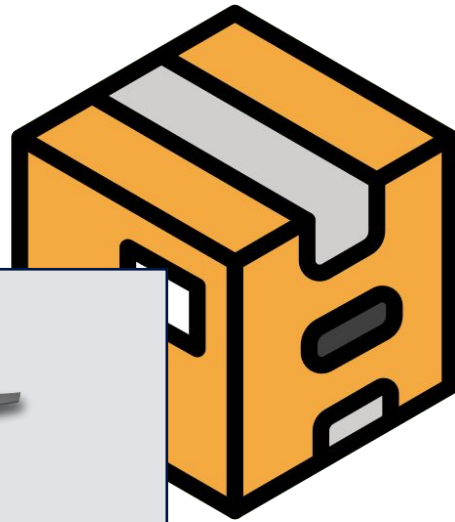
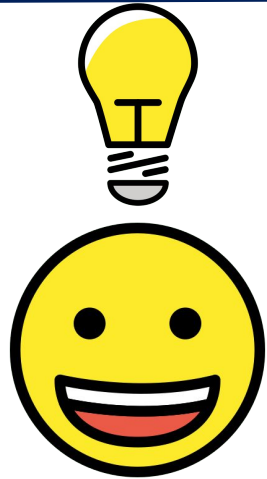
Age=56, Height=5'3", ABO=AB, ALT=40, Glucose=98, Creatinine=0.63, Sodium=182, Carbon Dioxide=25...

Age=41, Height=6'4", ABO=O, ALT=36, Glucose=84, Creatinine=0.79, Sodium=156, Carbon Dioxide=22...

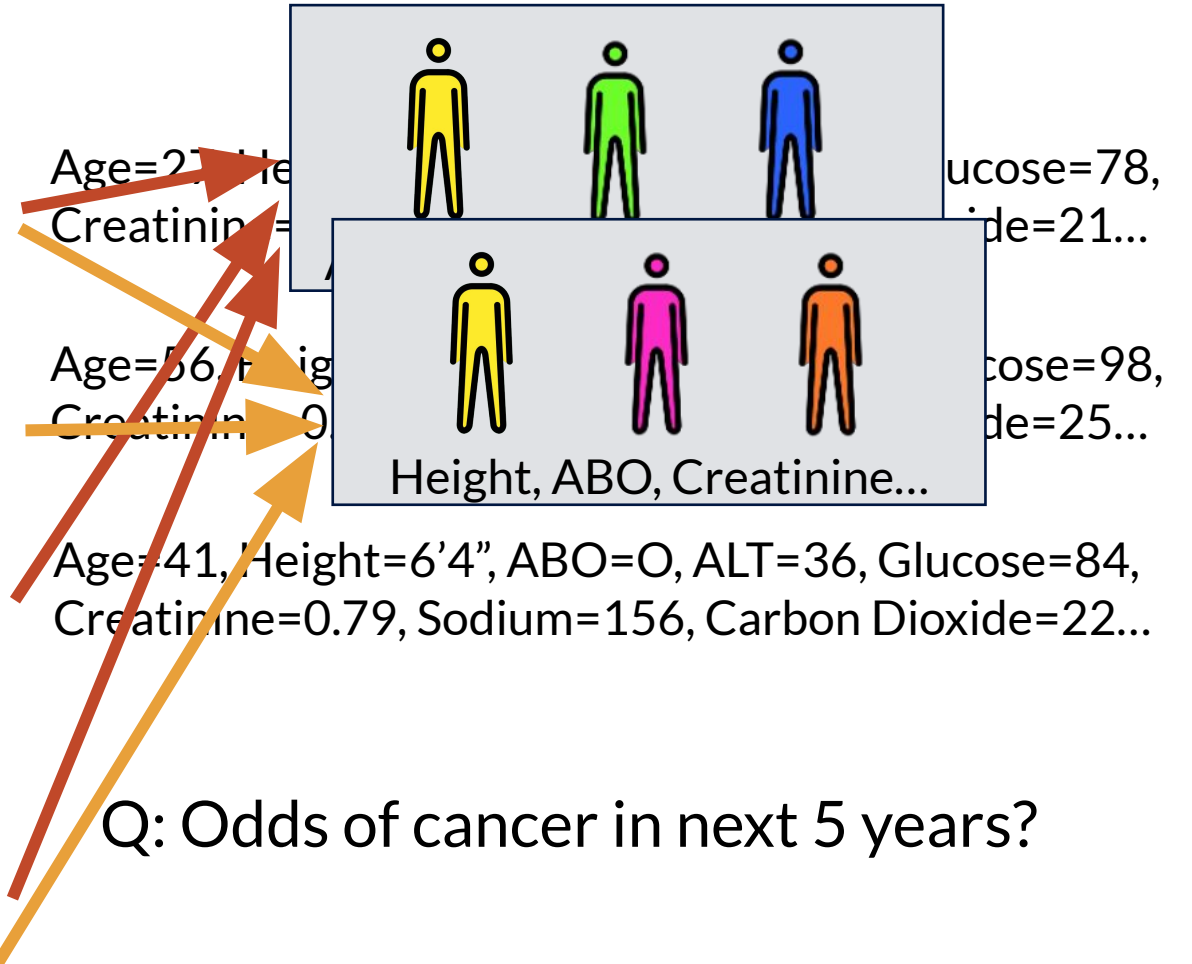
Q: Odds of cancer in next 5 years?

Age, ABO, Carbon Dioxide...

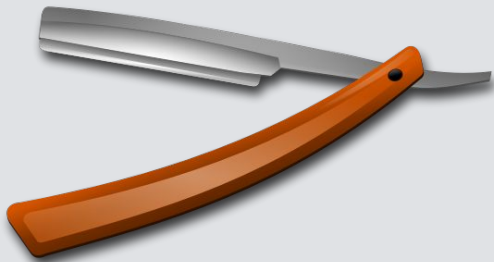
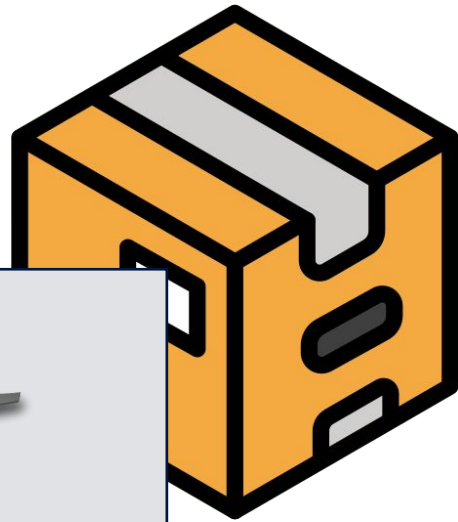
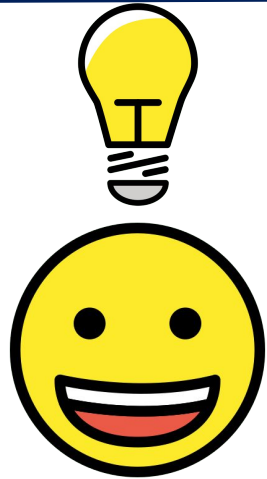
Too Much Data



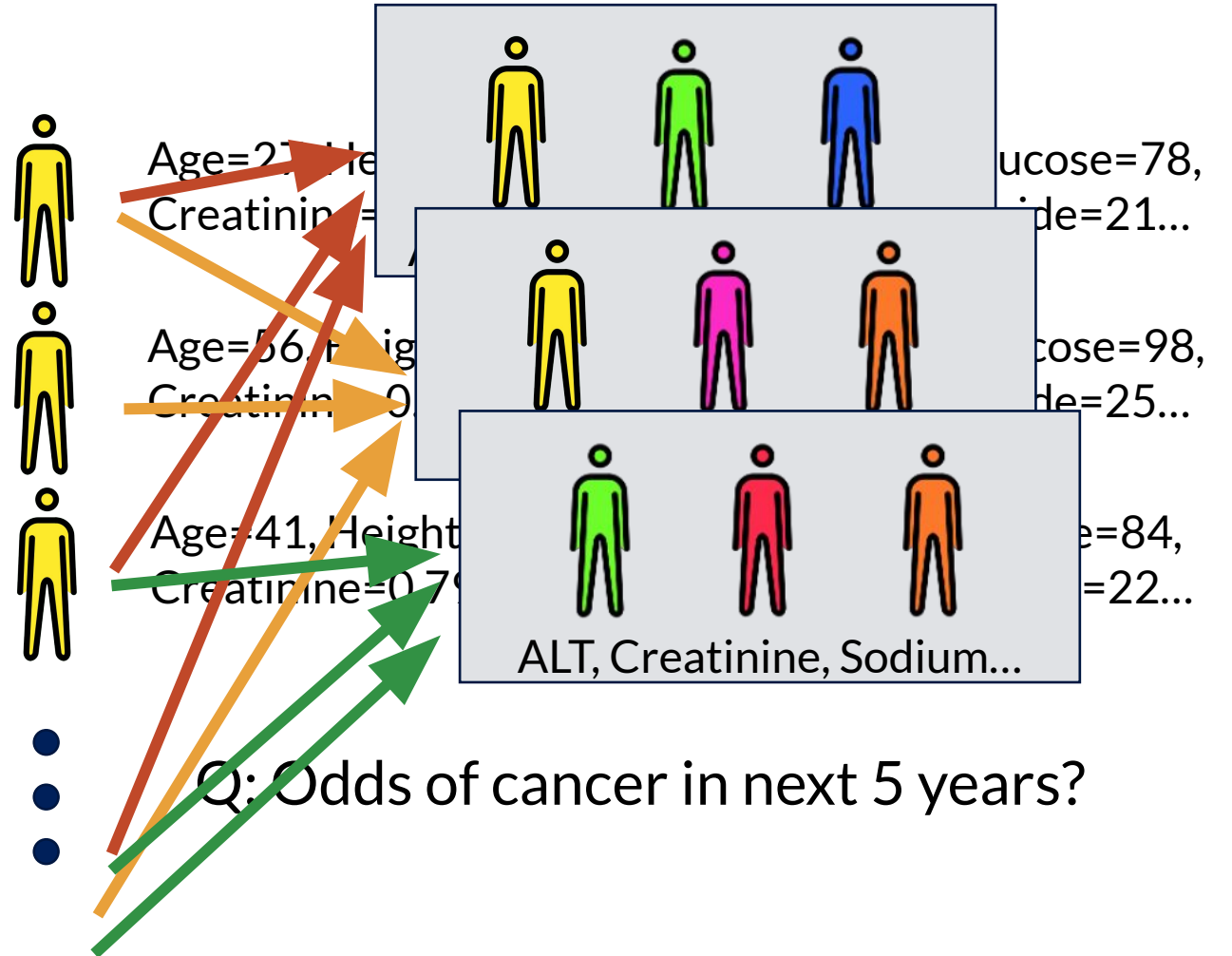
Occam's razor:
The true function should not be too complicated.



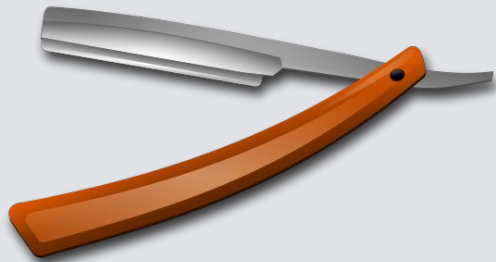
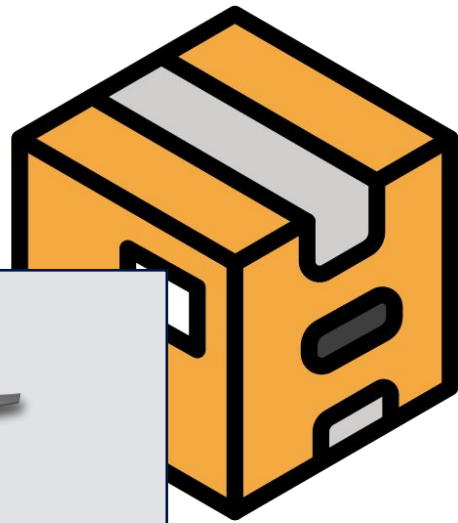
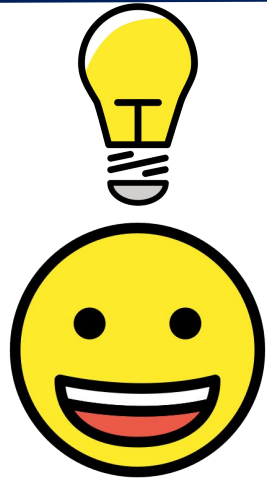
Too Much Data



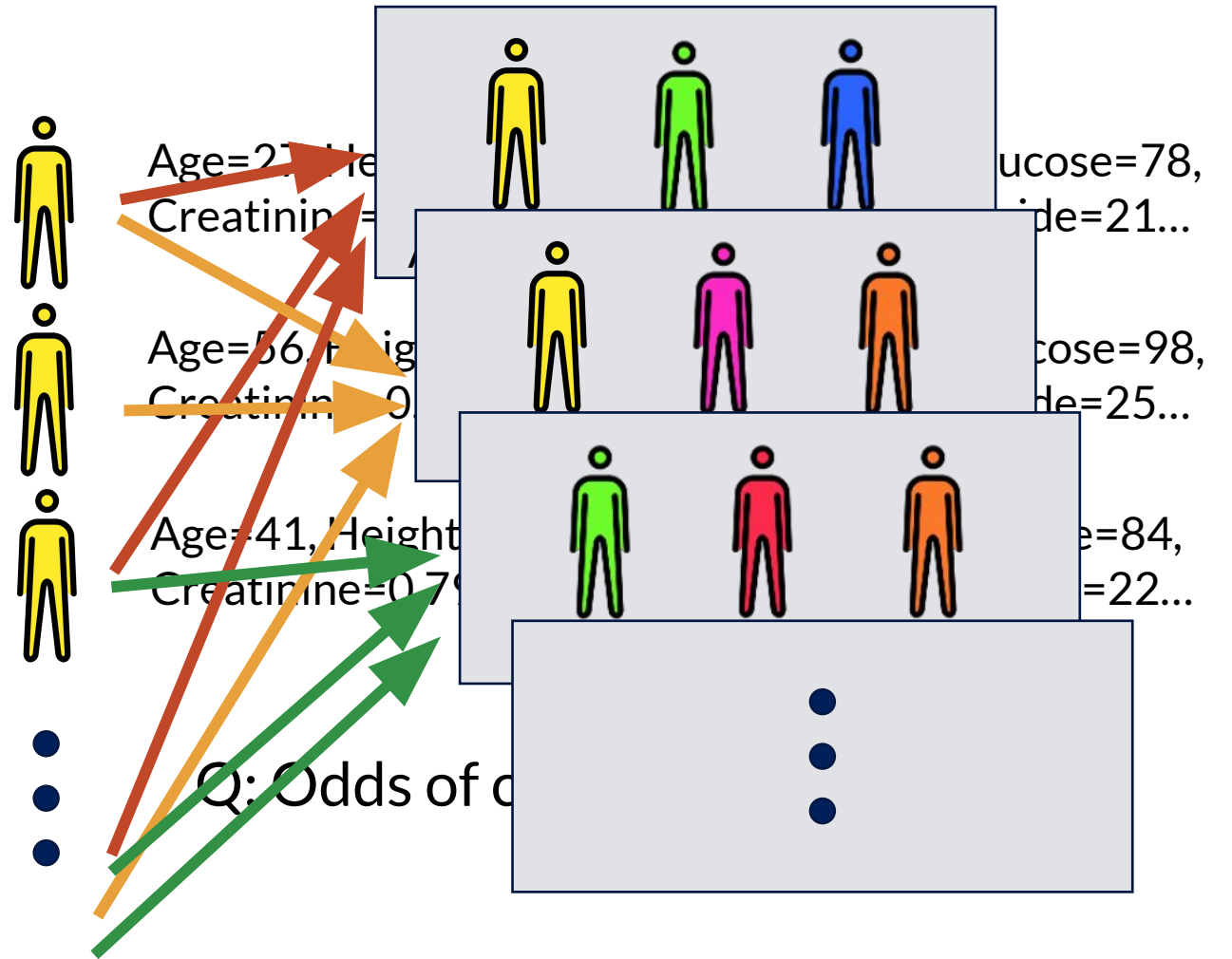
Occam's razor:
The true function should not be too complicated.



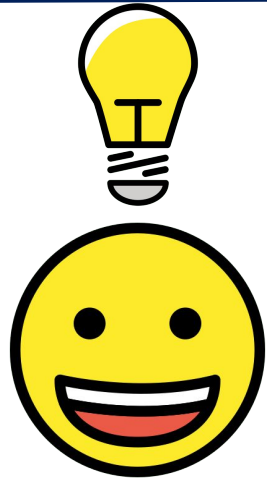
Too Much Data



Occam's razor:
The true function should not be too complicated.



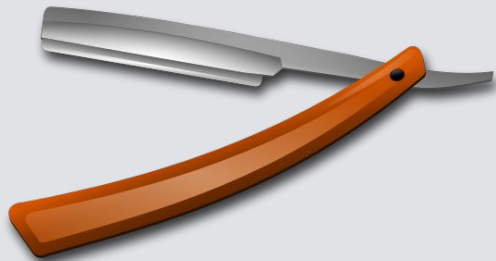
Too Much Data



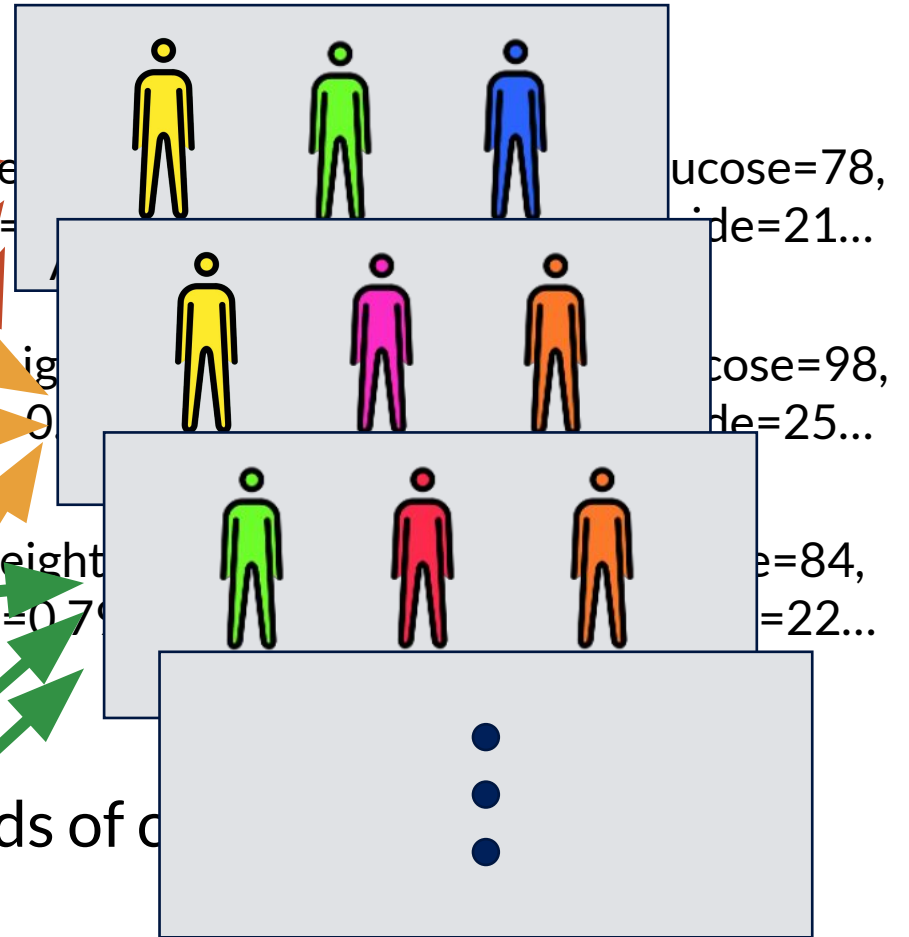
Few instances ✓

Few features ✓

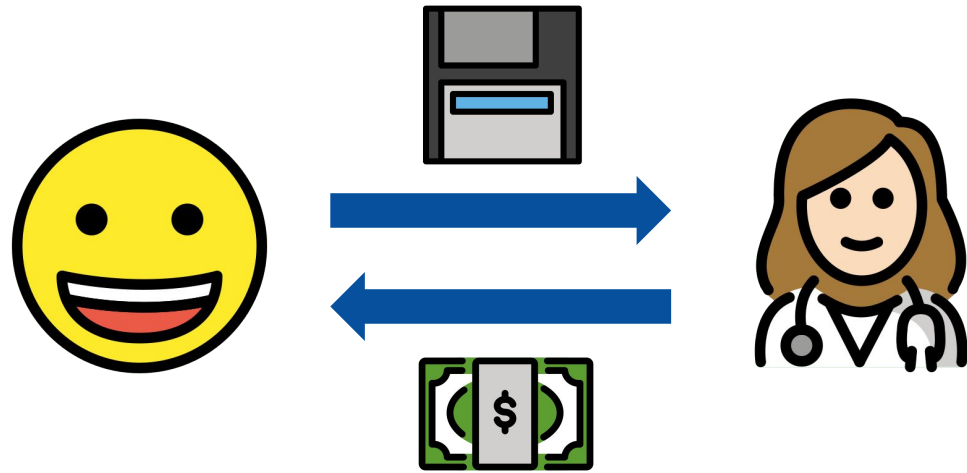
Simpler model ✓



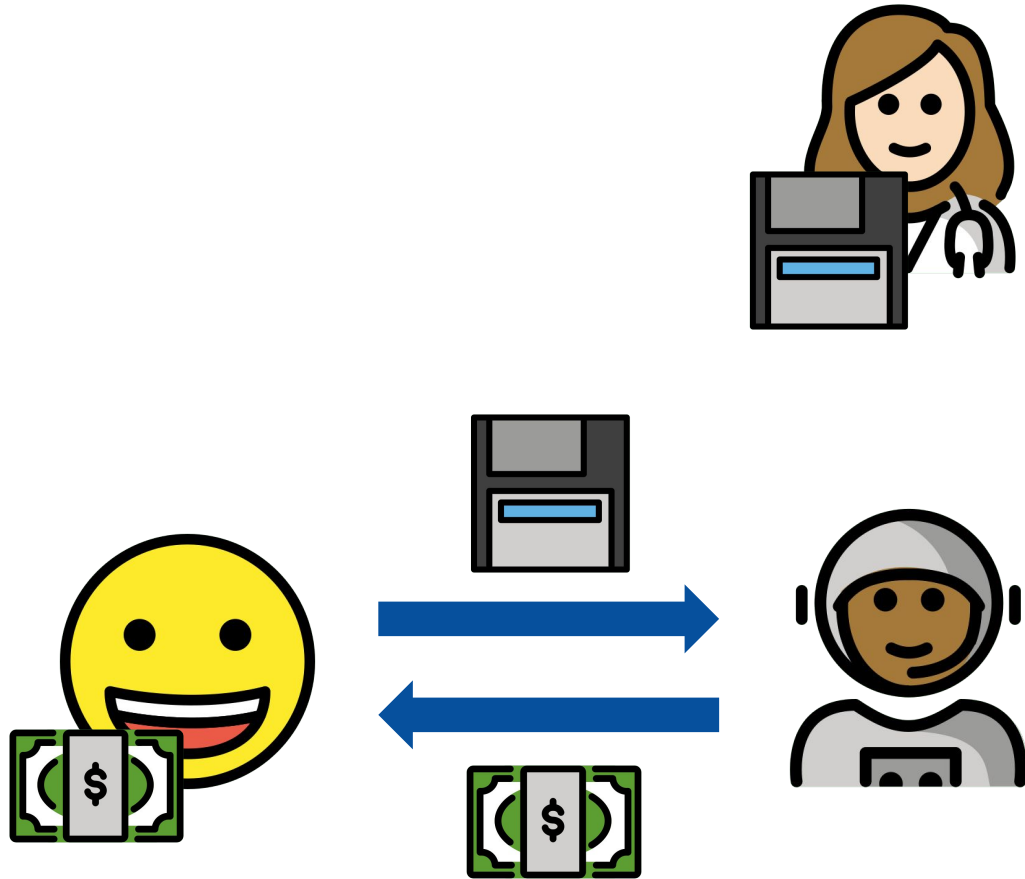
Occam's razor:
The true function should not be too complicated.



