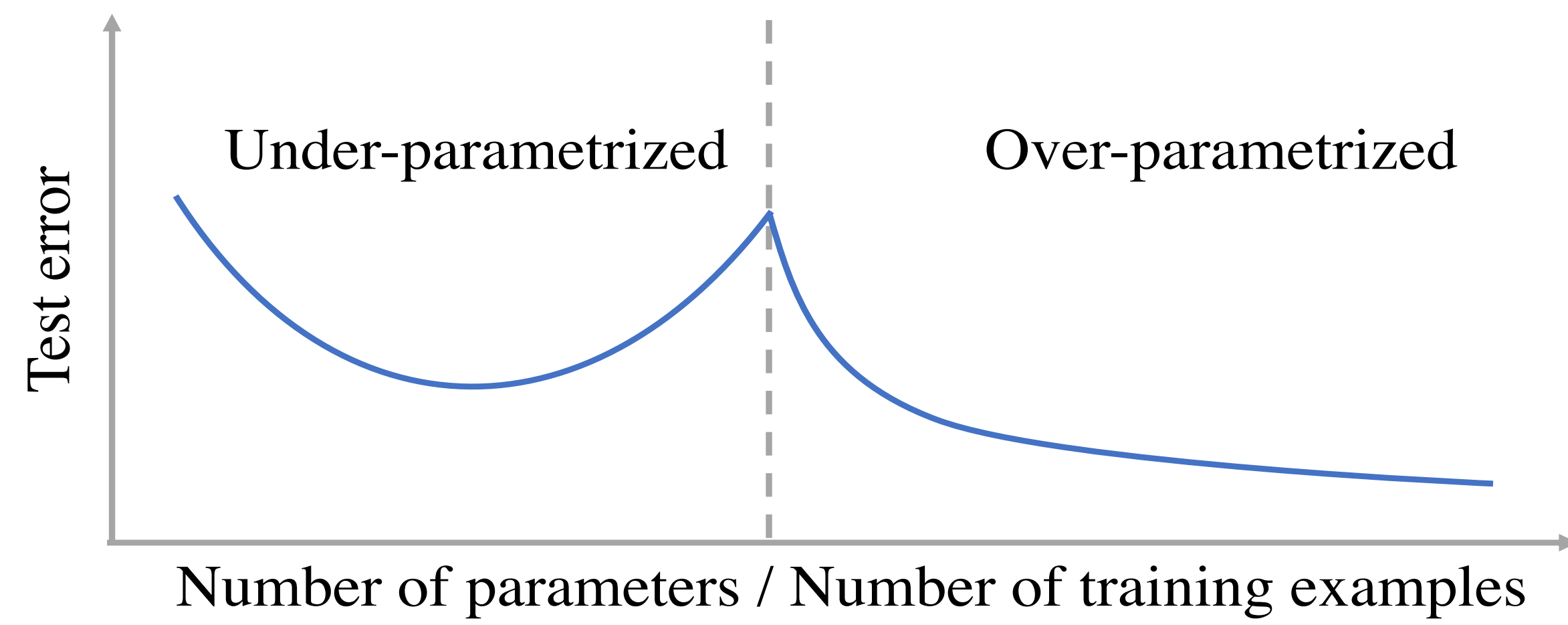
