

## Abstract

We study **sketched ridge regression ensembles** built from the general class of sketches **asymptotically free** from the data

- We precisely characterize **asymptotic risk**
- We prove that **generalized cross-validation (GCV)** provides consistent risk estimation for feature sketching ensembles
- We show that GCV also provides consistent **distribution estimation** enabling **prediction intervals**
- We employ an **ensemble trick** for efficiently estimating unsketched ridge regression risk

## Freely sketched ridge ensembles

**Given:** data  $(\mathbf{X}, \mathbf{y}) \in \mathbb{R}^{n \times p} \times \mathbb{R}$ , feature sketches  $\mathbf{S}_1, \dots, \mathbf{S}_K \in \mathbb{R}^{p \times q}$ , and the ensemble predictor at regularization level  $\lambda$

$$\hat{\beta}_\lambda^{\text{ens}} = \frac{1}{K} \sum_{k=1}^K \frac{1}{n} \mathbf{S}_k \left( \frac{1}{n} \mathbf{S}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{S}_k + \lambda \mathbf{I}_q \right)^{-1} \mathbf{S}_k^\top \mathbf{X}^\top \mathbf{y}$$

where  $\mathbf{S}_k \mathbf{S}_k^\top$  is **asymptotically free** from the data  $\frac{1}{n} \mathbf{X}^\top \mathbf{X}$

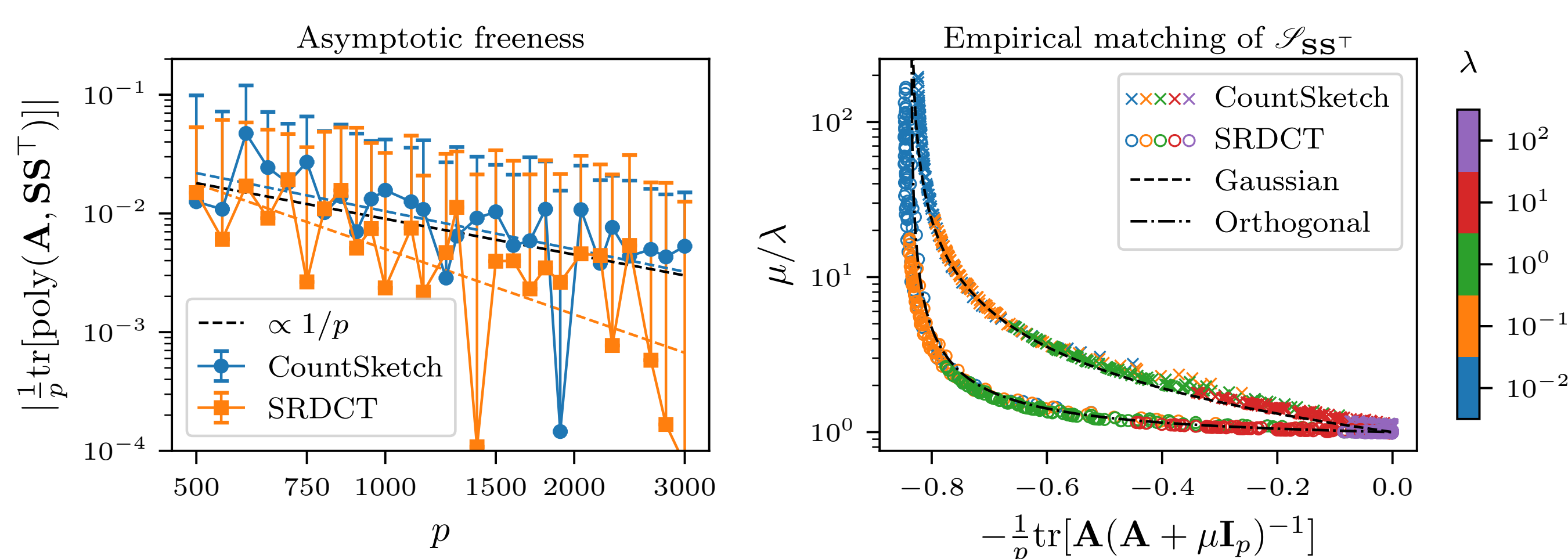
Two matrices  $\mathbf{A}$  and  $\mathbf{B}$  are almost surely **asymptotically free** if all mixed alternating products of centered polynomials are also centered:

$$\overline{\text{tr}}[p_1(\mathbf{A})p_2(\mathbf{B}) \dots p_{L-1}(\mathbf{A})p_L(\mathbf{B})] \xrightarrow{\text{a.s.}} 0.$$

Examples of known asymptotically free sketches:

- Independent:  $[\mathbf{S}_k]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$ , zero mean, bounded moments
- Rotationally invariant:  $\mathbf{S}_k = \mathbf{U}_k \mathbf{Q}_k$  with  $\mathbf{U}_k$  Haar-distributed
- Randomized Fourier transform:  $\mathbf{S}_k = \mathbf{D}_k \Phi_{DFT} \hat{\mathbf{S}}_k$

We provide **empirical support** for freeness for practical sketches.