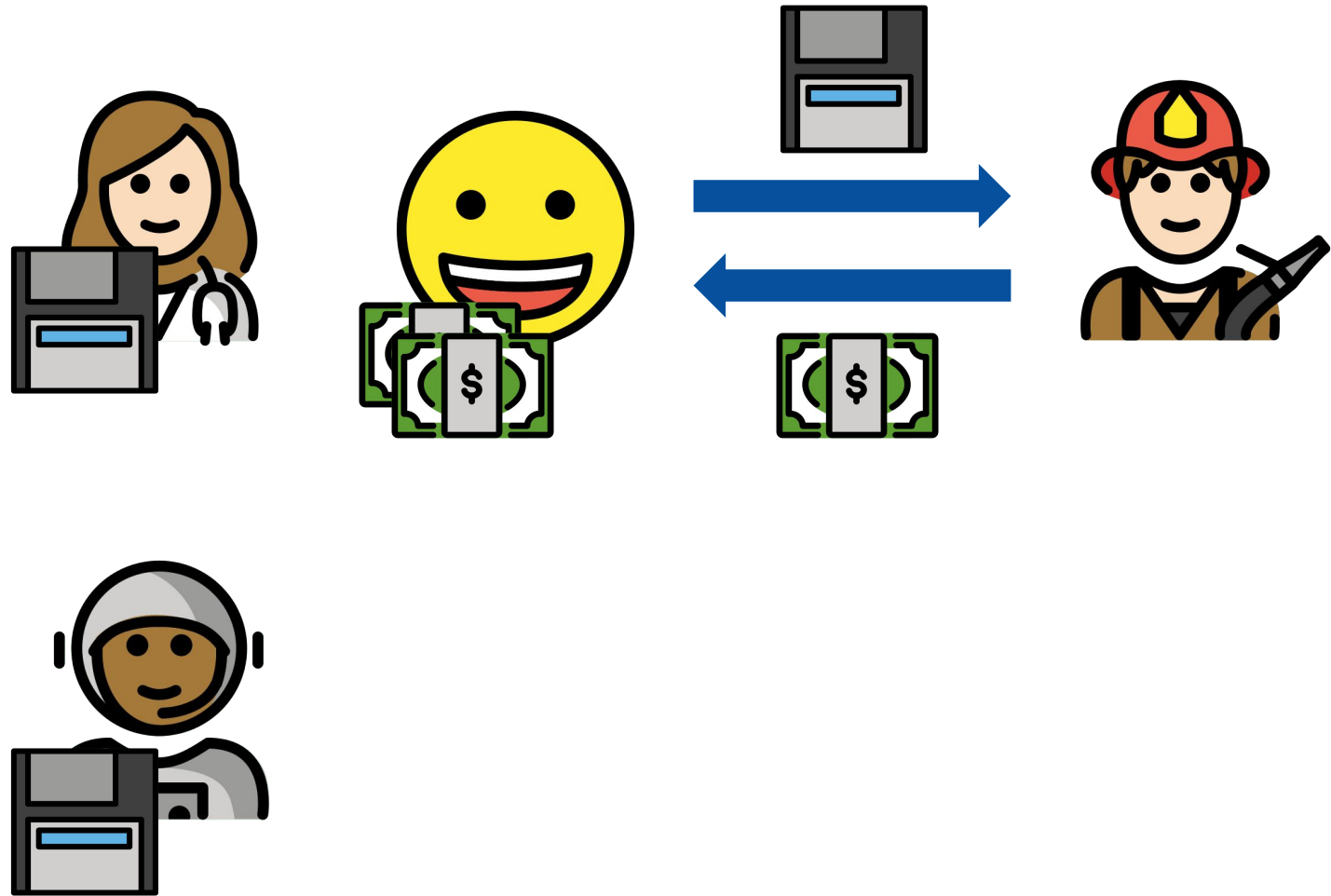
The State of Data Science



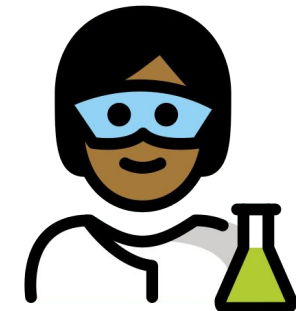
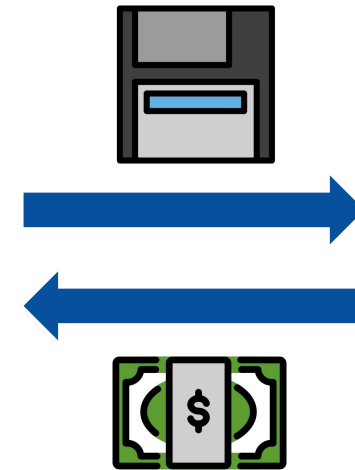
The State of Data Science



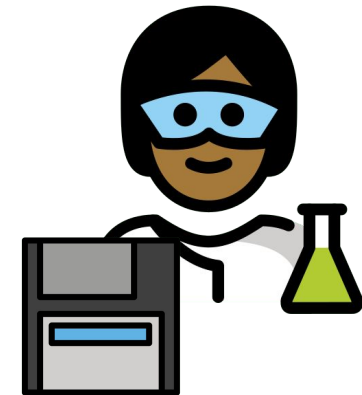
The State of Data Science



The State of Data Science

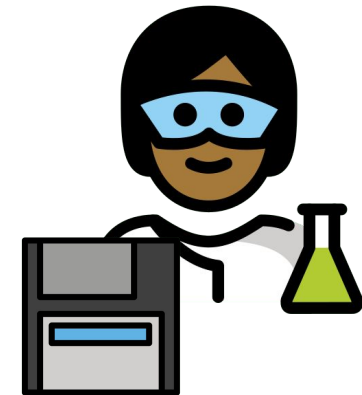
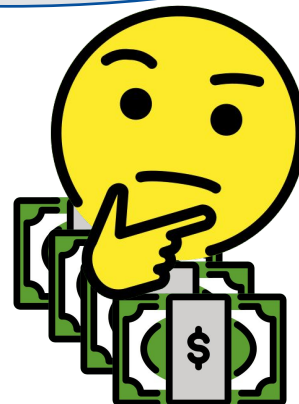


The State of Data Science



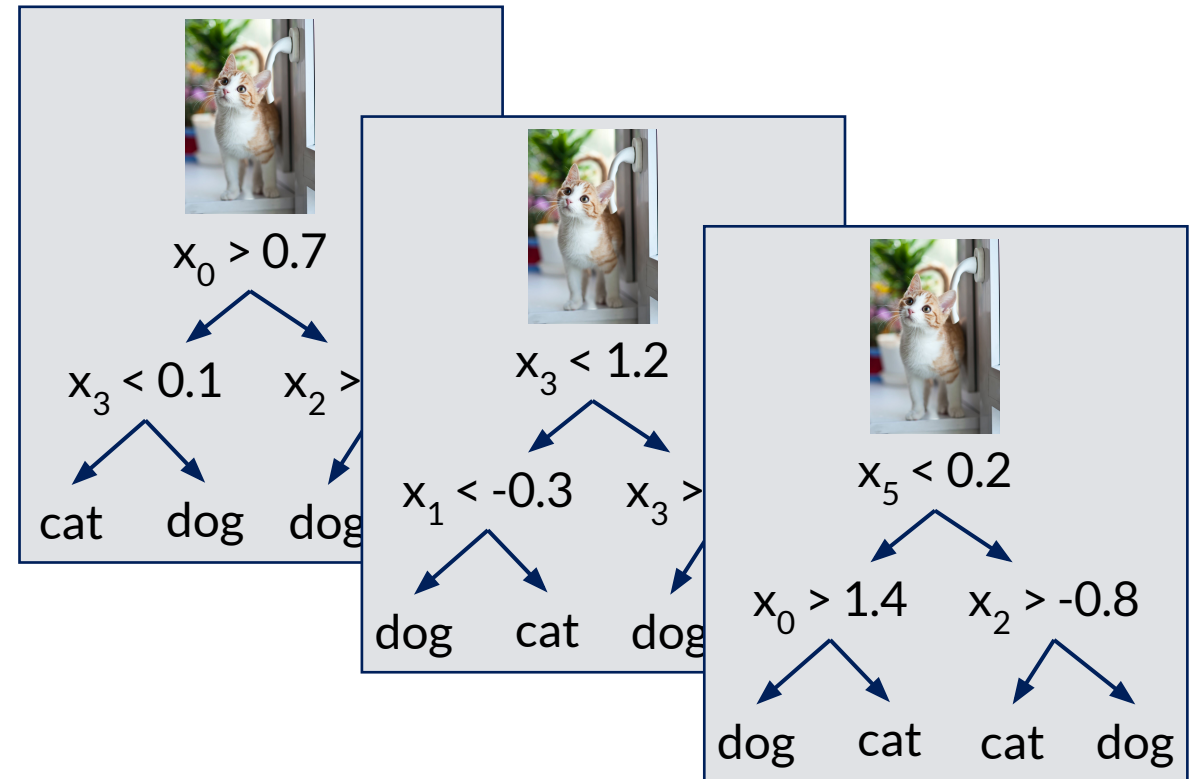
The State of Data Science

But what is this predictor?



Existing Practice

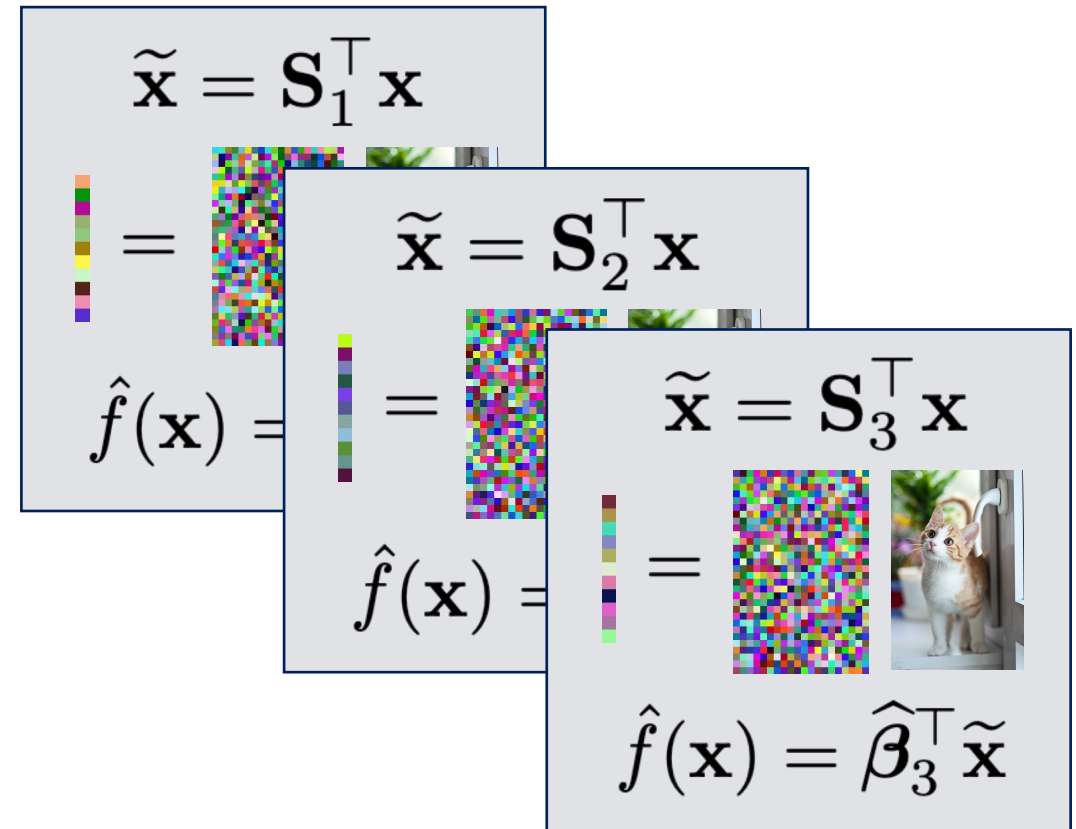
- Random forests [1]
 - Ensemble of **decision trees**
 - Trained on **bootstrap samples**
 - Branch on **random feature subsets**



[1] L Breiman. "Random forests." Machine Learning 45, 2001.

Existing Practice

- Random forests [1]
 - Ensemble of **decision trees**
 - Trained on **bootstrap samples**
 - Branch on **random feature subsets**
- Random projection ensembles [2]
 - Ensemble of **sketched regressors**
 - **Randomly** project **observations**
 - **Randomly** project **features**

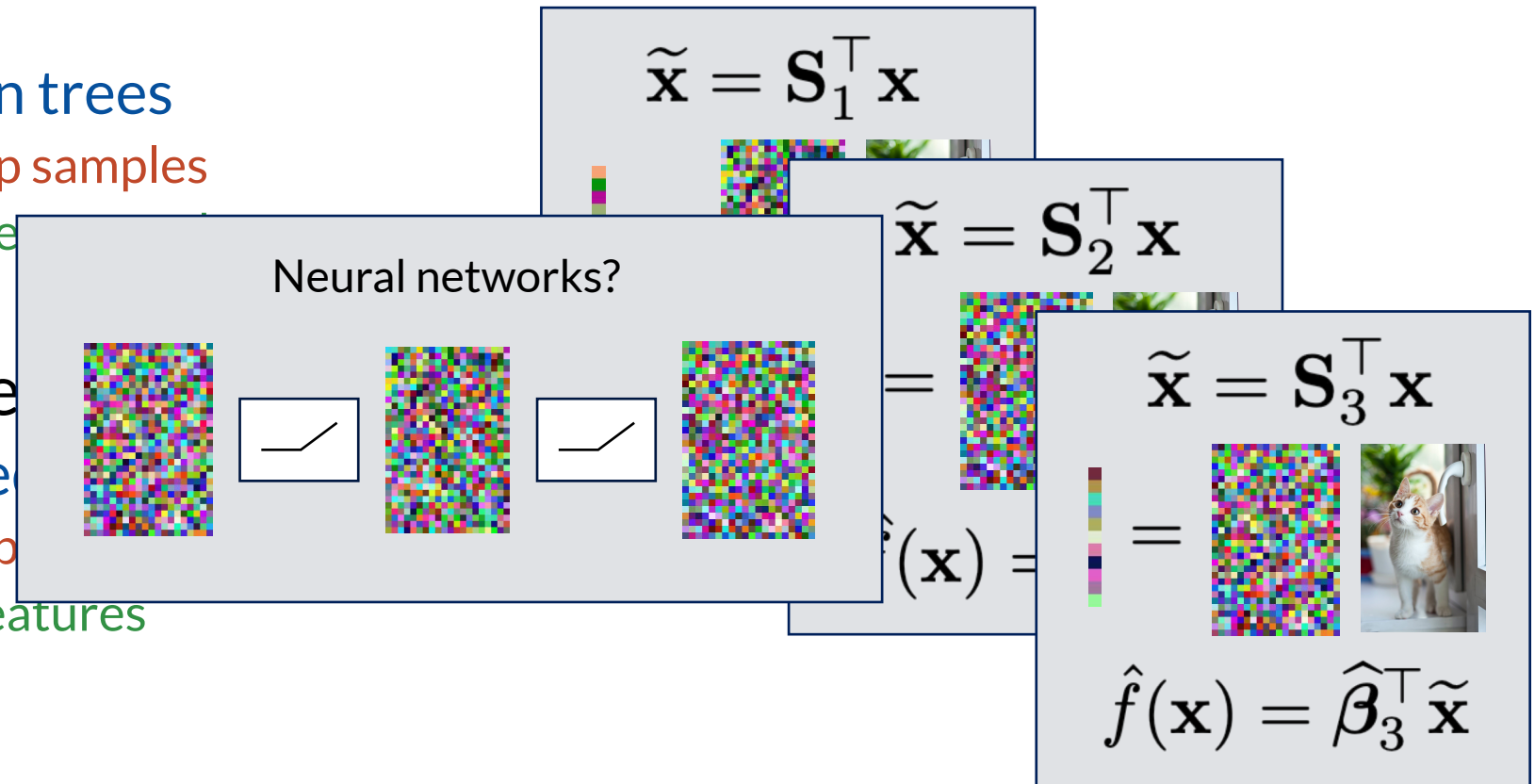


[1] L Breiman. "Random forests." Machine Learning 45, 2001.

[2] GA Thanei, C Heinze, N Meinshausen. "Random projections for large-scale regression." Big and Complex Data Analysis, 2017.

Existing Practice

- Random forests [1]
 - Ensemble of **decision trees**
 - Trained on **bootstrap samples**
 - Branch on **random features**
- Random projection ensemble
 - Ensemble of **sketches**
 - **Randomly** project objects
 - **Randomly** project features



[1] L Breiman. "Random forests." Machine Learning 45, 2001.

[2] GA Thanei, C Heinze, N Meinshausen. "Random projections for large-scale regression." Big and Complex Data Analysis, 2017.

Theoretical Questions

- **Prior work:** What **performance** do ensembles provably achieve?

Theoretical Questions

- **Prior work:** What **performance** do ensembles provably achieve?
 - Random forests are **difficult to analyze**
 - E.g., purely random forests [3] make analysis **tractable**
 - Simplified models [e.g., 2, 3] have proven **upper bounds**
 - **Consistent** estimation
 - Ensembles **strictly better** than individual models

[3] S Arlot, R Genuer. "Analysis of purely random forests bias." arXiv preprint arXiv:1407.3939.

Theoretical Questions

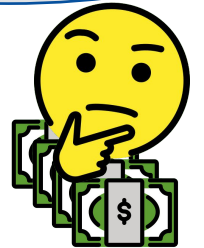
- **Prior work:** What **performance** do ensembles provably achieve?
 - Random forests are **difficult to analyze**
 - E.g., purely random forests [3] make analysis **tractable**
 - Simplified models [e.g., 2, 3] have proven **upper bounds**
 - **Consistent** estimation
 - Ensembles **strictly better** than individual models
- **New question:** What **predictions** do ensembles provably make?

[3] S Arlot, R Genuer. "Analysis of purely random forests bias." arXiv preprint arXiv:1407.3939.

Theoretical Questions

- **Prior work:** What **performance** do ensembles provably achieve?
 - Random forests are **difficult to analyze**
 - E.g., purely random forests [3] make analysis **tractable**
 - Simplified models [e.g., 2, 3] have proven **upper bounds**
 - **Consistent** estimation
 - Ensembles **strictly better** than individual models
- **New question:** What **predictions** do ensembles provably make?

But what is this predictor?

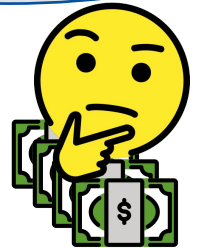


[3] S Arlot, R Genuer. "Analysis of purely random forests bias." arXiv preprint arXiv:1407.3939.

Theoretical Questions

- **Prior work:** What **performance** do ensembles provably achieve?
 - Random forests are **difficult to analyze**
 - E.g., purely random forests [3] make analysis **tractable**
 - Simplified models [e.g., 2, 3] have proven **upper bounds**
 - **Consistent** estimation
 - Ensembles **strictly better** than individual models
- **New question:** What **predictions** do ensembles provably make?
 - If there is **implicit regularization** in ensembles, can we make it **explicit**?
 - **Stronger** than consistency, **tighter** than upper bounds
 - Understand both **good/optimal** as well as **suboptimal** ensemble models

But what is this predictor?



[3] S Arlot, R Genuer. "Analysis of purely random forests bias." arXiv preprint arXiv:1407.3939.

Theoretical Questions

- **Prior work:** What **performance** do ensembles provably achieve?
 - Random forests are **difficult to analyze**
 - E.g., purely random forests [3] make analysis **tractable**
 - Simplified models [e.g., [1], [2]]
 - **Consistent** estimation
 - Ensembles **strictly** better
- **New question:** What **performance** can we provably make?
 - If there is **implicit regularization** in ensembles, can we make it **explicit**?
 - **Stronger** than consistency, **tighter** than upper bounds
 - Understand both **good/optimal** as well as **suboptimal** ensemble models

Spoiler:
Randomized least squares = ridge + noise
Randomized ensembles = ridge

But what is this predictor?



[3] S Arlot, R Genuer. "Analysis of purely random forests bias." arXiv preprint arXiv:1407.3939.