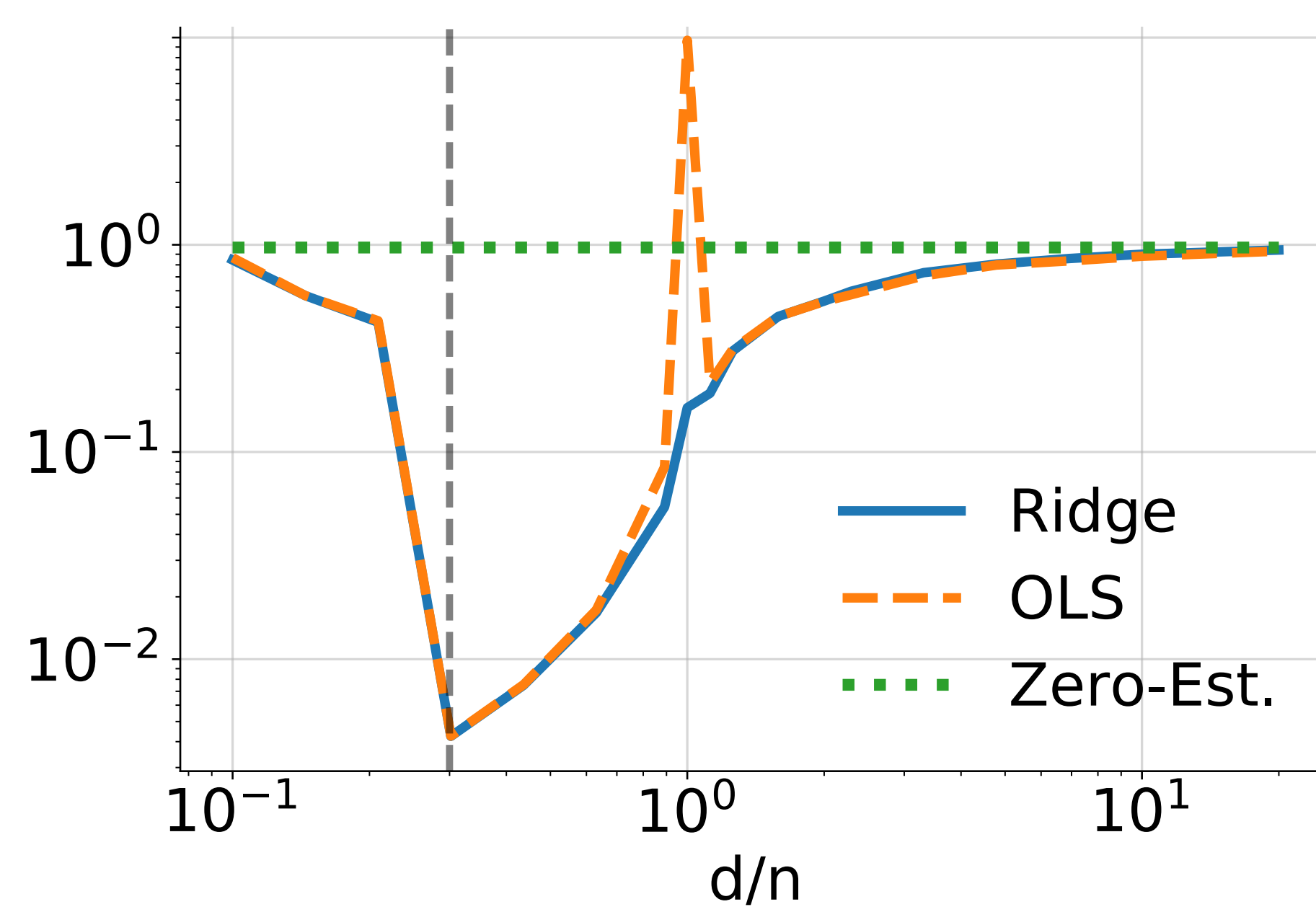