## Generalized cross-validation (GCV)

**Goal:** estimate the joint distribution of true labels and predictions  $(y_0, x_0^\top \hat{\beta}_\lambda^{\text{ens}})$  in order to estimate risk  $T(\hat{\beta}_\lambda^{\text{ens}}) = \mathbb{E} \left[ t(y_0, x_0^\top \hat{\beta}_\lambda^{\text{ens}}) \right]$ .

Ensemble predictions are **linear smoothers**  $\mathbf{X} \hat{\beta}_\lambda^{\text{ens}} = \mathbf{L}_\lambda^{\text{ens}} \mathbf{y}$  for

$$\mathbf{L}_\lambda^{\text{ens}} = \frac{1}{K} \sum_{k=1}^K \mathbf{L}_\lambda^k = \frac{1}{n} \mathbf{X} \mathbf{S}_k \left( \frac{1}{n} \mathbf{S}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{S}_k + \lambda \mathbf{I}_q \right)^{-1} \mathbf{S}_k^\top \mathbf{X}^\top,$$

giving us the **GCV-corrected empirical distribution**

$$\hat{P}_\lambda^{\text{ens}} = \frac{1}{n} \sum_{i=1}^n \delta \left\{ \left( y_i, \frac{x_i^\top \hat{\beta}_\lambda^{\text{ens}} - \frac{1}{n} \text{tr}[\mathbf{L}_\lambda^{\text{ens}}] y_i}{1 - \frac{1}{n} \text{tr}[\mathbf{L}_\lambda^{\text{ens}}]} \right) \right\}.$$

We plug in  $\hat{P}_\lambda^{\text{ens}}$  to obtain risk estimators

$$\hat{T}(\hat{\beta}_\lambda^{\text{ens}}) = \int t(y, z) d\hat{P}_\lambda^{\text{ens}}(y, z) \quad \text{and} \quad \hat{R}(\hat{\beta}_\lambda^{\text{ens}}) = \int (y-z)^2 d\hat{P}_\lambda^{\text{ens}}(y, z).$$

## Sketched ensemble risk

Free sketches satisfy an **asymptotic equivalence**:

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

where  $\mu \simeq \lambda \mathcal{L}_{\text{SS}^\top} \left( -\frac{1}{p} \text{tr} \left[ \mathbf{S}^\top \mathbf{A} \mathbf{S} (\mathbf{S}^\top \mathbf{A} \mathbf{S} + \lambda \mathbf{I}_q)^{-1} \right] \right)$ .

## Theoretical results

**Theorem 1.** For any free sketches  $\mathbf{S}_k$ ,

$$R(\hat{\beta}_\lambda^{\text{ens}}) \simeq R(\hat{\beta}_\mu^{\text{ridge}}) + \frac{\mu' \Delta}{K} \quad \text{and} \quad \hat{R}(\hat{\beta}_\lambda^{\text{ens}}) \simeq \hat{R}(\hat{\beta}_\mu^{\text{ridge}}) + \frac{\mu'' \Delta}{K},$$

where  $\hat{\beta}_\mu^{\text{ridge}} = \left( \frac{1}{n} \mathbf{X}^\top \mathbf{X} + \mu \mathbf{I}_p \right)^{-1} \frac{1}{n} \mathbf{X}^\top \mathbf{y}$ .