- Ground truth function $f: \mathcal{X} \rightarrow \mathbb{R}$
- Ensemble of estimators $\hat{f}_1, \dots, \hat{f}_K: \mathcal{X} \rightarrow \mathbb{R}$
- Squared error

$$R(\hat{f}_1, \dots, \hat{f}_K) \triangleq \mathbb{E} \left[R(\hat{f}_1, \dots, \hat{f}_K; x) \right]$$

$$R(\hat{f}_1, \dots, \hat{f}_K; x) \triangleq \mathbb{E} \left[\left(\frac{1}{K} \sum_{k=1}^K \hat{f}_k(x) - f(x) \right)^2 \right]$$

Early Work: Linear Regression Setting

- Training data: $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}$
 $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$

[4] DL, H Javadi, RG Baraniuk. “The implicit regularization of ordinary least squares ensembles.” AISTATS, 2020.

Early Work: Linear Regression Setting

- Training data: $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}$
 $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$
- Data model: $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$, $f(\mathbf{x}) = \boldsymbol{\beta}^{*\top} \mathbf{x}$, $\|\boldsymbol{\beta}^*\|_2 = 1$
 $y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$

[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

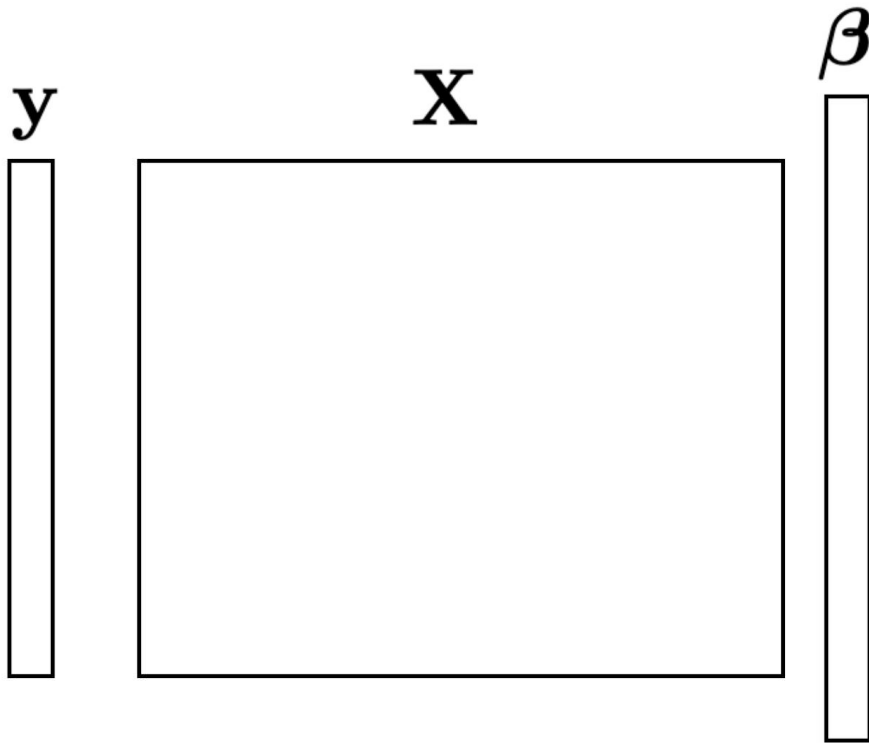
Early Work: Linear Regression Setting

- Training data: $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}$
 $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$
- Data model: $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$, $f(\mathbf{x}) = \boldsymbol{\beta}^{*\top} \mathbf{x}$, $\|\boldsymbol{\beta}^*\|_2 = 1$
 $y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$
- Estimators: $\hat{f}_k(\mathbf{x}) = \hat{\boldsymbol{\beta}}_k^\top \mathbf{x}$
 $\hat{\boldsymbol{\beta}}_k = \mathbf{S}_k \cdot \arg \min_{\mathbf{b}} \frac{1}{2} \|\mathbf{T}_k^\top (\mathbf{y} - \mathbf{X} \mathbf{S}_k \mathbf{b})\|_2^2$

[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

Random Subsampling

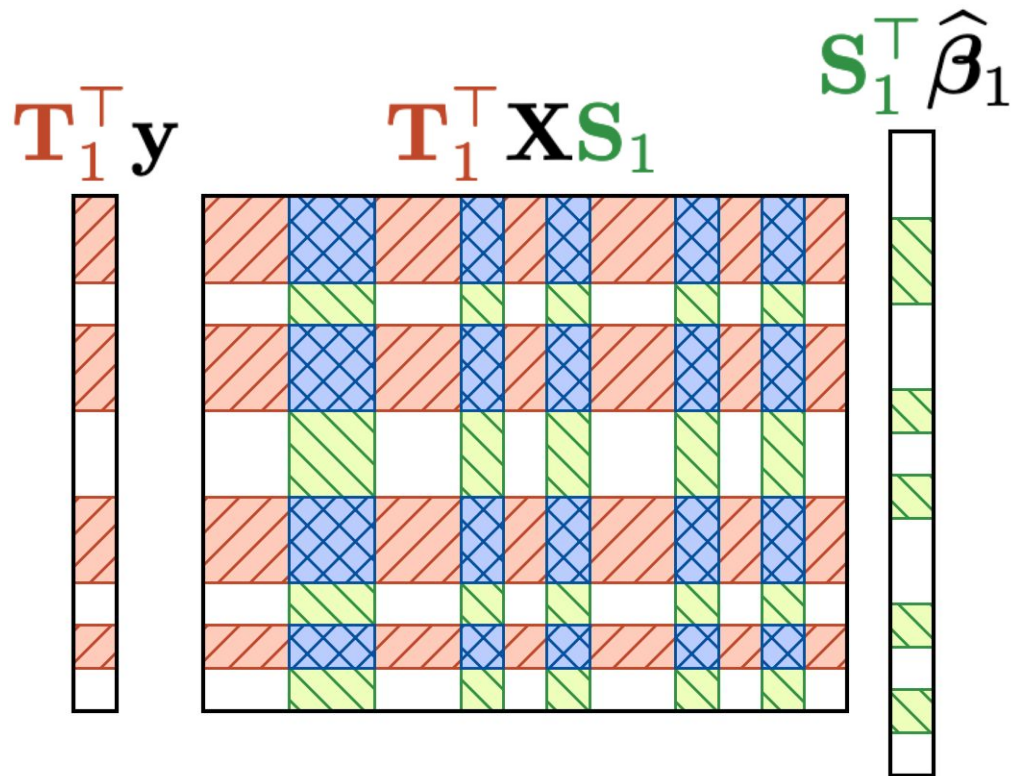
- Subsampling operators $\mathbf{S}_k \in \mathbb{R}^{p \times q}$, $\mathbf{T}_k \in \mathbb{R}^{n \times m}$, $m > q + 1$
 - Uniformly random columns of identity matrix



[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

Random Subsampling

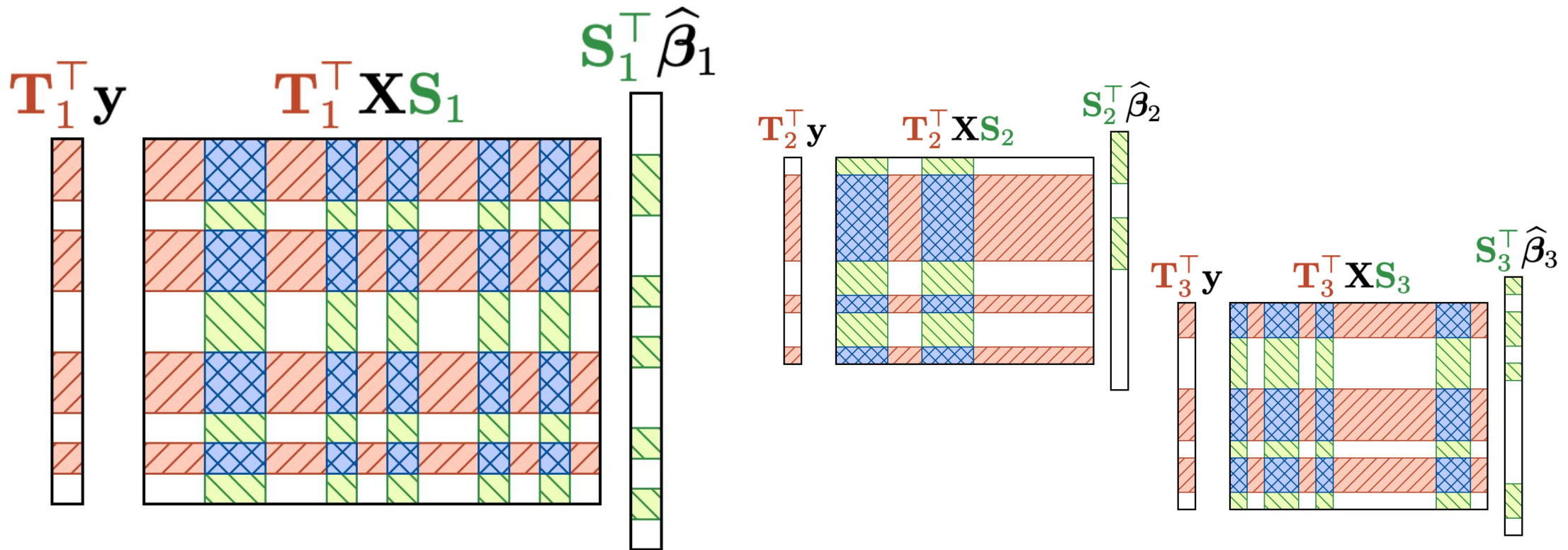
- Subsampling operators $\mathbf{S}_k \in \mathbb{R}^{p \times q}$, $\mathbf{T}_k \in \mathbb{R}^{n \times m}$, $m > q + 1$
 - Uniformly random columns of identity matrix



[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

Random Subsampling

- Subsampling operators $\mathbf{S}_k \in \mathbb{R}^{p \times q}$, $\mathbf{T}_k \in \mathbb{R}^{n \times m}$, $m > q + 1$
 - Uniformly random columns of identity matrix



[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

Turn the Crank...

$$R(\hat{f}_1, \dots, \hat{f}_K) = \frac{1}{K^2} \sum_{k, \ell=1}^K \mathbb{E} \left[\beta^{*\top} \left(\mathbf{I}_p - \mathbf{S}_k (\mathbf{T}_k^\top \mathbf{X} \mathbf{S}_k)^\dagger \mathbf{T}_k^\top \mathbf{X} \right)^\top \left(\mathbf{I}_p - \mathbf{S}_\ell (\mathbf{T}_\ell^\top \mathbf{X} \mathbf{S}_\ell)^\dagger \mathbf{T}_\ell^\top \mathbf{X} \right) \beta^* \right] \\ + \sigma^2 \mathbb{E} \left[\text{tr} \left[\mathbf{T}_k (\mathbf{S}_k^\top \mathbf{T}_k)^\dagger \mathbf{S}_k^\top \mathbf{S}_\ell (\mathbf{T}_\ell^\top \mathbf{X} \mathbf{S}_\ell)^\dagger \mathbf{T}_\ell^\top \right] \right]$$

Wishart matrix

Conditional expectation

Turn the Crank...

$$R(\hat{f}_1, \dots, \hat{f}_K) = \frac{1}{K^2} \sum_{k, \ell=1}^K \mathbb{E} \left[\beta^{*\top} \left(\mathbf{I}_p - \mathbf{S}_k (\mathbf{T}_k^\top \mathbf{X} \mathbf{S}_k)^\dagger \mathbf{T}_k^\top \mathbf{X} \right)^\top \left(\mathbf{I}_p - \mathbf{S}_\ell (\mathbf{T}_\ell^\top \mathbf{X} \mathbf{S}_\ell)^\dagger \mathbf{T}_\ell^\top \mathbf{X} \right) \beta^* \right] + \sigma^2 \mathbb{E} \left[\text{tr} \left[\mathbf{T}_k (\mathbf{S}_k^\top \mathbf{T}_k)^\dagger \mathbf{S}_k^\top \mathbf{S}_\ell (\mathbf{T}_\ell^\top \mathbf{X} \mathbf{S}_\ell)^\dagger \mathbf{T}_\ell^\top \right] \right]$$

Wishart matrix

Conditional expectation



Theorems

Error under Proportional Asymptotics

- Exact error expression for **finite** ensembles

Theorem 1. *In the limit as $(n, p, m, q) \rightarrow \infty$ with $p/n \rightarrow \gamma$, $m/n \rightarrow \eta$, $q/p \rightarrow \alpha$, if $\eta > \alpha\gamma$,*

$$R(\hat{f}_1, \dots, \hat{f}_K) = \frac{K-1}{K} \left(\frac{(1-\alpha)^2 + \sigma^2 \alpha^2 \gamma}{1-\alpha^2 \gamma} \right) + \frac{1}{K} \left(\frac{\eta(1-\alpha) + \sigma^2 \alpha \gamma}{\eta - \alpha \gamma} \right).$$

Error under Proportional Asymptotics

- Exact error expression for **finite** ensembles

Theorem 1. *In the limit as $(n, p, m, q) \rightarrow \infty$ with $p/n \rightarrow \gamma$, $m/n \rightarrow \eta$, $q/p \rightarrow \alpha$, if $\eta > \alpha\gamma$,*

$$R(\hat{f}_1, \dots, \hat{f}_K) = \frac{K-1}{K} \left(\frac{(1-\alpha)^2 + \sigma^2 \alpha^2 \gamma}{1-\alpha^2 \gamma} \right) + \frac{1}{K} \left(\frac{\eta(1-\alpha) + \sigma^2 \alpha \gamma}{\eta - \alpha \gamma} \right).$$

- **Infinite** ensemble error depends only on **feature subsampling**

$$R_{\text{ens}}^{\infty}(\alpha) \triangleq \lim_{K \rightarrow \infty} R(\hat{f}_1, \dots, \hat{f}_K) = \frac{(1-\alpha)^2 + \sigma^2 \alpha^2 \gamma}{1-\alpha^2 \gamma}$$

[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

Optimality Theorem

- Ridge regression: $\hat{f}_\lambda(\mathbf{x}) = \hat{\boldsymbol{\beta}}_\lambda^\top \mathbf{x}$, $\hat{\boldsymbol{\beta}}_\lambda = \arg \min_{\mathbf{b}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \frac{\lambda}{2} \|\mathbf{b}\|_2^2$

[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

Optimality Theorem

- Ridge regression: $\hat{f}_\lambda(\mathbf{x}) = \hat{\boldsymbol{\beta}}_\lambda^\top \mathbf{x}$, $\hat{\boldsymbol{\beta}}_\lambda = \arg \min_{\mathbf{b}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \frac{\lambda}{2} \|\mathbf{b}\|_2^2$
- Tuning feature subsampling is optimal

Theorem 2. *Under the conditions of Theorem 1 and if $\boldsymbol{\beta}^* \sim \mathcal{N}(\mathbf{0}, p^{-1}\mathbf{I})$,*

$$\inf_{\alpha < \gamma^{-1}} R_{\text{ens}}^\infty(\alpha) = \inf_{\lambda > 0} R(\hat{f}_\lambda).$$

[4] DL, H Javadi, RG Baraniuk. “The implicit regularization of ordinary least squares ensembles.” AISTATS, 2020.

Optimality Theorem

- Ridge regression: $\hat{f}_\lambda(\mathbf{x}) = \hat{\boldsymbol{\beta}}_\lambda^\top \mathbf{x}$, $\hat{\boldsymbol{\beta}}_\lambda = \arg \min_{\mathbf{b}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \frac{\lambda}{2} \|\mathbf{b}\|_2^2$
- Tuning feature subsampling is optimal

Theorem 2. *Under the conditions of Theorem 1 and if $\boldsymbol{\beta}^* \sim \mathcal{N}(\mathbf{0}, p^{-1}\mathbf{I})$,*

$$\inf_{\alpha < \gamma^{-1}} R_{\text{ens}}^\infty(\alpha) = \inf_{\lambda > 0} R(\hat{f}_\lambda).$$

Spoiler:

Randomized least squares = ridge + noise

Randomized ensembles = ridge

[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

Optimality Theorem

- Ridge regression: $\hat{f}_\lambda(\mathbf{x}) = \hat{\boldsymbol{\beta}}_\lambda^\top \mathbf{x}$, $\hat{\boldsymbol{\beta}}_\lambda = \arg \min_{\mathbf{b}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \frac{\lambda}{2} \|\mathbf{b}\|_2^2$
- Tuning **feature subsampling** is **optimal**

Theorem 2. Under the conditions of Theorem 1 and if $\boldsymbol{\beta}^* \sim \mathcal{N}(\mathbf{0}, p^{-1}\mathbf{I})$,

$$\inf_{\alpha < \gamma^{-1}} R_{\text{ens}}^\infty(\alpha) = \inf_{\lambda > 0} R(\hat{f}_\lambda).$$

- However, proof **sheds little insight**

$$\inf_{\alpha < \gamma^{-1}} R_{\text{ens}}^\infty(\alpha) = \frac{1}{2} \left(\frac{\gamma - 1}{\gamma} - \sigma^2 + \sqrt{\left(\sigma^2 - \frac{\gamma - 1}{\gamma} \right)^2 + 4\sigma^2} \right) = \inf_{\lambda > 0} R(\hat{f}_\lambda)$$

[4] DL, H Javadi, RG Baraniuk. "The implicit regularization of ordinary least squares ensembles." AISTATS, 2020.

Interpretation of Results

- Since ridge regression is **optimal**, ensemble is **optimal**
 - Ridge regression is the **minimum mean squared error (MMSE)** estimator

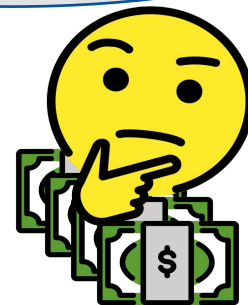
Interpretation of Results

- Since ridge regression is **optimal**, ensemble is **optimal**
 - Ridge regression is the **minimum mean squared error (MMSE)** estimator
- Does this imply ensemble converges to ridge regression?
 - For finite dimensions, MMSE is **unique**
 - For **infinite dimensions?**
 - Optimality theorem suggests convergence to ridge, but **not rigorous**

Interpretation of Results

- Since ridge regression is **optimal**, ensemble is **optimal**
 - Ridge regression is the **minimum mean squared error (MMSE)** estimator
- Does this imply ensemble converges to ridge regression?
 - For finite dimensions, MMSE
 - For **infinite dimension**
 - Optimality theorem says **not rigorous**

But what is this predictor?



Concurrent Observations

- Infinite ensemble with only feature subsampling:

$$\hat{f}_{\text{ens}}^{\infty}(\mathbf{x}) = \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K \hat{\beta}_k^{\top} \mathbf{x} = \mathbf{y}^{\top} \mathbf{X} \left(\lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K \mathbf{S}_k (\mathbf{S}_k^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{S}_k)^{-1} \mathbf{S}_k^{\top} \right) \mathbf{x}$$

Concurrent Observations

- Infinite ensemble with only feature subsampling:

$$\hat{f}_{\text{ens}}^{\infty}(\mathbf{x}) = \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K \hat{\beta}_k^{\top} \mathbf{x} = \mathbf{y}^{\top} \mathbf{X} \left(\lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K \mathbf{S}_k (\mathbf{S}_k^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{S}_k)^{-1} \mathbf{S}_k^{\top} \right) \mathbf{x}$$

- Determinantal Point Process pseudoinverse [7]:

Theorem (Mutny et al. 2020). If $\mathbf{M} \succ \mathbf{0}$ and $\mathbf{S} \sim \text{DPP}(\frac{1}{\lambda} \mathbf{M})$,

$$\mathbb{E} \left[\mathbf{S} (\mathbf{S}^{\top} \mathbf{M} \mathbf{S})^{-1} \mathbf{S}^{\top} \right] = (\mathbf{M} + \lambda \mathbf{I})^{-1},$$

where λ is the solution to $\mathbb{E} \left[\frac{q\mathbf{S}}{p} \right] = \text{tr}(\mathbf{M} (\mathbf{M} + \lambda \mathbf{I})^{-1})$.

[5] M Mutny, M Dereziński, A Krause. “Convergence analysis of block coordinate algorithms with determinantal sampling.” AISTATS, 2020.

Moving to Sketched Ensembles

- **Problem:** strict **isotropic** assumption on **data**
 - **Solution:** let subsampling operators be **isotropic sketches** instead

$$[\mathbf{S}]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{q}), \quad [\mathbf{T}]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{m})$$

Moving to Sketched Ensembles

- **Problem:** strict **isotropic** assumption on **data**
 - **Solution:** let subsampling operators be **isotropic sketches** instead

$$[\mathbf{S}]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{q}), \quad [\mathbf{T}]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{m})$$

- **Problem:** marginal error results
 - **Solution:** use asymptotic equivalences from **random matrix theory (RMT)**

Moving to Sketched Ensembles

- **Problem:** strict **isotropic** assumption on data
 - **Solution:** let subsampling operators be **isotropic sketches** instead

$$[\mathbf{S}]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{q}), \quad [\mathbf{T}]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{m})$$

- **Problem:** marginal error results
 - **Solution:** use asymptotic equivalences from **random matrix theory (RMT)**
- **Problem:** estimators use **ordinary least squares**, $m > q + 1$
 - **Solution:** consider **ridge regression** for arbitrary sketch sizes

Moving to Sketched Ensembles

- **Problem:** strict **isotropic** assumption on data
 - **Solution:** let subsampling operators be **isotropic sketches** instead

$$[\mathbf{S}]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{m}) \quad [\mathbf{T}]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{m})$$

- **Problem:** marginal error
 - **Solution:** use asymptotic

Spoiler:
Randomized least squares = ridge + noise
Randomized ensembles = ridge

matrix theory (RMT)

- **Problem:** estimators use **ordinary least squares**, $m > q + 1$
 - **Solution:** consider **ridge regression** for arbitrary sketch sizes

Asymptotic Equivalences

- Asymptotic equivalences [6]:

Definition. *Two sequences of matrices $\mathbf{A}_n, \mathbf{B}_n$ are **asymptotically equivalent**, written $\mathbf{A}_n \simeq \mathbf{B}_n$, if for every sequence Θ_n having uniformly bounded trace norm, almost surely*

$$\lim_{n \rightarrow \infty} \text{tr} [\Theta_n (\mathbf{A}_n - \mathbf{B}_n)] = 0.$$

[6] E Dobriban, Y Sheng. “Distributed linear regression by averaging.” *Annals of Statistics*, 2021.