Raaz Dwivedi, Chandan Singh,  
Bin Yu, Martin Wainwright

## INTRO: DOUBLE-DESCENT



- Bad estimators don't exhibit U-shaped bias-variance tradeoff, even in low-dimensions; why should we expect that from OLS?

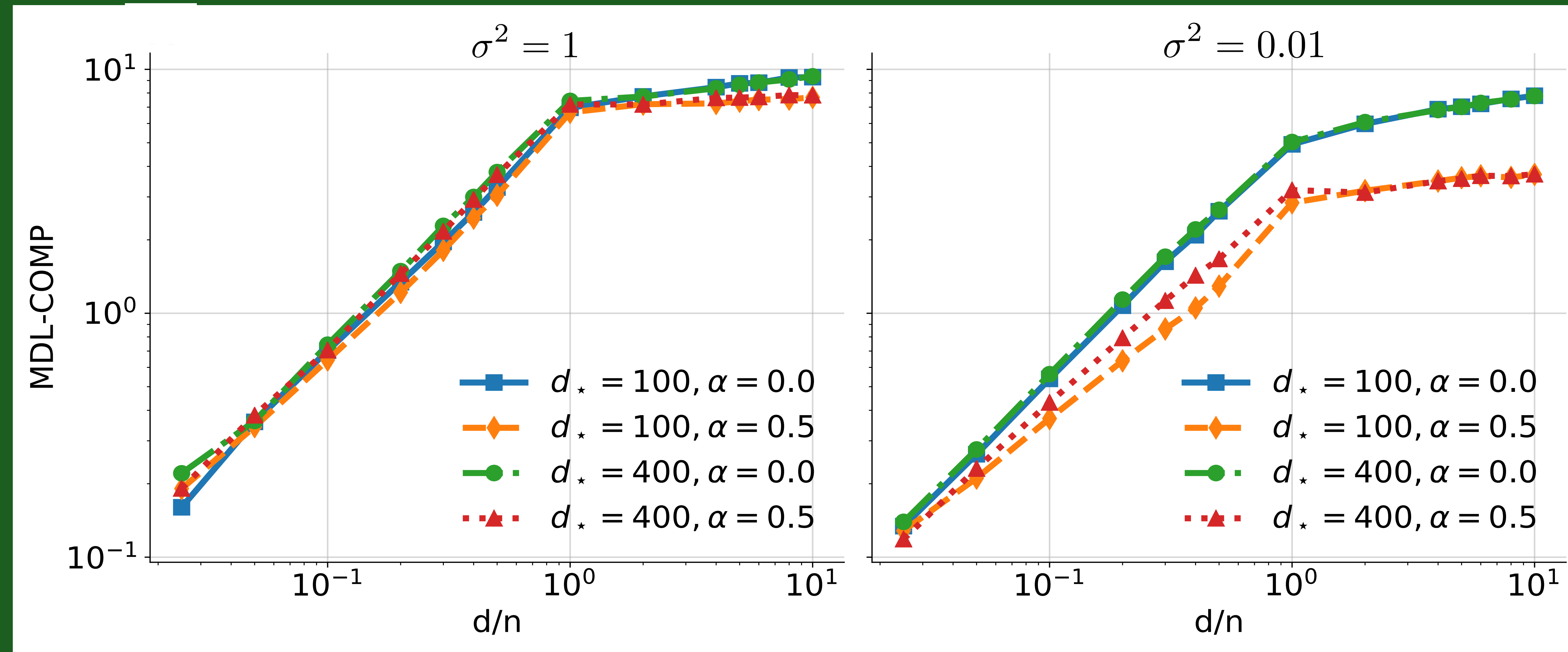


- Why is complexity = # parameters in overparameterized models?