**Theorem 2.** Under random data assumptions on  $\mathbf{X}$  and  $\mathbf{y}$ ,

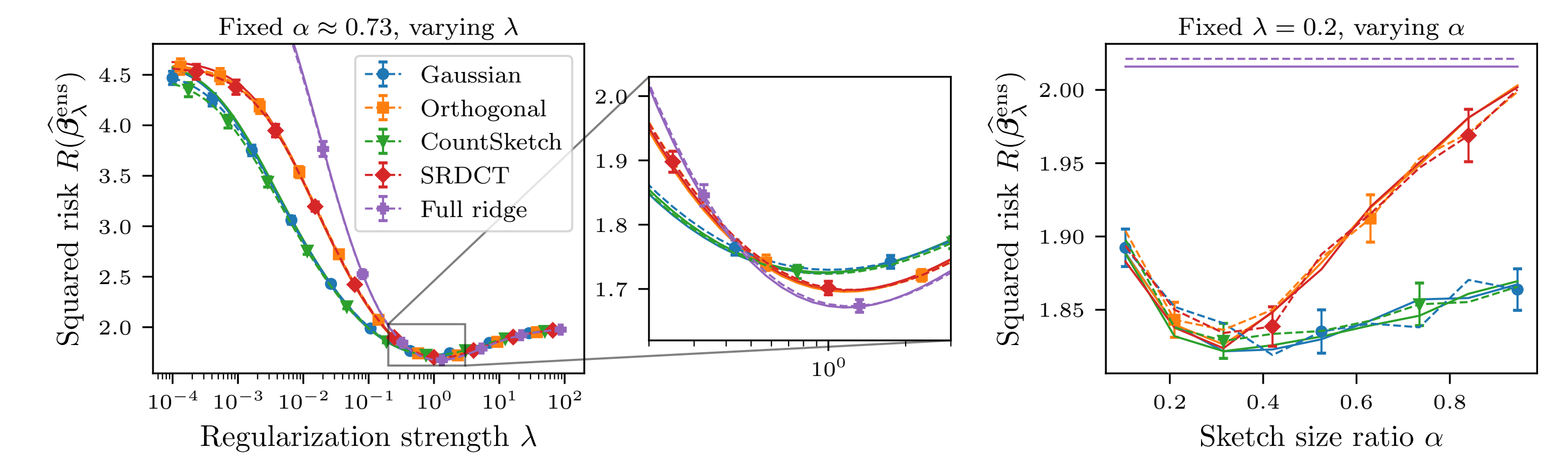
$$\mu' \simeq \mu'', \quad \text{and therefore} \quad \hat{R}(\hat{\beta}_\lambda^{\text{ens}}) \simeq R(\hat{\beta}_\lambda^{\text{ens}}).$$

**Theorem 3.** For any  $t$  pseudo-Lipshitz of order 2,

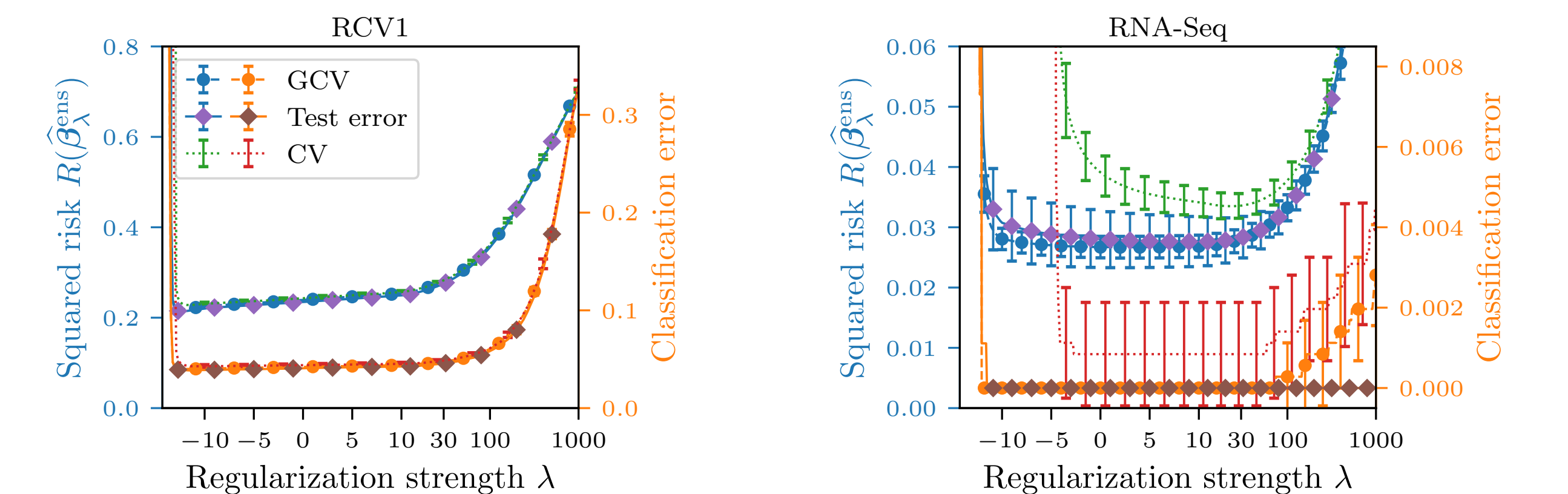
$$\hat{T}(\hat{\beta}_\lambda^{\text{ens}}) \simeq T(\hat{\beta}_\lambda^{\text{ens}}), \quad \text{and therefore} \quad \hat{P}_\lambda^{\text{ens}} \stackrel{2}{\simeq} P_\lambda^{\text{ens}}.$$