Asymptotic Equivalences

- Asymptotic equivalences [6]:

Definition. Two sequences of matrices $\mathbf{A}_n, \mathbf{B}_n$ are *asymptotically equivalent*, written $\mathbf{A}_n \simeq \mathbf{B}_n$, if for every sequence Θ_n having uniformly bounded trace norm, almost surely

$$\lim_{n \rightarrow \infty} \text{tr} [\Theta_n (\mathbf{A}_n - \mathbf{B}_n)] = 0.$$

- Admits a **calculus**:

- Addition $\mathbf{A}_n \simeq \mathbf{B}_n, \mathbf{C}_n \simeq \mathbf{D}_n \implies \mathbf{A}_n + \mathbf{C}_n \simeq \mathbf{B}_n + \mathbf{D}_n$
- Multiplication $\mathbf{A}_n \simeq \mathbf{B}_n \implies \mathbf{C}_n \mathbf{A}_n \mathbf{D}_n \simeq \mathbf{C}_n \mathbf{B}_n \mathbf{D}_n$
- Elements $\mathbf{A}_n \simeq \mathbf{B}_n \implies [\mathbf{A}_n]_{ij} - [\mathbf{B}_n]_{ij} \xrightarrow{\text{a.s.}} 0$
- Differentiation [7] $f(\mathbf{A}_n; z) \simeq g(\mathbf{B}_n; z) \implies f'(\mathbf{A}_n; z) \simeq g'(\mathbf{B}_n; z)$

[6] E Dobriban, Y Sheng. “Distributed linear regression by averaging.” *Annals of Statistics*, 2021.

[7] E Dobriban, Y Sheng. “WONDER: Weighted one-shot distributed ridge regression in high dimensions.” *JMLR*, 2020.

An Asymptotic Equivalence of Resolvents

Theorem (Rubio & Mestre 2011). *Let $\mathbf{Z} \in \mathbb{C}^{n \times p}$ be a random matrix consisting of i.i.d. random variables that have mean 0, variance 1, and finite absolute moment of order $8 + \delta$ for some $\delta > 0$. Let $\mathbf{\Sigma} \in \mathbb{C}^{p \times p}$ be a positive semidefinite matrix with operator norm uniformly bounded in p , and let $\mathbf{X} = \mathbf{Z}\mathbf{\Sigma}^{1/2}$. Then, for $z \in \mathbb{C}^+$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have*

$$\left(\frac{1}{n}\mathbf{X}^H\mathbf{X} - z\mathbf{I}_p\right)^{-1} \simeq (c(z)\mathbf{\Sigma} - z\mathbf{I}_p)^{-1}, \quad (32)$$

where $c(z)$ is the unique solution in \mathbb{C}^- to the fixed point equation

$$\frac{1}{c(z)} - 1 = \frac{1}{n} \operatorname{tr} \left[\mathbf{\Sigma} (c(z)\mathbf{\Sigma} - z\mathbf{I}_p)^{-1} \right]. \quad (33)$$

Furthermore, $\frac{1}{p} \operatorname{tr} \left[\mathbf{\Sigma} (c(z)\mathbf{\Sigma} - z\mathbf{I}_p)^{-1} \right]$ is a Stieltjes transform of a certain positive measure on $\mathbb{R}_{\geq 0}$ with total mass $\frac{1}{p} \operatorname{tr}[\mathbf{\Sigma}]$.

[8] F Rubio, X Mestre. "Spectral convergence for a general class of random matrices." Statistics & Probability Letters, 2011.

An Asymptotic Equivalence of Resolvents

Theorem (Rubio & Mestre 2011). Let $\mathbf{Z} \in \mathbb{C}^{n \times p}$ be a random matrix consisting of i.i.d. random variables that have mean 0, variance 1, and finite absolute moment of order $8 + \delta$ for some $\delta > 0$. Let $\Sigma \in \mathbb{C}^{p \times p}$ be a positive semidefinite matrix with operator norm uniformly bounded in p , and let $\mathbf{X} = \mathbf{Z}\Sigma^{1/2}$. Then, for $z \in \mathbb{C}^+$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

$$\left(\frac{1}{n}\mathbf{X}^H\mathbf{X} - z\mathbf{I}_p\right)^{-1} \simeq (c(z)\Sigma - z\mathbf{I}_p)^{-1}$$

What about real arguments?

where $c(z)$ is the unique solution in \mathbb{C}^- to the fixed point equation

$$\frac{1}{c(z)} - 1 = \frac{1}{n} \text{tr} \left[\Sigma (c(z)\Sigma - z\mathbf{I}_p)^{-1} \right]. \quad (33)$$

Furthermore, $\frac{1}{p} \text{tr} \left[\Sigma (c(z)\Sigma - z\mathbf{I}_p)^{-1} \right]$ is a Stieltjes transform of a certain positive measure on $\mathbb{R}_{\geq 0}$ with total mass $\frac{1}{p} \text{tr}[\Sigma]$.

[8] F Rubio, X Mestre. "Spectral convergence for a general class of random matrices." Statistics & Probability Letters, 2011.

Real-valued Asymptotic Equivalence

Theorem (Rubio & Mestre 2011). Let $\mathbf{Z} \in \mathbb{C}^{n \times p}$ be a random matrix consisting of i.i.d. random variables that have mean 0, variance 1, and finite absolute moment of order $8 + \delta$ for some $\delta > 0$. Let $\Sigma \in \mathbb{C}^{p \times p}$ be a positive semidefinite matrix with operator norm uniformly bounded in p , and let $\mathbf{X} = \mathbf{Z}\Sigma^{1/2}$. Then, for $z \in \mathbb{C}^+$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

$$\left(\frac{1}{n}\mathbf{X}^H\mathbf{X} - z\mathbf{I}_p\right)^{-1} \simeq (c(z)\Sigma - z\mathbf{I}_p)^{-1}, \quad (32)$$

where $c(z)$ is the unique solution in \mathbb{C}^- to the fixed point equation

$$\frac{1}{c(z)} - 1 = \frac{1}{n}\text{tr} \left[\Sigma (c(z)\Sigma - z\mathbf{I}_p)^{-1} \right]. \quad (33)$$

Furthermore, $\frac{1}{p}\text{tr} [\Sigma(c(z)\Sigma - z\mathbf{I}_p)^{-1}]$ is a Stieltjes transform of a certain positive measure on $\mathbb{R}_{>0}$ with total mass $\frac{1}{p}\text{tr}[\Sigma]$.

Real-valued Asymptotic Equivalence

Theorem 3. Let $\zeta_0, z_0 \in \mathbb{R}$ be the unique solutions, satisfying $\zeta_0 < \lambda_{\min}^+(\mathbf{\Sigma})$, to system of equations

$$1 = \frac{1}{n} \text{tr} \left[\mathbf{\Sigma}^2 (\mathbf{\Sigma} - \zeta_0 \mathbf{I}_p)^{-2} \right], \quad z_0 = \zeta_0 \left(1 - \frac{1}{n} \text{tr} \left[\mathbf{\Sigma} (\mathbf{\Sigma} - \zeta_0 \mathbf{I}_p)^{-1} \right] \right). \quad (34)$$

Then, for each $z \in \mathbb{R}$ satisfying $z < \liminf z_0$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

$$z \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq \zeta (\mathbf{\Sigma} - \zeta \mathbf{I}_p)^{-1}, \quad (35)$$

where $\zeta \in \mathbb{R}$ is the unique solution in $(-\infty, \zeta_0)$ to the fixed-point equation

$$z = \zeta \left(1 - \frac{1}{n} \text{tr} \left[\mathbf{\Sigma} (\mathbf{\Sigma} - \zeta \mathbf{I}_p)^{-1} \right] \right). \quad (36)$$

Furthermore, as $n, p \rightarrow \infty$, $|\zeta + \frac{1}{v(z)}| \xrightarrow{\text{a.s.}} 0$, where $v(z)$ is the companion Stieltjes transform of the spectrum of $\frac{1}{n} \mathbf{X}^H \mathbf{X}$ given by

$$v(z) = \frac{1}{n} \text{tr} \left[\left(\frac{1}{n} \mathbf{X} \mathbf{X}^H - z \mathbf{I}_n \right)^{-1} \right],$$

and $|z_0 - \lambda_{\min}^+(\frac{1}{n} \mathbf{X}^H \mathbf{X})| \xrightarrow{\text{a.s.}} 0$.

Theorem (Rubio & Mestre 2011). Let $\mathbf{Z} \in \mathbb{C}^{n \times p}$ be a random matrix consisting of i.i.d. random variables that have mean 0, variance 1, and finite absolute moment of order $8+\delta$ for some $\delta > 0$. Let $\mathbf{\Sigma} \in \mathbb{C}^{p \times p}$ be a positive semidefinite matrix with operator norm uniformly bounded in p , and let $\mathbf{X} = \mathbf{Z}\mathbf{\Sigma}^{1/2}$. Then, for $z \in \mathbb{C}^+$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

$$\left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)$$

where $c(z)$ is the u

Analytic continuation

Furthermore, $\frac{1}{p} \text{tr} \left[\mathbf{\Sigma} (c(z) \mathbf{\Sigma} - z \mathbf{I}_p)^{-1} \right]$ is a Stieltjes transform of a certain positive measure on $\mathbb{R}_{\geq 0}$ with total mass $\frac{1}{p} \text{tr}[\mathbf{\Sigma}]$.

Real-valued Asymptotic Equivalence

Theorem 3. Let $\zeta_0, z_0 \in \mathbb{R}$ be the unique solutions, satisfying $\zeta_0 < \lambda_{\min}^+(\Sigma)$, to system of equations

Limits of negative regularization

$$1 = \frac{1}{n} \operatorname{tr} \left[\Sigma^2 (\Sigma - \zeta_0 \mathbf{I}_p)^{-2} \right], \quad z_0 = \zeta_0 \left(1 - \frac{1}{n} \operatorname{tr} \left[\Sigma (\Sigma - \zeta_0 \mathbf{I}_p)^{-1} \right] \right). \quad (34)$$

Then, for each $z \in \mathbb{R}$ satisfying $z < \liminf z_0$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

$$z \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq \zeta (\Sigma - \zeta \mathbf{I}_p)^{-1}, \quad (35)$$

where $\zeta \in \mathbb{R}$ is the unique solution in $(-\infty, \zeta_0)$ to the fixed-point equation

$$z = \zeta \left(1 - \frac{1}{n} \operatorname{tr} \left[\Sigma (\Sigma - \zeta \mathbf{I}_p)^{-1} \right] \right). \quad (36)$$

Furthermore, as $n, p \rightarrow \infty$, $|\zeta + \frac{1}{v(z)}| \xrightarrow{\text{a.s.}} 0$, where $v(z)$ is the companion Stieltjes transform of the spectrum of $\frac{1}{n} \mathbf{X}^H \mathbf{X}$ given by

$$v(z) = \frac{1}{n} \operatorname{tr} \left[\left(\frac{1}{n} \mathbf{X} \mathbf{X}^H - z \mathbf{I}_n \right)^{-1} \right],$$

and $|z_0 - \lambda_{\min}^+(\frac{1}{n} \mathbf{X}^H \mathbf{X})| \xrightarrow{\text{a.s.}} 0$.

Theorem (Rubio & Mestre 2011). Let $\mathbf{Z} \in \mathbb{C}^{n \times p}$ be a random matrix consisting of i.i.d. random variables that have mean 0, variance 1, and finite absolute moment of order $8+\delta$ for some $\delta > 0$. Let $\Sigma \in \mathbb{C}^{p \times p}$ be a positive semidefinite matrix with operator norm uniformly bounded in p , and let $\mathbf{X} = \mathbf{Z}\Sigma^{1/2}$. Then, for $z \in \mathbb{C}^+$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

$$\left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq (c(z)\Sigma - z \mathbf{I}_p)^{-1}, \quad (32)$$

where $c(z)$ is the unique solution in \mathbb{C}^- to the fixed point equation

$$\frac{1}{c(z)} - 1 = \frac{1}{n} \operatorname{tr} \left[\Sigma (c(z)\Sigma - z \mathbf{I}_p)^{-1} \right]. \quad (33)$$

Furthermore, $\frac{1}{p} \operatorname{tr} [\Sigma (c(z)\Sigma - z \mathbf{I}_p)^{-1}]$ is a Stieltjes transform of a certain positive measure on $\mathbb{R}_{\geq 0}$ with total mass $\frac{1}{p} \operatorname{tr}[\Sigma]$.

Real-valued Asymptotic Equivalence

Theorem 3. Let $\zeta_0, z_0 \in \mathbb{R}$ be the unique solutions, satisfying $\zeta_0 < \lambda_{\min}^+(\Sigma)$, to system of equations

$$1 = \frac{1}{n} \text{tr} \left[\Sigma^2 (\Sigma - \zeta_0 \mathbf{I}_p)^{-2} \right], \quad z_0 = \zeta_0 \left(1 - \frac{1}{n} \text{tr} \left[\Sigma (\Sigma - \zeta_0 \mathbf{I}_p)^{-1} \right] \right). \quad (34)$$

Then, for each $z \in \mathbb{R}$ satisfying $z < \liminf z_0$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

Theorem (Rubio & Mestre 2011). Let $\mathbf{Z} \in \mathbb{C}^{n \times p}$ be a random matrix consisting of i.i.d. random variables that have mean 0, variance 1, and finite absolute moment of order $8+\delta$ for some $\delta > 0$. Let $\Sigma \in \mathbb{C}^{p \times p}$ be a positive semidefinite matrix with operator norm uniformly bounded in p , and let $\mathbf{X} = \mathbf{Z}\Sigma^{1/2}$. Then, for $z \in \mathbb{C}^+$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq$

$$\left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq (c(z) \Sigma - z \mathbf{I}_p)^{-1} \quad (33)$$

$$\frac{1}{c(z)} - 1 = \frac{1}{n} \text{tr} \left[\Sigma (c(z) \Sigma - z \mathbf{I}_p)^{-1} \right]. \quad (33)$$

Furthermore, $\frac{1}{p} \text{tr} [\Sigma (c(z) \Sigma - z \mathbf{I}_p)^{-1}]$ is a Stieltjes transform of a certain positive measure on $\mathbb{R}_{\geq 0}$ with total mass $\frac{1}{p} \text{tr} [\Sigma]$.

Reparameterization

where $\zeta \in \mathbb{R}$ is the unique solution in $(-\infty, \zeta_0)$ to the fixed-point equation

$$z \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq \zeta (\Sigma - \zeta \mathbf{I}_p)^{-1}, \quad (35)$$

$$z = \zeta \left(1 - \frac{1}{n} \text{tr} \left[\Sigma (\Sigma - \zeta \mathbf{I}_p)^{-1} \right] \right). \quad (36)$$

Furthermore, as $n, p \rightarrow \infty$, $|\zeta + \frac{1}{v(z)}| \xrightarrow{\text{a.s.}} 0$, where $v(z)$ is the companion Stieltjes transform of the spectrum of $\frac{1}{n} \mathbf{X}^H \mathbf{X}$ given by

$$v(z) = \frac{1}{n} \text{tr} \left[\left(\frac{1}{n} \mathbf{X} \mathbf{X}^H - z \mathbf{I}_n \right)^{-1} \right],$$

and $|z_0 - \lambda_{\min}^+(\frac{1}{n} \mathbf{X}^H \mathbf{X})| \xrightarrow{\text{a.s.}} 0$.

Real-valued Asymptotic Equivalence

Theorem 3. Let $\zeta_0, z_0 \in \mathbb{R}$ be the unique solutions, satisfying $\zeta_0 < \lambda_{\min}^+(\Sigma)$, to system of equations

$$1 = \frac{1}{n} \text{tr} \left[\Sigma^2 (\Sigma - \zeta_0 \mathbf{I}_p)^{-2} \right], \quad z_0 = \zeta_0 \left(1 - \frac{1}{n} \text{tr} \left[\Sigma (\Sigma - \zeta_0 \mathbf{I}_p)^{-1} \right] \right). \quad (34)$$

Then, for each $z \in \mathbb{R}$ satisfying $z < \liminf z_0$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

$$z \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq \zeta (\Sigma - \zeta \mathbf{I}_p)^{-1}, \quad (35)$$

where $\zeta \in \mathbb{R}$ is the unique solution in $(-\infty, \infty)$ of the equation $\zeta = \zeta \left(1 - \frac{1}{n} \text{tr} \left[\Sigma (\Sigma - \zeta \mathbf{I}_p)^{-1} \right] \right)$. Explicit form of implicit value

$$z = \zeta \left(1 - \frac{1}{n} \text{tr} \left[\Sigma (\Sigma - \zeta \mathbf{I}_p)^{-1} \right] \right). \quad (36)$$

Furthermore, as $n, p \rightarrow \infty$, $\left| \zeta + \frac{1}{v(z)} \right| \xrightarrow{\text{a.s.}} 0$, where $v(z)$ is the companion Stieltjes transform of the spectrum of $\frac{1}{n} \mathbf{X}^H \mathbf{X}$ given by

$$v(z) = \frac{1}{n} \text{tr} \left[\left(\frac{1}{n} \mathbf{X} \mathbf{X}^H - z \mathbf{I}_n \right)^{-1} \right]$$

and $|z_0 - \lambda_{\min}^+(\frac{1}{n} \mathbf{X}^H \mathbf{X})| \xrightarrow{\text{a.s.}} 0$.

Theorem (Rubio & Mestre 2011). Let $\mathbf{Z} \in \mathbb{C}^{n \times p}$ be a random matrix consisting of i.i.d. random variables that have mean 0, variance 1, and finite absolute moment of order $8+\delta$ for some $\delta > 0$. Let $\Sigma \in \mathbb{C}^{p \times p}$ be a positive semidefinite matrix with operator norm uniformly bounded in p , and let $\mathbf{X} = \mathbf{Z}\Sigma^{1/2}$. Then, for $z \in \mathbb{C}^+$, as $n, p \rightarrow \infty$ such that $0 < \liminf \frac{p}{n} \leq \limsup \frac{p}{n} < \infty$, we have

$$\left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq (c(z)\Sigma - z \mathbf{I}_p)^{-1}, \quad (32)$$

where $c(z)$ is the unique solution in \mathbb{C}^- to the fixed point equation

$$\frac{1}{c(z)} - 1 = \frac{1}{n} \text{tr} \left[\Sigma (c(z)\Sigma - z \mathbf{I}_p)^{-1} \right]. \quad (33)$$

Furthermore, $\frac{1}{p} \text{tr} \left[\Sigma (c(z)\Sigma - z \mathbf{I}_p)^{-1} \right]$ is a Stieltjes transform of a certain positive measure on $\mathbb{R}_{\geq 0}$ with total mass $\frac{1}{p} \text{tr}[\Sigma]$.

Obtaining a Sketching Equivalence

$$z \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq \zeta (\boldsymbol{\Sigma} - \zeta \mathbf{I}_p)^{-1}$$

Obtaining a Sketching Equivalence

$$z \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq \zeta \left(\boldsymbol{\Sigma} - \zeta \mathbf{I}_p \right)^{-1}$$



$$\mathbf{X} = \sqrt{q} \mathbf{S}^H \mathbf{A}^{1/2}$$

$$\lambda = -z$$

$$\mu = -\zeta$$

Obtaining a Sketching Equivalence

$$z \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq \zeta \left(\boldsymbol{\Sigma} - \zeta \mathbf{I}_p \right)^{-1}$$



$$\mathbf{X} = \sqrt{q} \mathbf{S}^H \mathbf{A}^{1/2}$$

$$\lambda = -z$$

$$\mu = -\zeta$$

$$\mathbf{I}_p - \lambda \left(\mathbf{A}^{1/2} \mathbf{S} \mathbf{S}^H \mathbf{A}^{1/2} + \lambda \mathbf{I}_p \right)^{-1} \simeq \mathbf{I}_p - \mu \left(\mathbf{A} + \mu \mathbf{I}_p \right)^{-1}$$

$$\implies \mathbf{A}^{1/2} \mathbf{S} \left(\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q \right)^{-1} \mathbf{S}^H \mathbf{A}^{1/2} \simeq \mathbf{A}^{1/2} \left(\mathbf{A} + \mu \mathbf{I}_p \right)^{-1} \mathbf{A}^{1/2}$$

Obtaining a Sketching Equivalence

$$z \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} - z \mathbf{I}_p \right)^{-1} \simeq \zeta \left(\boldsymbol{\Sigma} - \zeta \mathbf{I}_p \right)^{-1}$$

Uniform convergence



$$\begin{aligned} \mathbf{X} &= \sqrt{q} \mathbf{S}^H \mathbf{A}^{1/2} \\ \lambda &= -z \\ \mu &= -\zeta \end{aligned}$$

$$\begin{aligned} \mathbf{I}_p - \lambda \left(\mathbf{A}^{1/2} \mathbf{S} \mathbf{S}^H \mathbf{A}^{1/2} + \lambda \mathbf{I}_p \right)^{-1} &\simeq \mathbf{I}_p - \mu \left(\mathbf{A} + \mu \mathbf{I}_p \right)^{-1} \\ \Rightarrow \cancel{\mathbf{A}^{1/2}} \mathbf{S} \left(\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q \right)^{-1} \mathbf{S}^H \cancel{\mathbf{A}^{1/2}} &\simeq \cancel{\mathbf{A}^{1/2}} \left(\mathbf{A} + \mu \mathbf{I}_p \right)^{-1} \cancel{\mathbf{A}^{1/2}} \end{aligned}$$



$$\hat{f}_{\text{ens}}(\mathbf{x}) = \mathbf{y}^\top \mathbf{X} \left(\frac{1}{K} \sum_{k=1}^K \mathbf{S}_k \left(\mathbf{S}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{S}_k \right)^{-1} \mathbf{S}_k^\top \right) \mathbf{x}$$

First-order Sketching Equivalence

Theorem 4. For each $\lambda > \limsup \lambda_0$, as $q, p \rightarrow \infty$,

$$\mathbf{S}(\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \mathbf{S}^H \simeq (\mathbf{A} + \mu \mathbf{I}_p)^{-1},$$

where μ is the unique solution in (μ_0, ∞) to the fixed point equation

$$\lambda = \mu \left(1 - \frac{1}{q} \text{tr} \left[\mathbf{A} (\mathbf{A} + \mu \mathbf{I}_p)^{-1} \right] \right).$$

First-order Sketching Equivalence

Theorem 4. For each $\lambda > \limsup \lambda_0$, as $q, p \rightarrow \infty$,

$$\mathbf{S}(\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \mathbf{S}^H \simeq (\mathbf{A} + \mu \mathbf{I}_p)^{-1},$$

where μ is the unique solution in (μ_0, ∞) to the fixed point equation

$$\lambda = \mu \left(1 - \frac{1}{q} \text{tr} \left[\mathbf{A} (\mathbf{A} + \mu \mathbf{I}_p)^{-1} \right] \right).$$

- That is, **sketching + ridge** = **another ridge without sketching**.