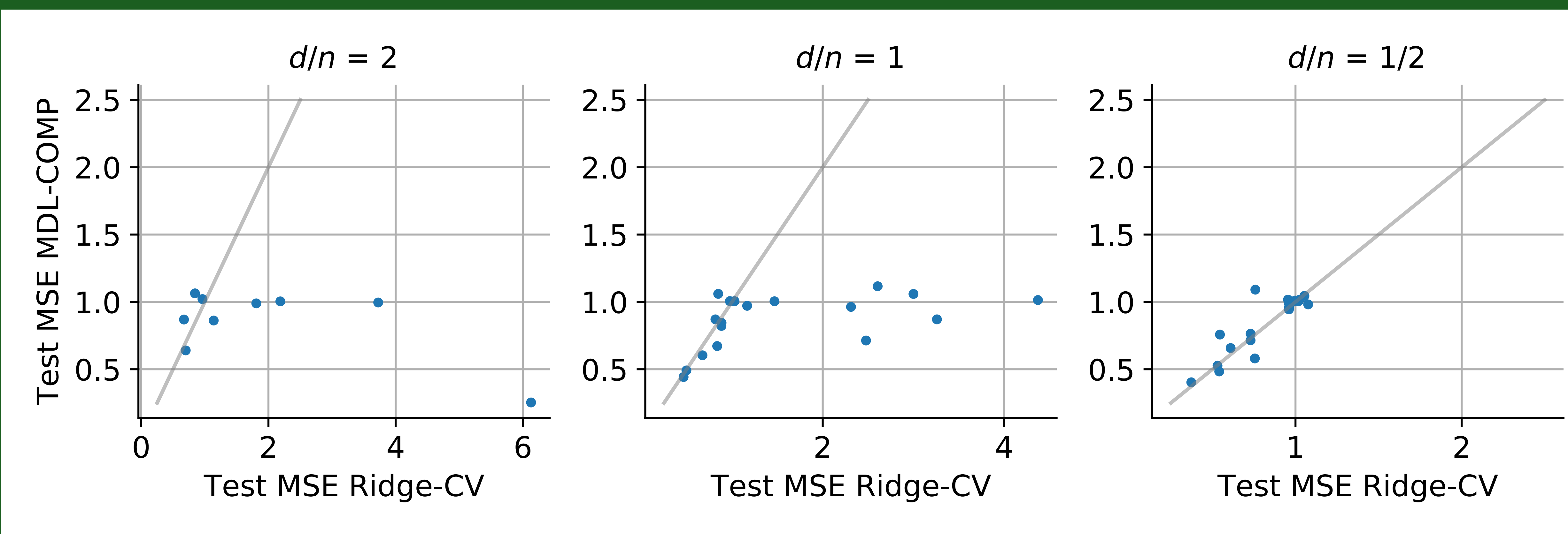
## REVISIT COMPLEXITY

1. Build on the normalized maximum likelihood (NML) principle in MDL: luckiness NML codes induced by ridge estimators to account for degeneracy in high-dimensions
2. Define MDL-COMP as optimal excess codelength over the LNML codes
3. Find a measure that depends on covariate design, and signal strength and is not mere parameter count

## MDL-Complexity can grow logarithmic with parameters for overparameterized models



& inform hyper-parameter tuning in limited-data settings, using just a single fold computation



Bottom plot: 19 Real datasets from Penn ML Benchmark (CV = cross-validation)  
Top plot: Gaussian synthetic data

## RESULTS

- MDL-COMP has a non-linear scaling
- Informs in-sample generalization for linear and kernel methods
- Provides competitive performance to cross-validation (CV) in out-of-sample

### MDL-COMP Expressions

• Linear models:  $y_i = x_i^\top \theta^* + \mathcal{N}(0, \sigma^2)$

$$\text{MDL-COMP} = \frac{1}{n} \sum_{i=1}^d \log \left( \rho_i + \frac{\sigma^2}{w_i^2} \right)$$

$w = \mathbf{U}^\top \theta^*$ ,  $\mathbf{U}, \{\rho_i\} \sim$  eigenvectors and eigenvalues of  $\mathbf{X}^\top \mathbf{X}$

- Gaussian design:  $d_* = \# \text{true feats}$ ,  $d = \# \text{fitted feats}$ ,  $n = \# \text{samples}$ ; with  $d_* = n$ 

$$\approx \begin{cases} \frac{d}{n} \log(1 + d_*/\|\theta^*\|^2) & \text{when } d < n \\ \log[d(1/\|\theta^*\|^2 + 1/n)] & \text{when } d > n. \end{cases}$$
- For the plots on the top:  $\theta^* \sim U(\mathbb{S}^{d_*-1})$ 

$$x_i \sim \mathcal{N}(0, \text{diag}(1, 2^{-\alpha}, \dots, d^{-\alpha}))$$

- Kernel methods:  $y_i = f^*(x_i) + \mathcal{N}(0, \sigma^2)$

$$\text{MDL-COMP} = \inf_{\lambda} \left( \frac{\lambda \|f^*\|_{\mathbb{H}}^2}{2n\sigma^2} + \frac{1}{2n} \sum_{i=1}^n \log \left( 1 + \frac{\rho_i}{\lambda} \right) \right)$$

$\{\rho_i\} \sim$  eigenvalues of  $\mathbf{K} = (\text{kernel}(x_i, x_j))_{i,j=1}^n$

### MDL-COMP based hyper-parameter tuning

- Kernel ridge regression:

$$\min_{\lambda} \frac{\|\mathbf{K}\hat{\theta} - \mathbf{y}\|^2 + \lambda \hat{\theta}^\top \mathbf{K} \hat{\theta}}{\sigma^2} + \sum_{i=1}^d \log \left( 1 + \frac{\rho_i}{\lambda} \right)$$

$\hat{\theta} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$ ,  $\{\rho_i\} \sim$  eigenvalues of  $\mathbf{K}$

- For linear model:  $\mathbf{K} = \mathbf{X}\mathbf{X}^\top$