## Empirical results

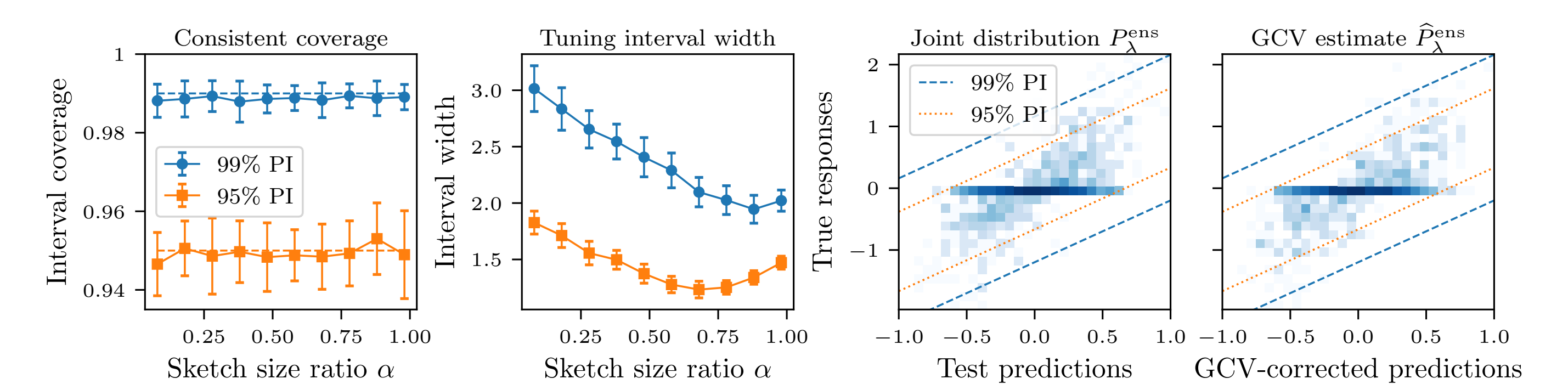
### Consistency across sketching families



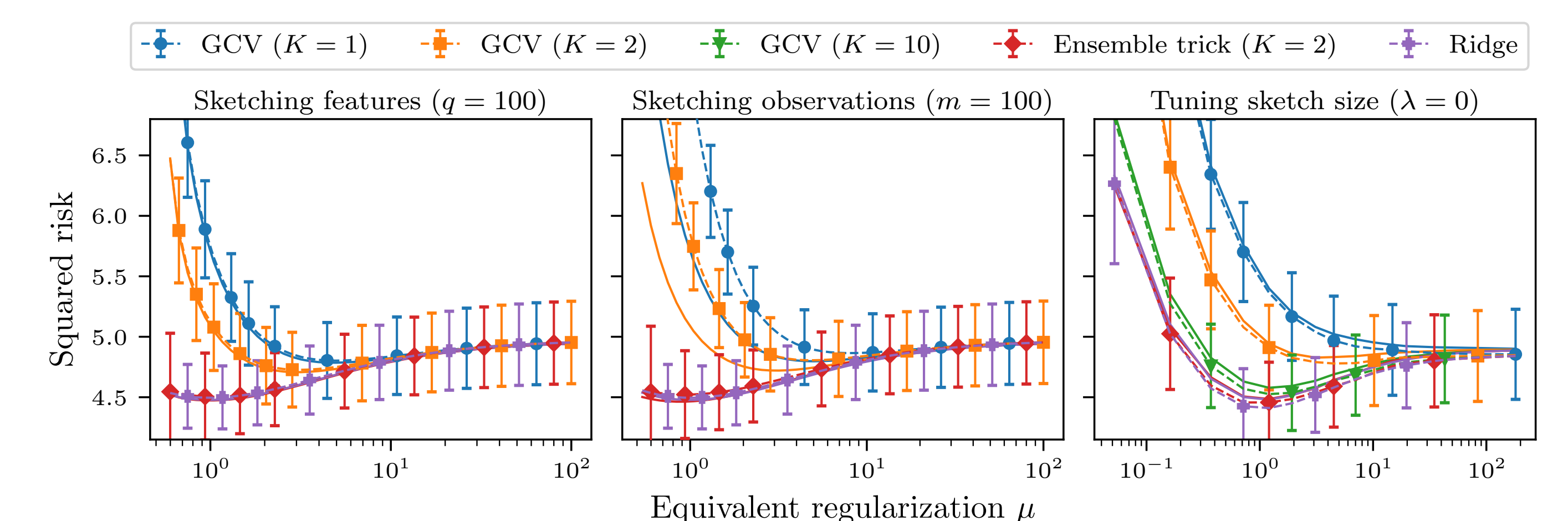
### Risk estimation for real data



### Consistent prediction intervals



### Fast estimation of unsketched risk using ensemble trick



## Acknowledgements

This collaboration was partially supported by Office of Naval Research MURI grant N00014-20-1-2787. DL was supported by Army Research Office grant 2003514594.