First-order Sketching Equivalence

Theorem 4. For each $\lambda > \limsup \lambda_0$, as $q, p \rightarrow \infty$,

$$\mathbf{S}(\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \mathbf{S}^H \simeq (\mathbf{A} + \mu \mathbf{I}_p)^{-1}$$

where μ is the unique solution in (μ_0, ∞) to the fixed point

$$\mathbb{E} \left[\frac{q\mathbf{S}}{p} \right] = \text{tr}(\mathbf{M} (\mathbf{M} + \mu \mathbf{I})^{-1})$$

$$\lambda = \mu \left(1 - \frac{1}{q} \text{tr} \left[\mathbf{A} (\mathbf{A} + \mu \mathbf{I}_p)^{-1} \right] \right).$$

Same as DPP when $\lambda = 0$

- That is, sketching + ridge = another ridge without sketching.



First-order Sketching Equivalence

Theorem 4. For each $\lambda > \limsup \lambda_0$, as $q, p \rightarrow \infty$,

$$\mathbf{S}(\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \mathbf{S}^H \simeq (\mathbf{A} + \mu \mathbf{I}_p)^{-1},$$

where μ is the unique solution in (μ_0, ∞) to the fixed point equation

$$\lambda = \mu \left(1 - \frac{1}{q} \text{tr} \left[\mathbf{A} (\mathbf{A} + \mu \mathbf{I}_p)^{-1} \right] \right).$$

• That is, sketching +

Spoiler:

Randomized least squares = ridge + noise

Randomized ensembles = ridge

sketching.



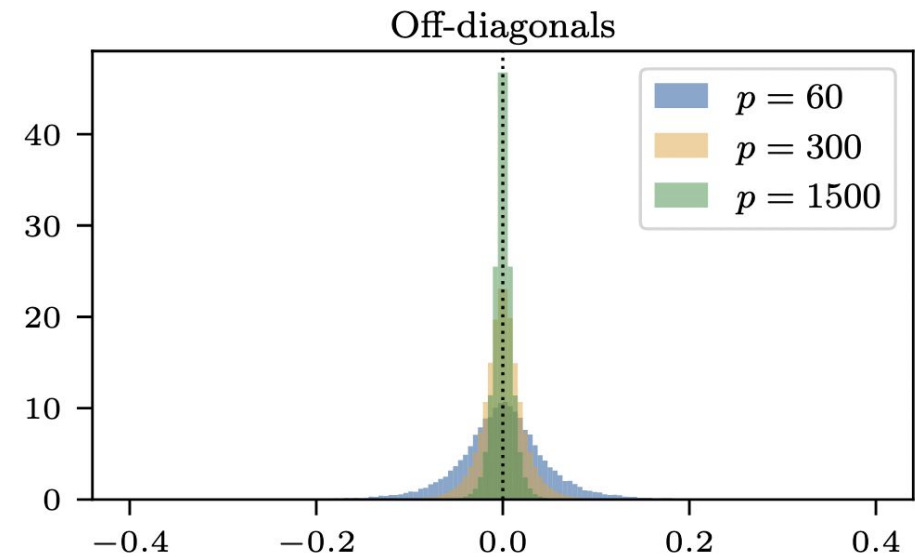
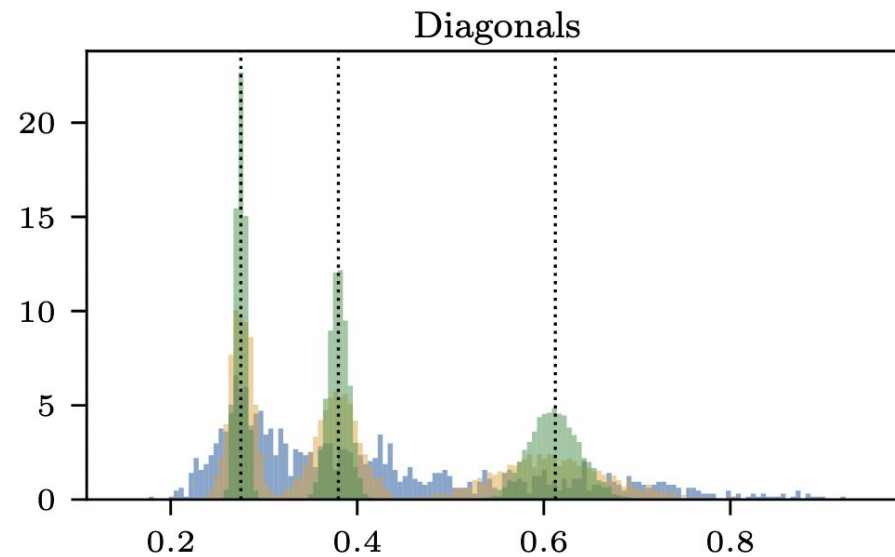
Element-wise Convergence

$$\mathbf{A} = \text{diag}(0 \dots, 1 \dots, 2 \dots)$$

$$\alpha = 0.8$$

$$\lambda = 1$$

$$\mu \approx 1.63$$



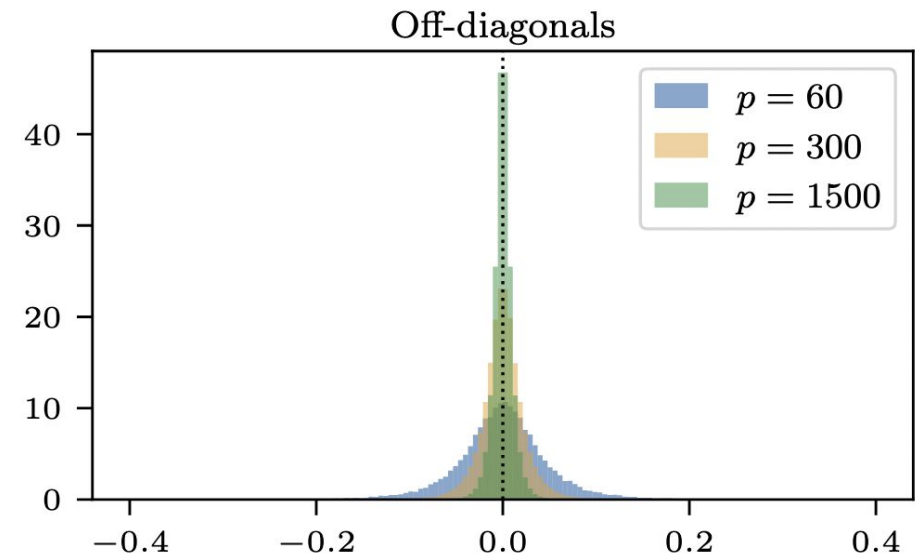
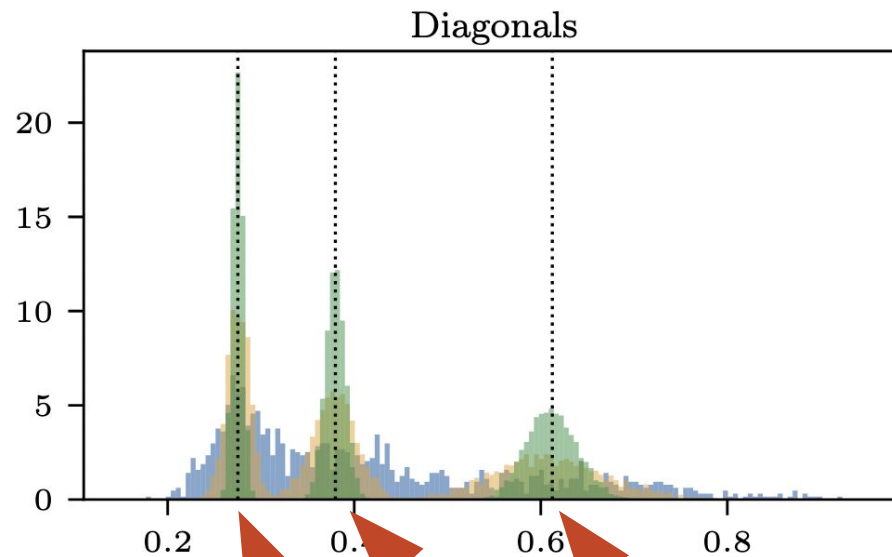
Element-wise Convergence

$$\mathbf{A} = \text{diag}(0 \dots, 1 \dots, 2 \dots)$$

$$\alpha = 0.8$$

$$\lambda = 1$$

$$\mu \approx 1.63$$

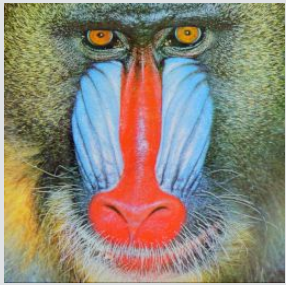


$$\frac{1}{a+\mu} = \frac{1}{3.63}, \frac{1}{2.63}, \frac{1}{1.63}$$

A Blurring Analogy

Original Image

256x256



A Blurring Analogy

$S^H AS$
Downsample



64x64





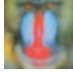




32x32



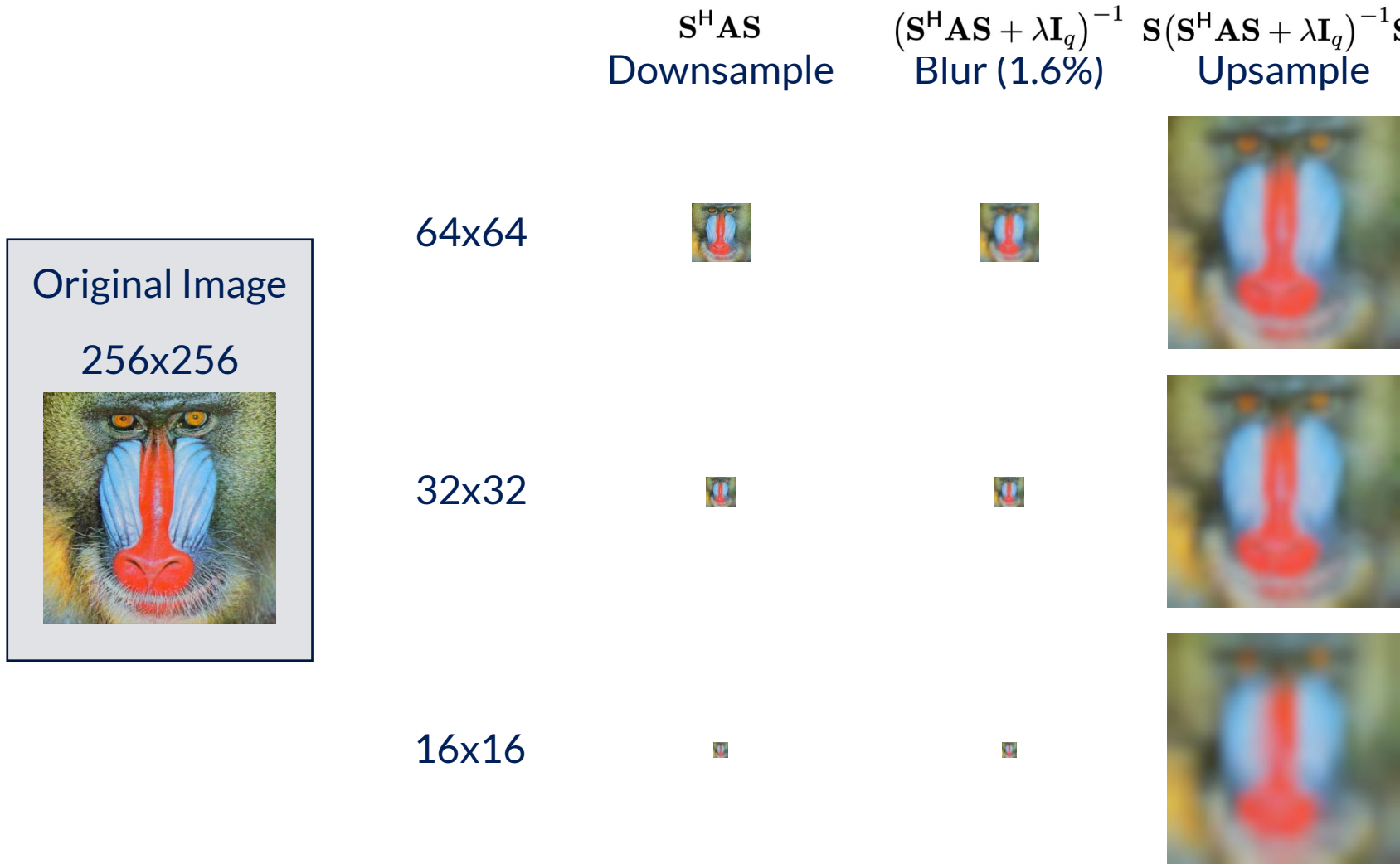
16x16



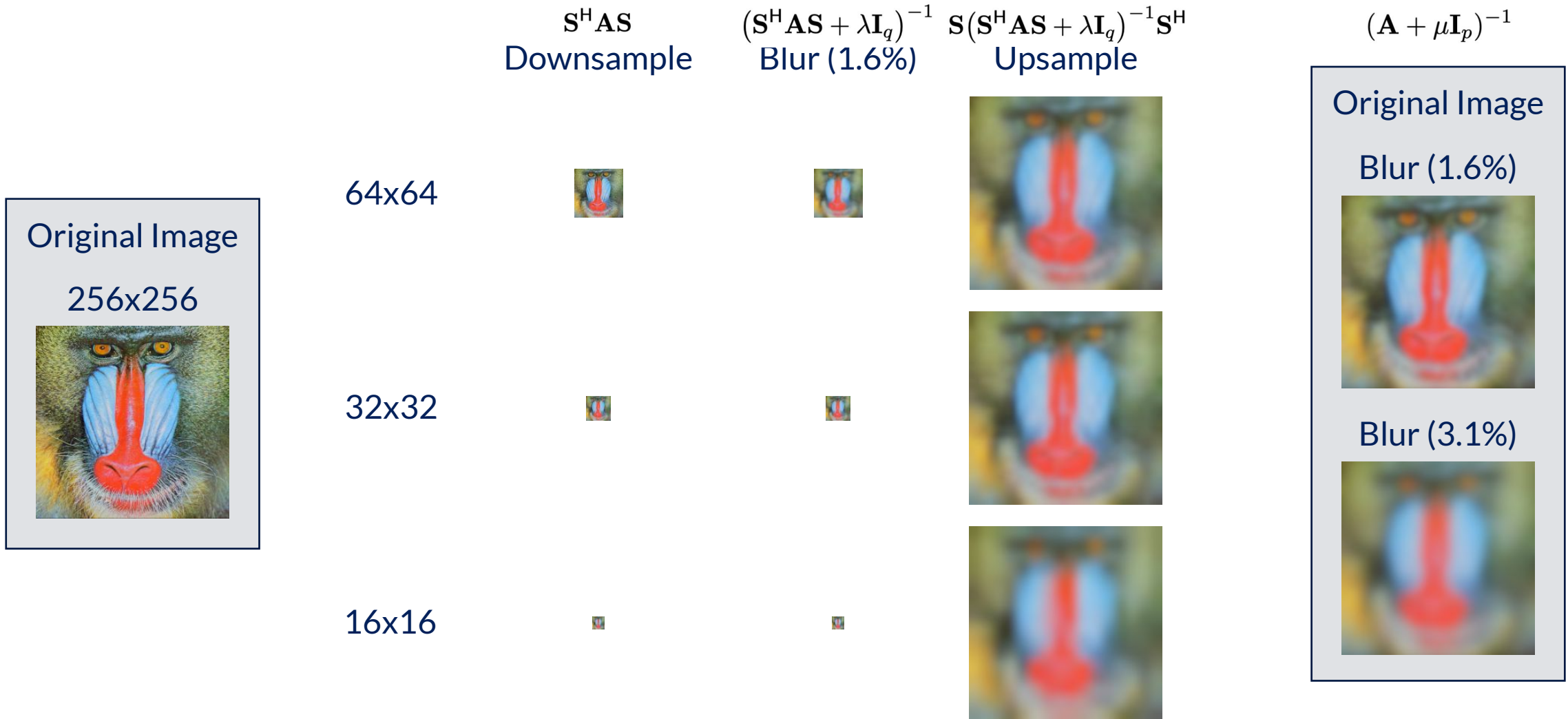
A Blurring Analogy

	$S^H AS$ Downsample	$(S^H AS + \lambda I_q)^{-1}$ Blur (1.6%)
Original Image 256x256 		
64x64		
32x32		
16x16		

A Blurring Analogy

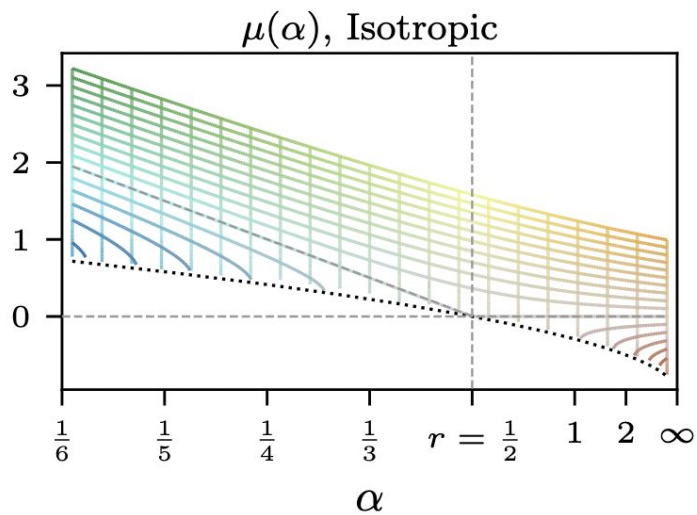
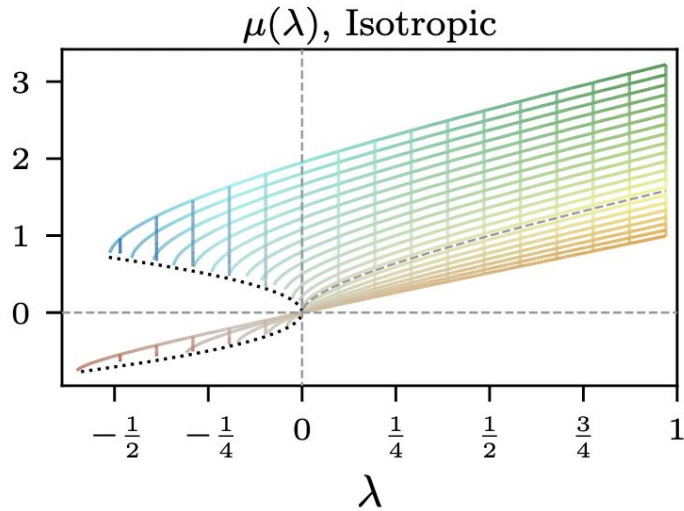


A Blurring Analogy

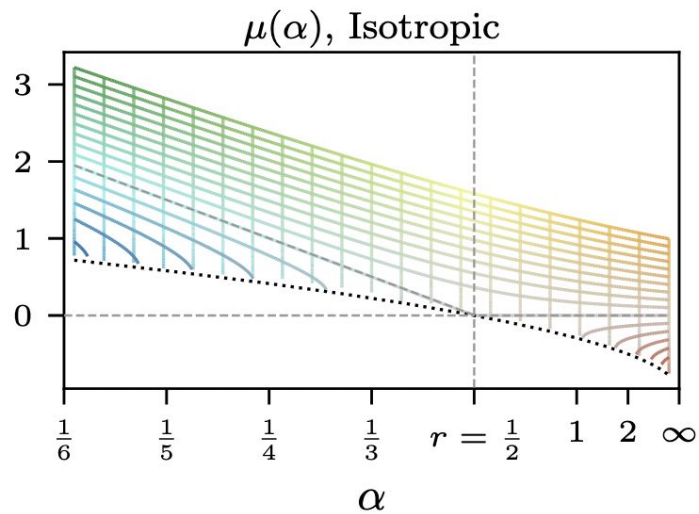
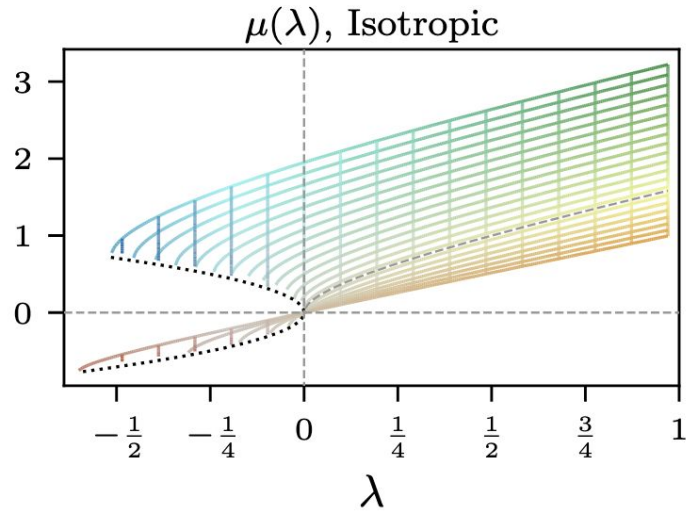


Implicit Regularization Behavior

- **Example:** isotropic spectrum with $r(\mathbf{A}) = \frac{1}{2}$

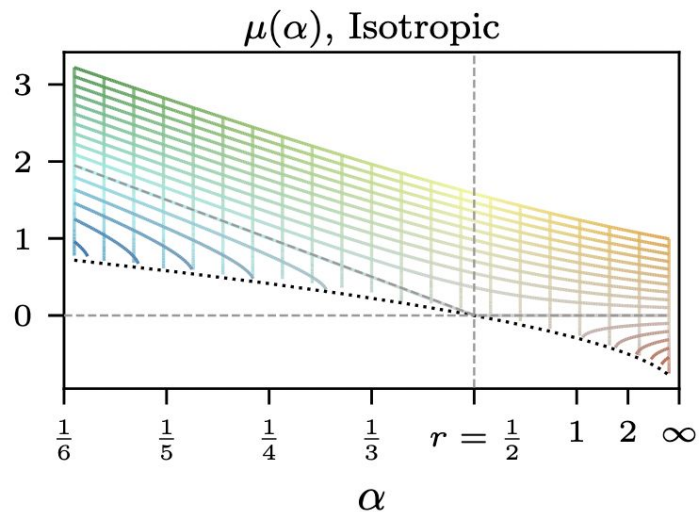
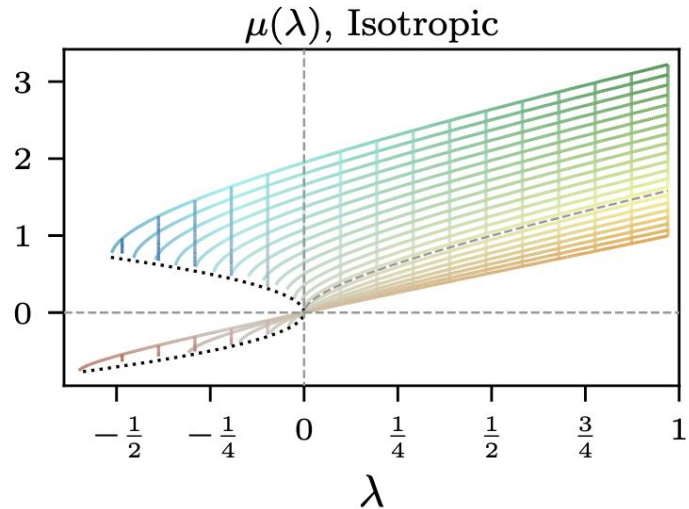


Implicit Regularization Behavior



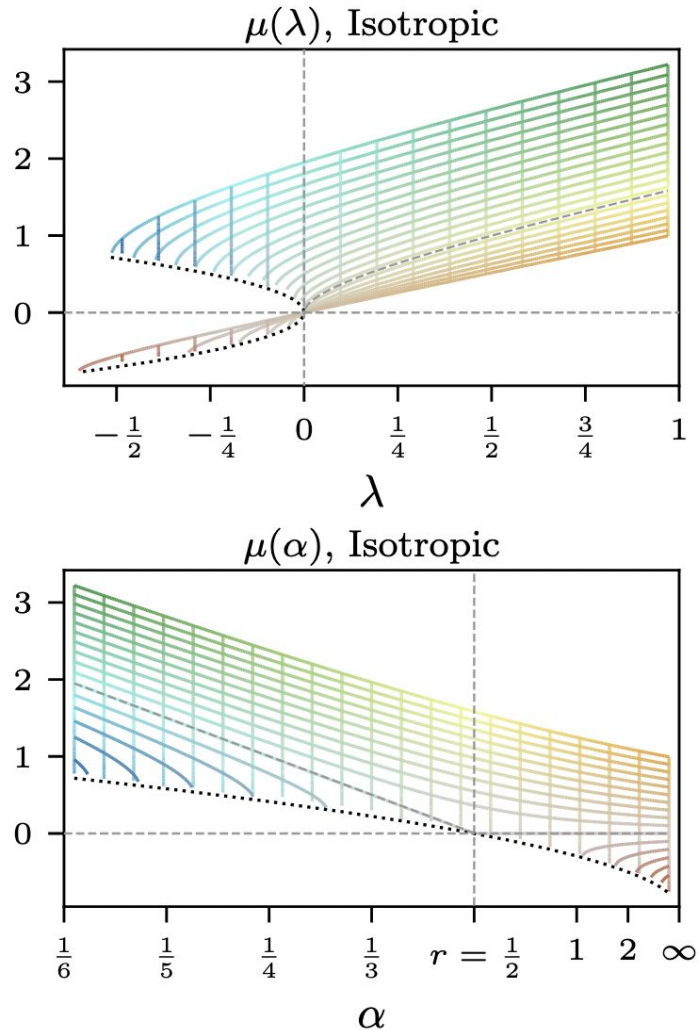
- **Example:** isotropic spectrum with $r(\mathbf{A}) = \frac{1}{2}$
- $\lambda \mapsto \mu$ is **increasing** and concave
- $\alpha \mapsto \mu$ is **decreasing** unless $\alpha > r(\mathbf{A})$ and $\lambda < 0$

Implicit Regularization Behavior



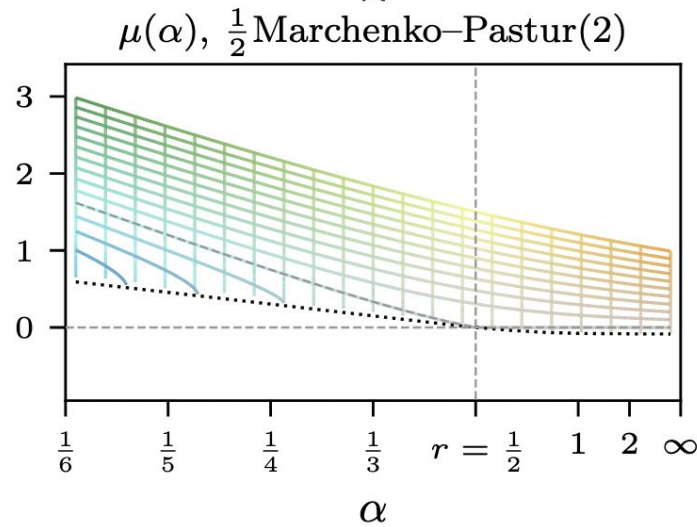
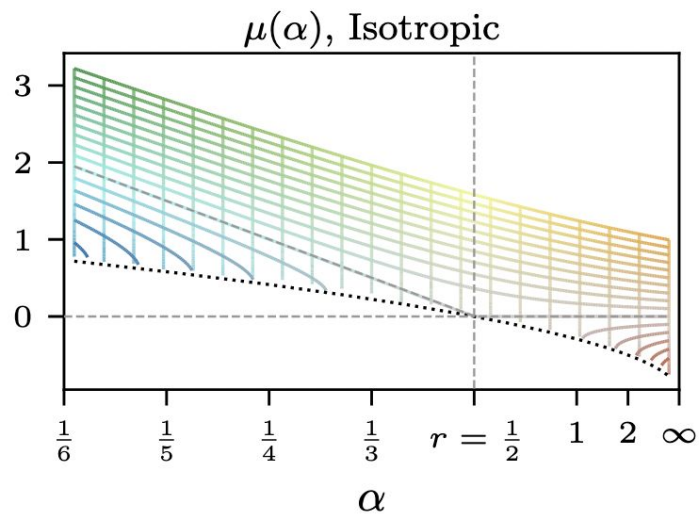
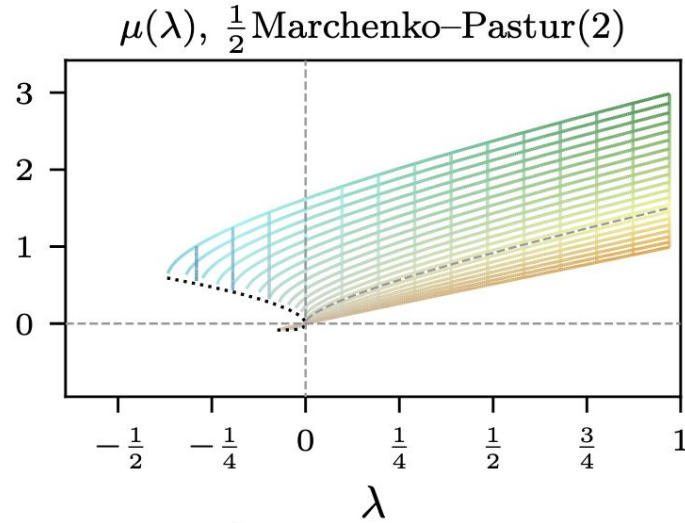
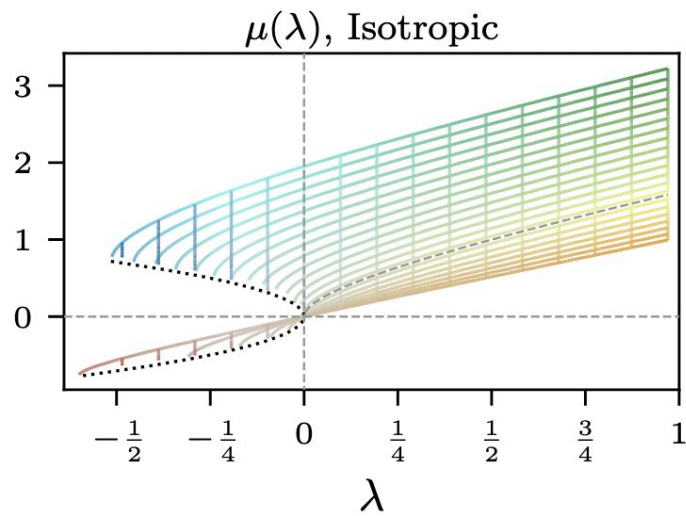
- **Example:** isotropic spectrum with $r(\mathbf{A}) = \frac{1}{2}$
- $\lambda \mapsto \mu$ is **increasing** and concave
- $\alpha \mapsto \mu$ is **decreasing** unless $\alpha > r(\mathbf{A})$ and $\lambda < 0$
- $\mu \geq \lambda$ unless $\alpha > r(\mathbf{A})$ and $\lambda < 0$
- $\mu \leq \lambda + \frac{1}{q} \text{tr}[\mathbf{A}]$

Implicit Regularization Behavior



- **Example:** isotropic spectrum with $r(\mathbf{A}) = \frac{1}{2}$
- $\lambda \mapsto \mu$ is **increasing** and concave
- $\alpha \mapsto \mu$ is **decreasing** unless $\alpha > r(\mathbf{A})$ and $\lambda < 0$
- $\mu \geq \lambda$ unless $\alpha > r(\mathbf{A})$ and $\lambda < 0$
- $\mu \leq \lambda + \frac{1}{q} \text{tr}[\mathbf{A}]$
- $\text{sign}(\mu) = \text{sign}(\lambda)$ if $\alpha > r(\mathbf{A})$
- else $\mu \geq 0$

Implicit Regularization Behavior



$\mu(\lambda)$ is concave and increasing unless $\alpha > r(\mathbf{A})$ and $\lambda < 0$
 $\mu(\lambda)$ is convex and decreasing unless $\alpha > r(\mathbf{A})$ and $\lambda < 0$
 $\mu(\lambda)$ is convex and decreasing if $\alpha > r(\mathbf{A})$

First-Order is Not Enough

- First-order equivalence is similar to **expectation** equivalence

First-Order is Not Enough

- First-order equivalence is similar to **expectation** equivalence
- $\mathbb{E}[X] = \mathbb{E}[Y]$ **does not imply** that $\mathbb{E}[X^2] = \mathbb{E}[Y^2]$
- Similarly, **products of equivalences** **do not compose** if not independent

First-Order is Not Enough

- First-order equivalence is similar to **expectation** equivalence
- $\mathbb{E}[X] = \mathbb{E}[Y]$ **does not imply** that $\mathbb{E}[X^2] = \mathbb{E}[Y^2]$
- Similarly, **products of equivalences do not compose** if not independent
- **Solution:** derivative rule of asymptotic equivalence

$$\frac{d}{dz} (\mathbf{A} - z\mathbf{I})^{-1} = -(\mathbf{A} - z\mathbf{I})^{-2}$$

Second-order Sketching Equivalence

Theorem 5. *If $\Psi \in \mathbb{C}^{p \times p}$ is a deterministic or random positive semidefinite matrix independent of \mathbf{S} with $\|\Psi\|_{\text{op}}$ uniformly bounded in p , then*

$$\mathbf{S}(\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \mathbf{S}^H \Psi \mathbf{S} (\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \mathbf{S}^H \simeq (\mathbf{A} + \mu \mathbf{I}_p)^{-1} (\Psi + \mu' \mathbf{I}_p) (\mathbf{A} + \mu \mathbf{I}_p)^{-1},$$

where

$$\mu' = \frac{\frac{1}{q} \text{tr} \left[\mu^3 (\mathbf{A} + \mu \mathbf{I}_p)^{-1} \Psi (\mathbf{A} + \mu \mathbf{I}_p)^{-1} \right]}{\lambda + \frac{1}{q} \text{tr} \left[\mu^2 \mathbf{A} (\mathbf{A} + \mu \mathbf{I}_p)^{-2} \right]} \geq 0.$$

Second-order Sketching Equivalence

Theorem 5. If $\Psi \in \mathbb{C}^{p \times p}$ is a deterministic or random positive semidefinite matrix independent of \mathbf{S} with $\|\Psi\|_{\text{op}}$ uniformly bounded in p , then

$$\mathbf{S}(\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \mathbf{S}^H \Psi \mathbf{S} (\mathbf{S}^H \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \mathbf{S}^H \simeq (\mathbf{A} + \mu \mathbf{I}_p)^{-1} (\Psi + \mu' \mathbf{I}_p) (\mathbf{A} + \mu \mathbf{I}_p)^{-1},$$

where

$$\mu' = \frac{\frac{1}{q} \text{tr} \left[\mu^3 (\mathbf{A} + \mu \mathbf{I}_p)^{-1} \Psi (\mathbf{A} + \mu \mathbf{I}_p)^{-1} \right]}{\lambda + \frac{1}{q} \text{tr} \left[\mu^2 \mathbf{A} (\mathbf{A} + \mu \mathbf{I}_p)^{-2} \right]} \geq 0.$$

- That is, second-order adds an isotropic **inflation factor**
- Depends significantly on the choice of Ψ

Second-order Sketching Equivalence

Theorem 5. If $\Psi \in \mathbb{C}^{p \times p}$ is a deterministic or random positive semidefinite matrix independent of \mathbf{S} with $\|\Psi\|_{\text{op}}$ uniformly bounded in p , then

where

$$\mathbf{S}(\mathbf{S}^H \mathbf{A} + \mu \mathbf{I}_p)^{-1} \mathbf{S} \approx (\mathbf{A} + \mu \mathbf{I}_p)^{-1} (\Psi + \mu' \mathbf{I}_p) (\mathbf{A} + \mu \mathbf{I}_p)^{-1},$$

Spoiler:

Randomized least squares = ridge + noise

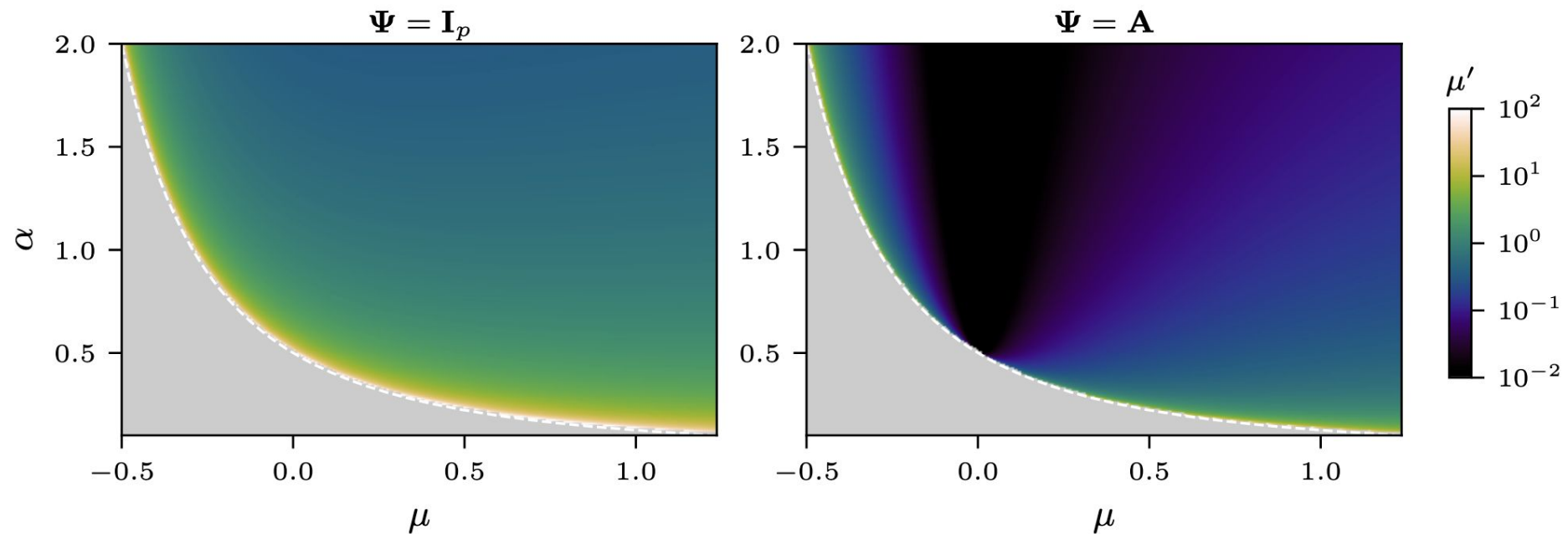
Randomized ensembles = ridge

$$\frac{\Psi (\mathbf{A} + \mu \mathbf{I}_p)^{-1}}{\lambda + \frac{1}{q} \text{tr} [\mu^2 \mathbf{A} (\mathbf{A} + \mu \mathbf{I}_p)^{-2}]} \geq 0.$$

- That is, second-order adds an isotropic **inflation factor**
- Depends significantly on the choice of Ψ

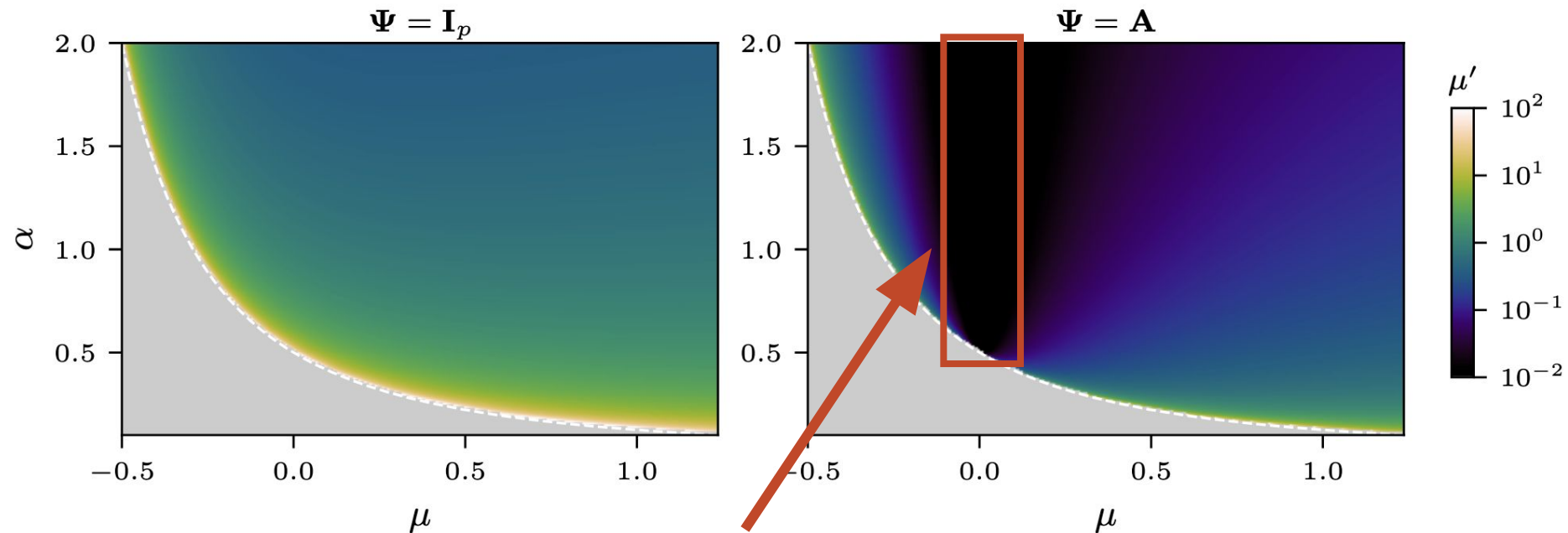
One Incredible Regime

- **Example:** isotropic spectrum with $r(\mathbf{A}) = \frac{1}{2}$



One Incredible Regime

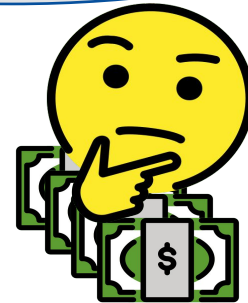
- **Example:** isotropic spectrum with $r(\mathbf{A}) = \frac{1}{2}$



- If $\Psi \in \text{Range}(\mathbf{A})$, $\alpha > r(\mathbf{A})$, and $\mu = \lambda = 0$, there is no inflation
 - Agrees with classical sketching results: sketch larger than rank
 - Sketching is ideal for benign overfitting

Application: Ridge Regression

How about a machine learning problem?



Sketched Ridge Regression

- Primal (observations) and dual (features) sketching:

$$\hat{\beta}_P \triangleq \arg \min_{\mathbf{b}} \frac{1}{n} \|\mathbf{T}^\top (\mathbf{y} - \mathbf{X}\mathbf{b})\|_2^2 + \lambda \|\mathbf{b}\|_2^2$$

$$\hat{\beta}_D \triangleq \mathbf{S} \cdot \arg \min_{\mathbf{b}} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\mathbf{S}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_2^2$$

$$\hat{\beta}_{PD} \triangleq \mathbf{S} \cdot \arg \min_{\mathbf{b}} \frac{1}{n} \|\mathbf{T}^\top (\mathbf{y} - \mathbf{X}\mathbf{S}\mathbf{b})\|_2^2 + \lambda \|\mathbf{b}\|_2^2$$

Preparing for Equivalences

- Express in terms of the **sketched pseudoinverse**:

$$\widehat{\boldsymbol{\beta}}_{\psi} = \frac{1}{\sqrt{n}} \mathbf{X}_{\psi}^{\dagger} \mathbf{y}, \quad \psi \in \{P, D, PD\}$$

$$\mathbf{X}_{P}^{\dagger} \triangleq \frac{1}{\sqrt{n}} \left(\frac{1}{n} \mathbf{X}^{\top} \mathbf{T} \mathbf{T}^{\top} \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^{\top} \mathbf{T} \mathbf{T}^{\top}$$

$$\mathbf{X}_{D}^{\dagger} \triangleq \frac{1}{\sqrt{n}} \mathbf{S} \left(\frac{1}{n} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{S} + \lambda \mathbf{I} \right)^{-1} \mathbf{S}^{\top} \mathbf{X}^{\top}$$

$$\mathbf{X}_{PD}^{\dagger} \triangleq \frac{1}{\sqrt{n}} \mathbf{S} \left(\frac{1}{n} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{T} \mathbf{T}^{\top} \mathbf{X} \mathbf{S} + \lambda \mathbf{I} \right)^{-1} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{T} \mathbf{T}^{\top}$$

First-order Data Pseudoinverse Equivalence

Theorem 6. *If $\left\| \frac{1}{\sqrt{n}} \mathbf{X} \right\|_{\text{op}}$ is uniformly bounded in p , then as $m, n, q, p \rightarrow \infty$,*

$$\mathbf{X}_{\psi}^{\dagger} \simeq \frac{1}{\sqrt{n}} \left(\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} + \mu_{\psi} \mathbf{I} \right)^{-1} \mathbf{X}^{\text{H}},$$

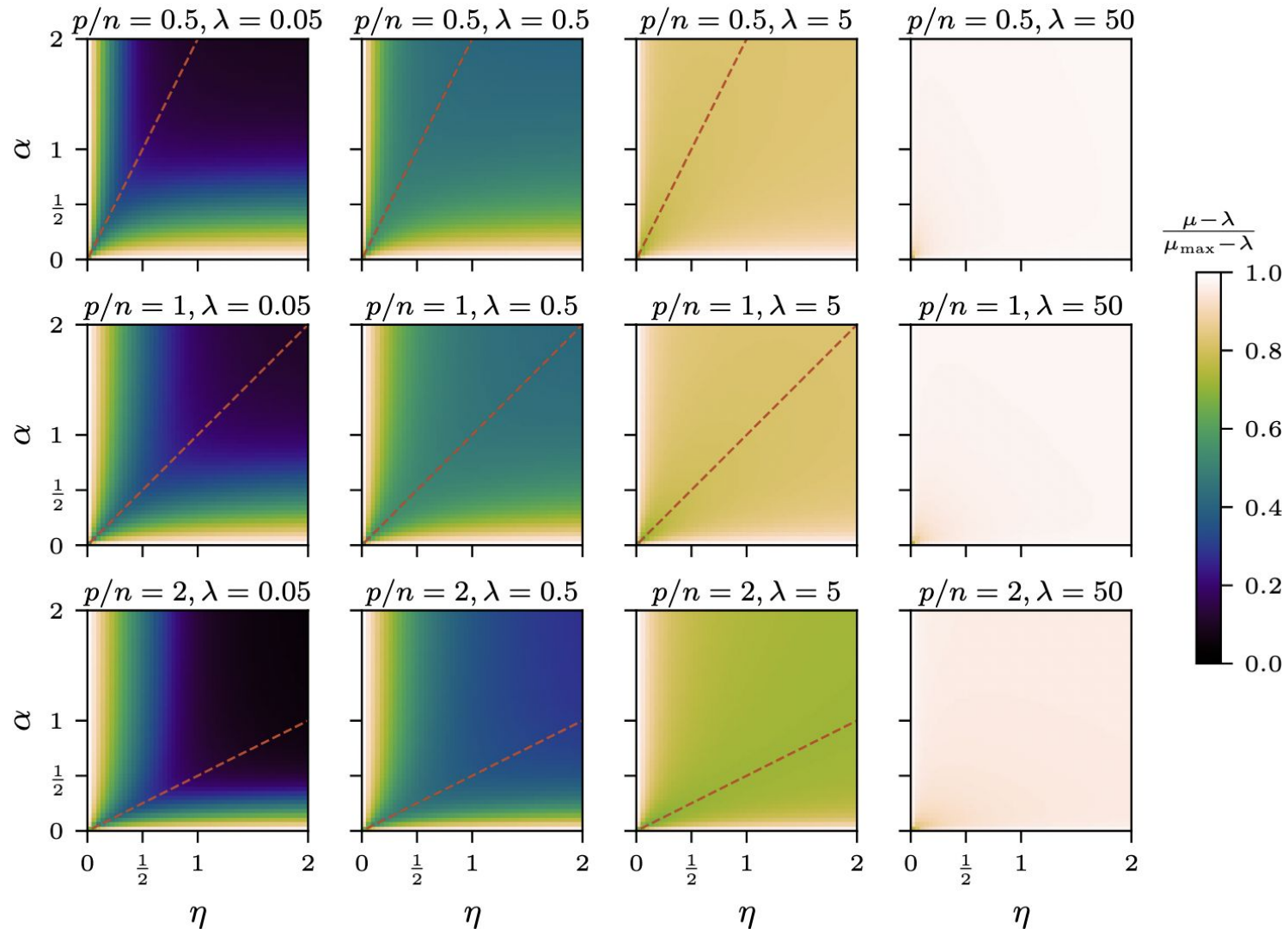
where μ_{ψ} are the most positive solutions to the equations

$$\lambda = \mu_{\text{P}} \left(1 - \frac{1}{m} \text{tr} \left[\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} \left(\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} + \mu_{\text{P}} \mathbf{I} \right)^{-1} \right] \right),$$

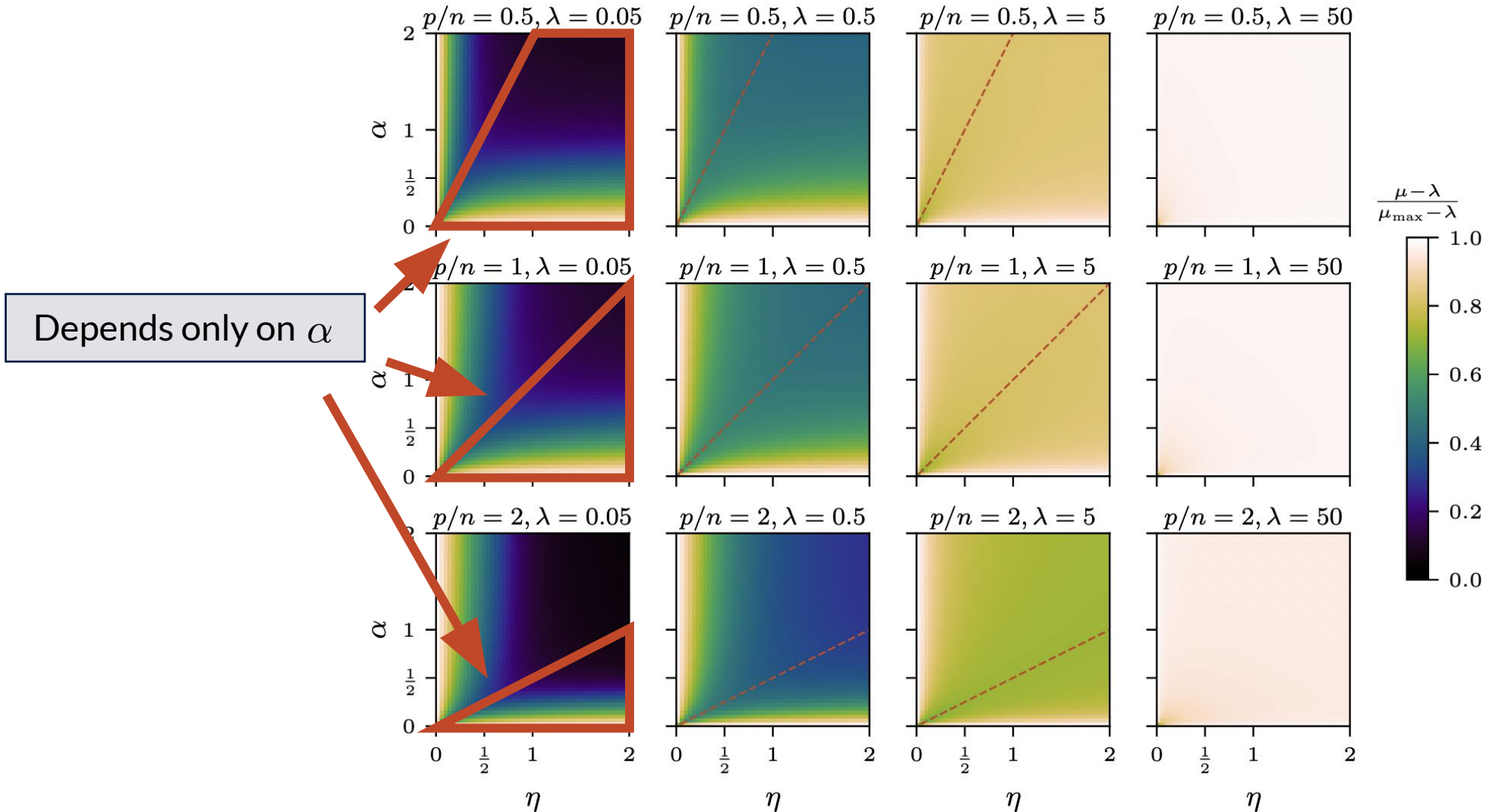
$$\lambda = \mu_{\text{D}} \left(1 - \frac{1}{q} \text{tr} \left[\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} \left(\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} + \mu_{\text{D}} \mathbf{I} \right)^{-1} \right] \right),$$

$$\frac{\lambda}{\theta} - 1 = \frac{m}{q} \left(\frac{\theta}{\mu_{\text{PD}}} - 1 \right), \quad \theta = \mu_{\text{PD}} \left(1 - \frac{1}{m} \text{tr} \left[\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} \left(\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} + \mu_{\text{PD}} \mathbf{I} \right)^{-1} \right] \right).$$

Sharp Phase Transition for Joint Sketching



Sharp Phase Transition for Joint Sketching



Sketching Makes Same Predictions as Ridge

Corollary 7. *If $\|\frac{1}{\sqrt{n}}\mathbf{y}\|_2$ is uniformly bounded in p and $\mathbf{w} \in \mathbb{C}^p$ is independent of \mathbf{S} and \mathbf{T} such that $\|\mathbf{w}\|_2$ is uniformly bounded in p , then for any continuous function $f: \mathbb{C} \rightarrow \mathbb{C}$, as $p \rightarrow \infty$,*

$$f(\mathbf{w}^H \hat{\boldsymbol{\beta}}_\psi) - f(\mathbf{w}^H \hat{\boldsymbol{\beta}}_{\mu_\psi}^{\text{ridge}}) \xrightarrow{\text{a.s.}} 0.$$

Sketching Makes Same Predictions as Ridge

Corollary 7. *If $\|\frac{1}{\sqrt{n}}\mathbf{y}\|_2$ is uniformly bounded in p and $\mathbf{w} \in \mathbb{C}^p$ is independent of \mathbf{S} and \mathbf{T} such that $\|\mathbf{w}\|_2$ is uniformly bounded in p , then for any continuous function $f: \mathbb{C} \rightarrow \mathbb{C}$, as $p \rightarrow \infty$,*

$$f(\mathbf{w}^H \hat{\boldsymbol{\beta}}_\psi) - f(\mathbf{w}^H \hat{\boldsymbol{\beta}}_{\mu_\psi}^{\text{ridge}}) \xrightarrow{\text{a.s.}} 0.$$

- Sketching asymptotically makes the **same prediction** as the **equivalent ridge** on **any single test point**

Sketching Makes Same Predictions as Ridge

Corollary 7. *If $\|\frac{1}{\sqrt{n}}\mathbf{y}\|_2$ is uniformly bounded in p and $\mathbf{w} \in \mathbb{C}^p$ is independent of \mathbf{S} and \mathbf{T} such that $\|\mathbf{w}\|_2$ is uniformly bounded in p , then for any continuous function $f: \mathbb{C} \rightarrow \mathbb{C}$, as $p \rightarrow \infty$,*

$$f(\mathbf{w}^H \hat{\boldsymbol{\beta}}_{\psi}) - f(\mathbf{w}^H \hat{\boldsymbol{\beta}}_{\mu_{\psi}}^{\text{ridge}}) \xrightarrow{\text{a.s.}} 0.$$

- Sketching asymptotically makes the **same prediction** as the **equivalent ridge** on **any single test point**

Why this qualifier?



Sketching Makes Same Predictions as Ridge

Corollary 7. If $\|\frac{1}{\sqrt{n}}\mathbf{y}\|_2$ is uniformly bounded in p and $\mathbf{w} \in \mathbb{C}^p$ is independent of \mathbf{Q} and \mathbf{T} such that $\|\mathbf{w}\|_2$ is uniformly bounded in p , then for any continuous function $f: \mathbb{C} \rightarrow \mathbb{C}$, as $p \rightarrow \infty$,

$$f(\mathbf{w}^H \hat{\boldsymbol{\beta}}_\psi) - f(\mathbf{w}^H \hat{\boldsymbol{\beta}}^{\text{ridge}}) \rightarrow 0$$

- Sketching asymptotically predicts the same prediction as the equivalent ridge on any single point.

Pointwise convergence does not imply uniform convergence

Why this qualifier?



Quadratic Metrics of Sketched Ensembles

- Ensemble of independent sketches: $\hat{\beta}_{\psi}^{\text{ens}} = \frac{1}{K} \sum_{k=1}^K \hat{\beta}_{\psi}^{(k)}$
- Quadratic error metrics: $\mathcal{E}_{\Psi}(\beta, \beta') = (\beta - \beta')^{\text{H}} \Psi (\beta - \beta')$
 - Includes test risk

Quadratic Metrics of Sketched Ensembles

Theorem 8. If $\Psi \in \mathbb{C}^{p \times p}$ is a positive semidefinite matrix and $\beta' \in \mathbb{C}^p$ a vector such that $\|\Psi\|_{\text{op}}$ and $\|\beta'\|_2$ are uniformly bounded in p and (Ψ, β) is independent of $(\mathbf{S}_k, \mathbf{T}_k)_{k=1}^K$, then for $\psi \in \{P, D\}$,

$$\begin{aligned} \mathcal{E}_\Psi(\hat{\beta}_P^{\text{ens}}, \beta') - \left(\mathcal{E}_\Psi(\hat{\beta}_{\mu_P}^{\text{ridge}}, \beta') + \frac{\mu'_P}{Kn} \mathbf{y}^H (\frac{1}{n} \mathbf{X} \mathbf{X}^H + \mu_P \mathbf{I})^{-2} \mathbf{y} \right) &\xrightarrow{\text{a.s.}} 0, \\ \mathcal{E}_\Psi(\hat{\beta}_D^{\text{ens}}, \beta') - \left(\mathcal{E}_\Psi(\hat{\beta}_{\mu_D}^{\text{ridge}}, \beta') + \frac{\mu'_D}{Kn} \mathbf{y}^H \mathbf{X} (\frac{1}{n} \mathbf{X}^H \mathbf{X} + \mu_D \mathbf{I})^{-2} \mathbf{X}^H \mathbf{y} \right) &\xrightarrow{\text{a.s.}} 0, \end{aligned}$$

where

$$\begin{aligned} \mu'_P &= \frac{\frac{1}{m} \text{tr} \left[\mu_P^3 (\frac{1}{n} \mathbf{X} \mathbf{X}^H + \mu_P \mathbf{I})^{-1} \frac{1}{n} \mathbf{X} \Psi \mathbf{X}^H (\frac{1}{n} \mathbf{X} \mathbf{X}^H + \mu_P \mathbf{I})^{-1} \right]}{\lambda + \frac{1}{m} \text{tr} \left[\mu_P^2 \frac{1}{n} \mathbf{X} \mathbf{X}^H (\frac{1}{n} \mathbf{X} \mathbf{X}^H + \mu_P \mathbf{I})^{-2} \right]}, \\ \mu'_D &= \frac{\frac{1}{q} \text{tr} \left[\mu_D^3 (\frac{1}{n} \mathbf{X}^H \mathbf{X} + \mu_D \mathbf{I})^{-1} \Psi (\frac{1}{n} \mathbf{X}^H \mathbf{X} + \mu_D \mathbf{I})^{-1} \right]}{\lambda + \frac{1}{q} \text{tr} \left[\mu_D^2 \frac{1}{n} \mathbf{X}^H \mathbf{X} (\frac{1}{n} \mathbf{X}^H \mathbf{X} + \mu_D \mathbf{I})^{-2} \right]}. \end{aligned}$$

Quadratic Metrics of Sketched Ensembles

Theorem 8. If $\Psi \in \mathbb{C}^{p \times p}$ is a positive semidefinite matrix and $\beta' \in \mathbb{C}^p$ a vector such that $\|\Psi\|_{\text{op}}$ and $\|\beta'\|_2$ are uniformly bounded in p and (Ψ, β) is independent of $(\mathbf{S}_k, \mathbf{T}_k)_{k=1}^K$, then for $\psi \in \{P, D\}$,

$$\mathcal{E}_{\Psi}(\hat{\beta}_{\psi}^{\text{ens}}, \beta') \xrightarrow{\text{a.s.}} \mathcal{E}_{\Psi}(\hat{\beta}_{\mu_{\psi}}^{\text{ridge}}, \beta') + O\left(\frac{\mu'_{\psi}}{K}\right)$$

where

$$\mu'_P = \frac{\frac{1}{m} \text{tr} \left[\mu_P^3 \left(\frac{1}{n} \mathbf{X} \mathbf{X}^H + \mu_P \mathbf{I} \right)^{-1} \frac{1}{n} \mathbf{X} \Psi \mathbf{X}^H \left(\frac{1}{n} \mathbf{X} \mathbf{X}^H + \mu_P \mathbf{I} \right)^{-1} \right]}{\lambda + \frac{1}{m} \text{tr} \left[\mu_P^2 \frac{1}{n} \mathbf{X} \mathbf{X}^H \left(\frac{1}{n} \mathbf{X} \mathbf{X}^H + \mu_P \mathbf{I} \right)^{-2} \right]},$$

$$\mu'_D = \frac{\frac{1}{q} \text{tr} \left[\mu_D^3 \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} + \mu_D \mathbf{I} \right)^{-1} \Psi \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} + \mu_D \mathbf{I} \right)^{-1} \right]}{\lambda + \frac{1}{q} \text{tr} \left[\mu_D^2 \frac{1}{n} \mathbf{X}^H \mathbf{X} \left(\frac{1}{n} \mathbf{X}^H \mathbf{X} + \mu_D \mathbf{I} \right)^{-2} \right]}.$$

Quadratic Metrics of Sketched Ensembles

Theorem 8. If $\Psi \in \mathbb{C}^{p \times p}$ is a positive semidefinite matrix and $\beta' \in \mathbb{C}^p$ a vector such that $\|\Psi\|_{\text{op}}$ and $\|\beta'\|_2$ are uniformly bounded in p and (Ψ, β) is independent of $(\mathbf{S}_k, \mathbf{T}_k)_{k=1}^K$, then for $\psi \in \{P, D\}$,

$$\mathcal{E}_{\Psi}(\hat{\beta}_{\psi}^{\text{ens}}, \beta') \xrightarrow{\text{a.s.}} \mathcal{E}_{\Psi}(\hat{\beta}_{\mu_{\psi}}^{\text{ridge}}, \beta') + O\left(\frac{\mu'_{\psi}}{K}\right)$$

where

$$\mu'_P = \frac{\frac{1}{m} \text{tr} \left[\mu_P^3 \left(\frac{1}{n} \mathbf{X} \mathbf{X}^{\text{H}} + \mu_P \mathbf{I} \right)^{-1} \right]}{\lambda + \frac{1}{m} \text{tr} \left[\mu_P^2 \frac{1}{n} \mathbf{X} \mathbf{X}^{\text{H}} \left(\frac{1}{n} \mathbf{X} \mathbf{X}^{\text{H}} + \mu_P \mathbf{I} \right)^{-1} \right]}$$

$$\mu'_D = \frac{\frac{1}{q} \text{tr} \left[\mu_D^3 \left(\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} + \mu_D \mathbf{I} \right)^{-1} \Psi \left(\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} + \mu_D \mathbf{I} \right)^{-1} \right]}{\lambda + \frac{1}{q} \text{tr} \left[\mu_D^2 \frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} \left(\frac{1}{n} \mathbf{X}^{\text{H}} \mathbf{X} + \mu_D \mathbf{I} \right)^{-1} \right]}$$

Well, ain't this a geometrical oddity.

$O\left(\frac{\mu'_{\psi}}{K}\right)$ from everywhere!



Quadratic Conclusions

$$\mathcal{E}_{\Psi}(\hat{\beta}_{\psi}^{\text{ens}}, \beta') \xrightarrow{\text{a.s.}} \mathcal{E}_{\Psi}(\hat{\beta}_{\mu_{\psi}}^{\text{ridge}}, \beta') + O\left(\frac{\mu'_{\psi}}{K}\right)$$

- For infinite K , sketched ensemble = ridge



Quadratic Conclusions

$$\mathcal{E}_{\Psi}(\hat{\beta}_{\psi}^{\text{ens}}, \beta') \xrightarrow{\text{a.s.}} \mathcal{E}_{\Psi}(\hat{\beta}_{\mu_{\psi}}^{\text{ridge}}, \beta') + O\left(\frac{\mu'_{\psi}}{K}\right)$$

- For infinite K , sketched ensemble = ridge



Spoiler:

Randomized least squares = ridge + noise

Randomized ensembles = ridge

Quadratic Conclusions

$$\mathcal{E}_{\Psi}(\hat{\beta}_{\psi}^{\text{ens}}, \beta') \xrightarrow{\text{a.s.}} \mathcal{E}_{\Psi}(\hat{\beta}_{\mu_{\psi}}^{\text{ridge}}, \beta') + O\left(\frac{\mu'_{\psi}}{K}\right)$$

- For infinite K , sketched ensemble = ridge
- For finite K , sketched ensemble is worse than ridge



Quadratic Conclusions

$$\mathcal{E}_{\Psi}(\hat{\beta}_{\psi}^{\text{ens}}, \beta') \xrightarrow{\text{a.s.}} \mathcal{E}_{\Psi}(\hat{\beta}_{\mu_{\psi}}^{\text{ridge}}, \beta') + O\left(\frac{\mu'_{\psi}}{K}\right)$$

- For infinite K , sketched ensemble = ridge
- For finite K , sketched ensemble is worse than ridge
 - Unless $\Psi \in \text{Range}(\mathbf{A})$, $\alpha > r(\mathbf{A})$, and $\mu = \lambda = 0!$



A Sketched Ensemble Efficiency Experiment

- Setup: fixed $O(p^2n)$ budget ensembles with $K = \lfloor \frac{1}{\alpha^2} \rfloor$, $r(\Sigma) = \frac{1}{2}$

A Sketched Ensemble Efficiency Experiment

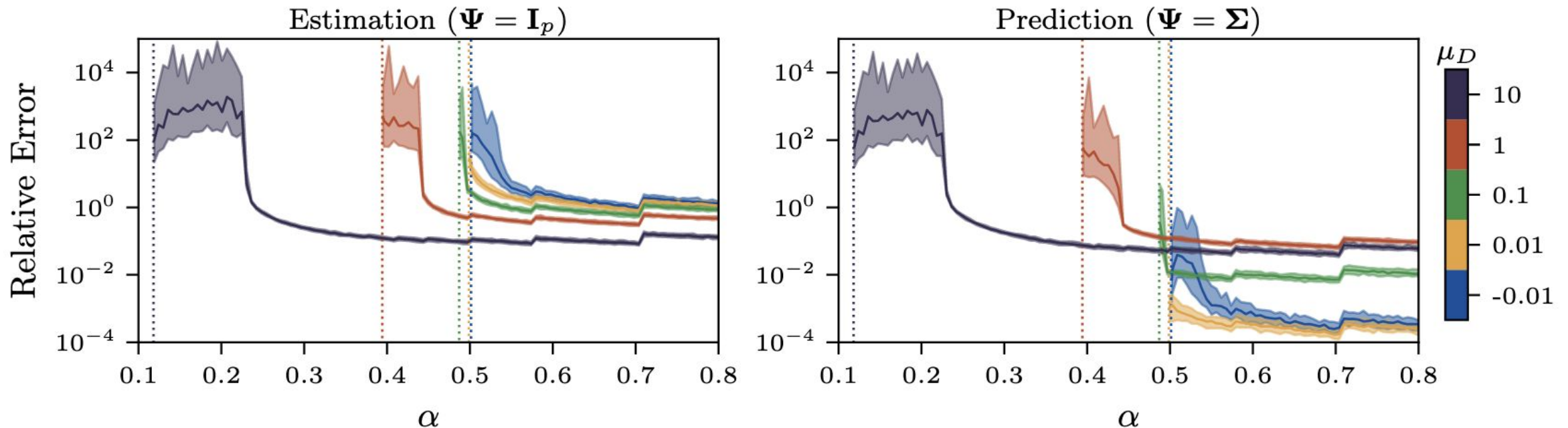
- Setup: **fixed** $O(p^2n)$ **budget** ensembles with $K = \lfloor \frac{1}{\alpha^2} \rfloor$, $r(\Sigma) = \frac{1}{2}$
- Fixed target μ_D , varying α , with λ uniquely determined

A Sketched Ensemble Efficiency Experiment

- Setup: fixed $O(p^2n)$ budget ensembles with $K = \lfloor \frac{1}{\alpha^2} \rfloor$, $r(\Sigma) = \frac{1}{2}$
- Fixed target μ_D , varying α , with λ uniquely determined
- Error: relative error $\frac{\mathcal{E}_\Psi(\hat{\beta}_D^{\text{ens}}, \hat{\beta}_{\mu_D}^{\text{ridge}})}{\mathcal{E}_\Psi(\mathbf{0}, \hat{\beta}_{\mu_D}^{\text{ridge}})}$

A Sketched Ensemble Efficiency Experiment

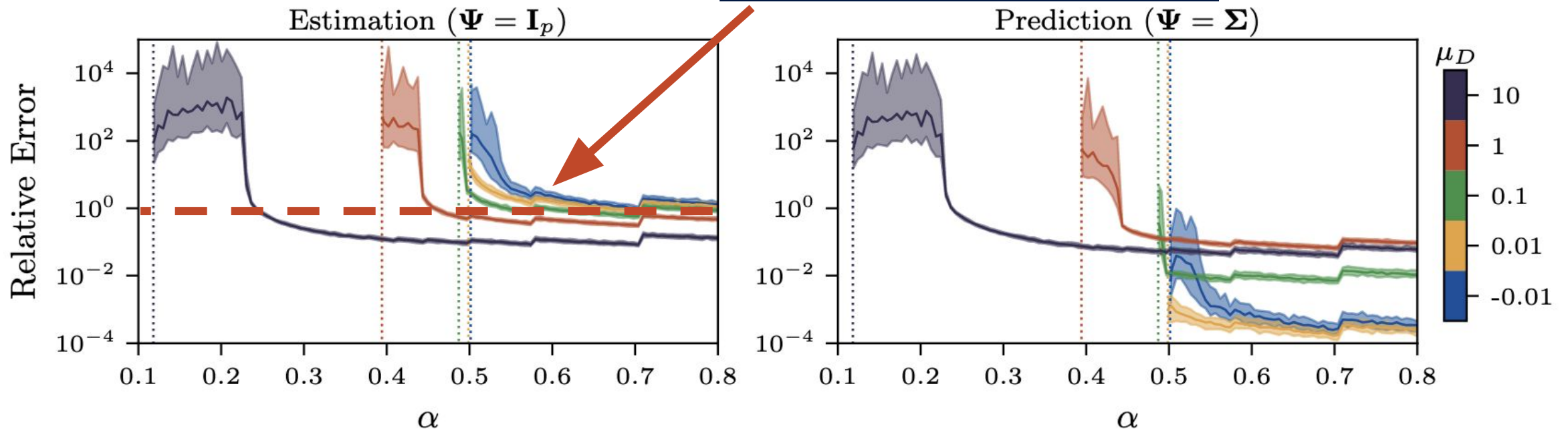
- Setup: fixed $O(p^2n)$ budget ensembles with $K = \lfloor \frac{1}{\alpha^2} \rfloor$, $r(\Sigma) = \frac{1}{2}$
- Fixed target μ_D , varying α , with λ uniquely determined
- Error: relative error $\frac{\mathcal{E}_\Psi(\hat{\beta}_D^{\text{ens}}, \hat{\beta}_{\mu_D}^{\text{ridge}})}{\mathcal{E}_\Psi(\mathbf{0}, \hat{\beta}_{\mu_D}^{\text{ridge}})}$



A Sketched Ensemble Efficiency Experiment

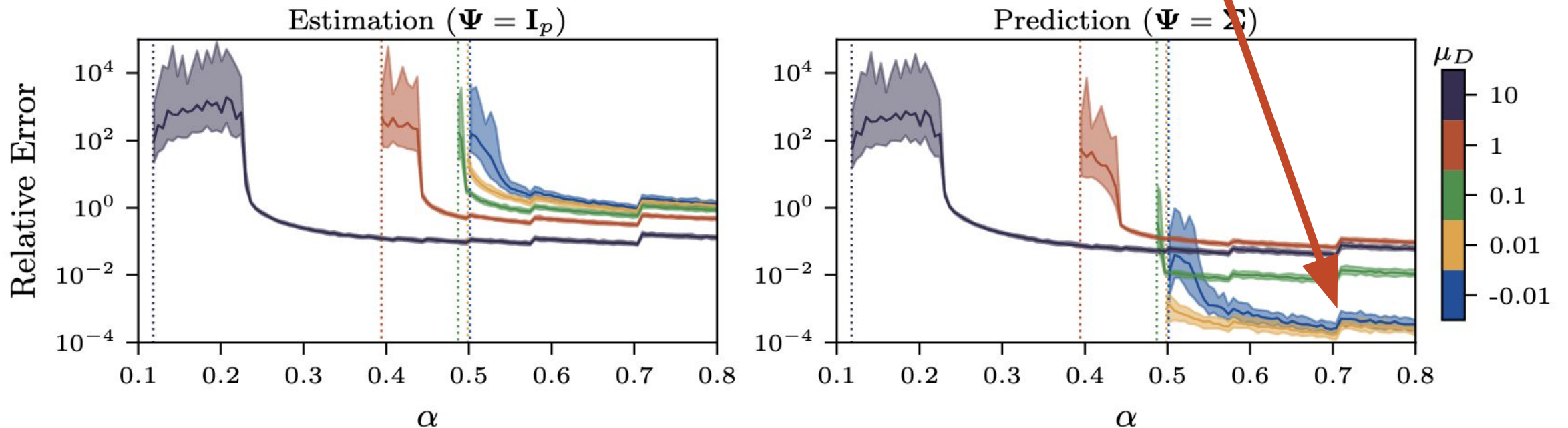
- Setup: fixed $O(p^2n)$ budget ensembles with $K = \lfloor \frac{1}{\alpha^2} \rfloor$, $r(\Sigma) = \frac{1}{2}$
- Fixed target μ_D , varying α , with λ uniquely determined
- Error: relative error $\frac{\mathcal{E}_\Psi(\hat{\beta}_D^{\text{ens}}, \hat{\beta}_{\mu_D}^{\text{ridge}})}{\mathcal{E}_\Psi(\mathbf{0}, \hat{\beta}_{\mu_D}^{\text{ridge}})}$

Non-vanishing estimation error

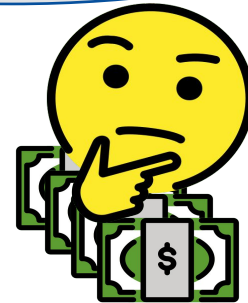


A Sketched Ensemble Efficiency Experiment

- Setup: fixed $O(p^2n)$ budget ensembles with $K = \lfloor \frac{1}{\alpha^2} \rfloor$, $r(\Sigma) = \frac{1}{2}$
- Fixed target μ_D , varying α , with λ uniquely determined
- Error: relative error $\frac{\mathcal{E}_\Psi(\hat{\beta}_D^{\text{ens}}, \hat{\beta}_{\mu_D}^{\text{ridge}})}{\mathcal{E}_\Psi(\mathbf{0}, \hat{\beta}_{\mu_D}^{\text{ridge}})}$



What if I use a better/faster
sketch than i.i.d.?



A Free Sketching Conjecture

Conjecture 9. *Let $\mathbf{S} \in \mathbb{C}^{p \times q}$ be a norm-preserving sketch such that $\mathbf{S}\mathbf{S}^H$ and \mathbf{A} converge almost surely to operators that are free with respect to the average trace $\frac{1}{p}\text{tr}[\cdot]$. Then there exists a monotonic mapping $\lambda \mapsto \gamma$ such that*

$$\mathbf{S}(\mathbf{S}^H\mathbf{A}\mathbf{S} + \lambda\mathbf{I}_q)^{-1}\mathbf{S}^H \simeq (\mathbf{A} + \gamma\mathbf{I}_p)^{-1}.$$

A Free Sketching Conjecture

Conjecture 9. *Let $\mathbf{S} \in \mathbb{C}^{p \times q}$ be a norm-preserving sketch such that $\mathbf{S}\mathbf{S}^H$ and \mathbf{A} converge almost surely to operators that are free with respect to the average trace $\frac{1}{p}\text{tr}[\cdot]$. Then there exists a monotonic mapping $\lambda \mapsto \gamma$ such that*

$$\mathbf{S}(\mathbf{S}^H\mathbf{A}\mathbf{S} + \lambda\mathbf{I}_q)^{-1}\mathbf{S}^H \simeq (\mathbf{A} + \gamma\mathbf{I}_p)^{-1}.$$

- Includes i.i.d. sketching and more
 - Orthogonal sketching
 - Efficient sketches like CountSketch, FJLT, SRHT?

A Free Sketching Conjecture

Conjecture 9. *Let $\mathbf{S} \in \mathbb{C}^{p \times q}$ be a norm-preserving sketch such that $\mathbf{S}\mathbf{S}^H$ and \mathbf{A} converge almost surely to operators that are free with respect to the average trace $\frac{1}{p}\text{tr}[\cdot]$. Then there exists a monotonic mapping $\lambda \mapsto \gamma$ such that*

$$\mathbf{S}(\mathbf{S}^H\mathbf{A}\mathbf{S} + \lambda\mathbf{I}_q)^{-1}\mathbf{S}^H \simeq (\mathbf{A} + \gamma\mathbf{I}_p)^{-1}.$$

- Includes i.i.d. sketching and more
 - Orthogonal sketching
 - Efficient sketches like CountSketch, FJLT, SRHT?
- Spectrum of sketch controls $\lambda \mapsto \gamma$

A Free Sketching Conjecture

Conjecture 9. *Let $\mathbf{S} \in \mathbb{C}^{p \times q}$ be a norm-preserving sketch such that $\mathbf{S}\mathbf{S}^H$ and \mathbf{A} converge almost surely to operators that are free with respect to the average trace $\frac{1}{p}\text{tr}[\cdot]$. Then there exists a monotonic mapping $\lambda \mapsto \gamma$ such that*

$$\mathbf{S}(\mathbf{S}^H\mathbf{A}\mathbf{S} + \lambda\mathbf{I}_q)^{-1}\mathbf{S}^H \simeq (\mathbf{A} + \gamma\mathbf{I}_p)^{-1}.$$

- Includes i.i.d. sketching and more
 - Orthogonal sketching
 - Efficient sketches like CountSketch, FJLT, SRHT?
- Spectrum of sketch controls $\lambda \mapsto \gamma$
- Higher order equivalences naturally follow

A Free Sketching Conjecture

Conjecture 9. *Let $\mathbf{S} \in \mathbb{C}^{p \times q}$ be a norm-preserving sketch such that $\mathbf{S}\mathbf{S}^H$ and \mathbf{A} converge almost surely to operators that are free with respect to the average trace $\frac{1}{p}\text{tr}[\cdot]$. Then there exists a monotonic mapping $\lambda \mapsto \gamma$ such that*

$$\mathbf{S}(\mathbf{S}^H\mathbf{A}\mathbf{S} + \lambda\mathbf{I}_q)^{-1}\mathbf{S}^H \simeq (\mathbf{A} + \gamma\mathbf{I}_p)^{-1}.$$

- Includes i.i.d. sketching and more
 - Orthogonal sketching
 - Efficient sketches like CountSketch, FJLT, SRHT?
- Spectrum of sketch controls $\lambda \mapsto \gamma$
- Higher order equivalences naturally follow

Orthogonal Sketching

Conjecture 10. For $q \leq p$ let $\sqrt{\frac{q}{p}}\mathbf{Q} \in \mathbb{C}^{p \times q}$ be a Haar-distributed matrix with orthonormal columns. Then

$$\mathbf{Q}(\mathbf{Q}^H \mathbf{A} \mathbf{Q} + \lambda \mathbf{I}_q)^{-1} \mathbf{Q}^H \simeq (\mathbf{A} + \gamma \mathbf{I}_p)^{-1},$$

where γ is the most positive solution to

$$\frac{1}{p} \text{tr} \left[(\mathbf{A} + \gamma \mathbf{I}_p)^{-1} \right] (\gamma - \alpha \lambda) = 1 - \alpha.$$

Furthermore, for $\mu > 0$ applied to the same $(\mathbf{A}, \alpha, \lambda)$, we have $\gamma < \mu$.

Orthogonal Sketching

Conjecture 10. For $q \leq p$ let $\sqrt{\frac{q}{p}}\mathbf{Q} \in \mathbb{C}^{p \times q}$ be a Haar-distributed matrix with orthonormal columns. Then

$$\mathbf{Q}(\mathbf{Q}^H \mathbf{A} \mathbf{Q} + \lambda \mathbf{I}_q)^{-1} \mathbf{Q}^H \simeq (\mathbf{A} + \gamma \mathbf{I}_p)^{-1},$$

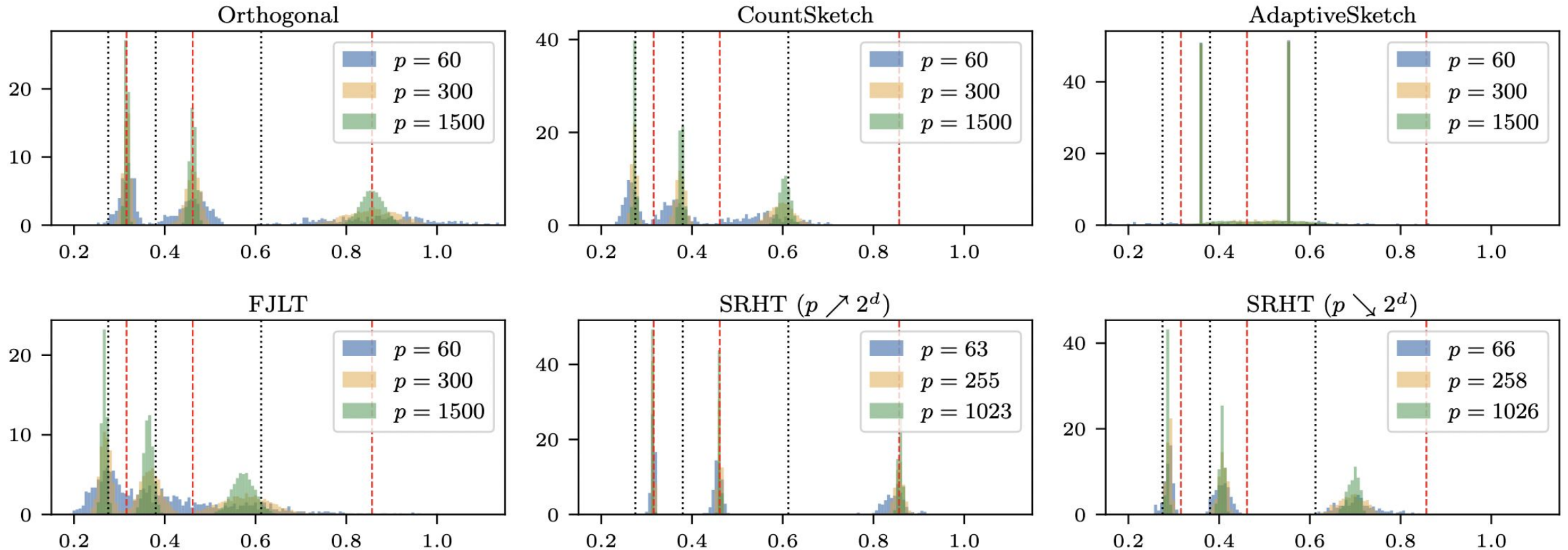
where γ is the most positive solution to

$$\frac{1}{p} \text{tr} \left[(\mathbf{A} + \gamma \mathbf{I}_p)^{-1} \right] (\gamma - \alpha \lambda) = 1 - \alpha.$$

Furthermore, for $\mu > 0$ applied to the same $(\mathbf{A}, \alpha, \lambda)$, we have $\gamma < \mu$.

- Same form as i.i.d. sketching, but with less regularization

Equivalence for Sketches Used in Practice?



- Early work:
 - Hints of **deep connection** between ensembles and ridge
 - Tuned ensembles with subsampling achieve same risk as **optimal ridge**

- Early work:

- Hints of **deep connection** between ensembles and ridge
 - Tuned ensembles with subsampling achieve same risk as **optimal ridge**

- Current work:

- **Asymptotic equivalence** between **random projections** and ridge
 - **Ridge equivalence** on a weak level **even for single learners**
 - Convergence in quadratic metrics to ridge regression for **ensembles**
 - Sufficiently large sketches enable **accurate ridgeless regression** even without ensembles

- Early work:

- Hints of **deep connection** between ensembles and ridge
 - Tuned ensembles with subsampling achieve same risk as **optimal ridge**

- Current work:

- **Asymptotic equivalence** between **random projections** and ridge
 - **Ridge equivalence** on a weak level **even for single learners**
 - Convergence in quadratic metrics to ridge regression for **ensembles**
 - Sufficiently large sketches enable **accurate ridgeless regression** even without ensembles

- Future work:

- **More asymptotic equivalences**
 - Generalized cross-validation with sketching
 - **General linear models** via leave-one-dimension-out
 - Asymptotics of **PCA**

Questions?

Anisotropic Sketching

Corollary 11. *Let \mathbf{W} be an invertible $p \times p$ positive semidefnite matrix, either deterministic or random but independent of \mathbf{S} with $\limsup \|\mathbf{W}\|_{\text{op}} < \infty$. Let $\tilde{\mathbf{S}} = \mathbf{W}^{1/2}\mathbf{S}$. Then for each $\lambda > -\liminf \lambda_{\min}^+(\tilde{\mathbf{S}}^\top \mathbf{A} \tilde{\mathbf{S}})$ as $p, q \rightarrow \infty$ such that $0 < \liminf \frac{q}{p} \leq \limsup \frac{q}{p} < \infty$,*

$$\tilde{\mathbf{S}}(\tilde{\mathbf{S}}^\top \mathbf{A} \tilde{\mathbf{S}} + \lambda \mathbf{I}_q)^{-1} \tilde{\mathbf{S}}^\top \simeq (\mathbf{A} + \mu \mathbf{W}^{-1})^{-1},$$

where μ most positive solution to

$$\lambda = \mu \left(1 - \frac{1}{q} \text{tr} \left[\mathbf{A} (\mathbf{A} + \mu \mathbf{W}^{-1})^{-1} \right] \right).$$

Equivalence for PCA

Theorem 12. *Let $\mathbf{X} = \mathbf{Z}\Sigma^{1/2}$ and $\Sigma \in \mathbb{C}^{p \times p}$ have eigenvalue decomposition $\mathbf{U}\mathbf{D}\mathbf{U}^H$, and let $\Pi_{\mathcal{A}}$ be the projection operator of the principal eigenspace corresponding to a set of eigenvalues \mathcal{A} of the matrix $\frac{1}{n}\mathbf{X}^H\mathbf{X}$. Then there exists a family of measures μ_{σ^2} for all $\sigma^2 \geq 0$ such that for any $\mathcal{A} \subseteq \mathbb{R}_{\geq 0}$, in the limit as $p \rightarrow \infty$,*

$$\Pi_{\mathcal{A}} \simeq \mathbf{U}\Lambda\mathbf{U}^H,$$

where Λ is a diagonal matrix defined for by

$$[\Lambda]_{ii} = \mu_{[\Sigma]_{ii}}(\mathcal{A}).$$