

Generalization Error Rates in Kernel Regression: The Crossover from the Noiseless to Noisy Regime

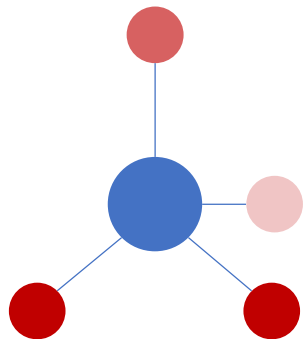
Hugo Cui,¹ Bruno Loureiro,² Florent Krzakala,² and Lenka Zdeborová¹

¹*SPOC, EPFL*

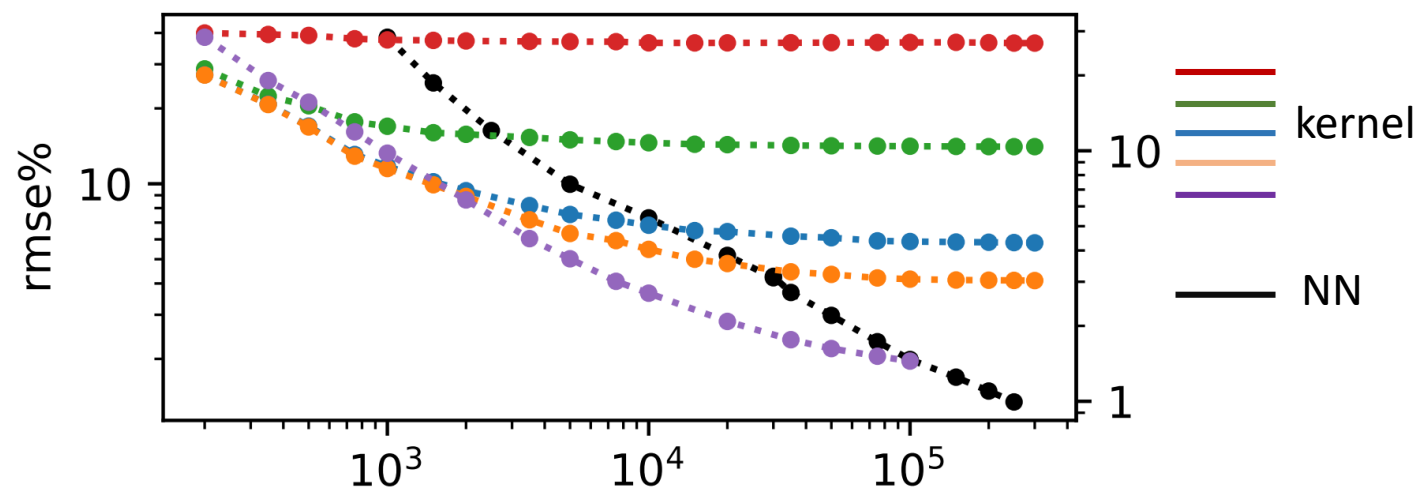
²*IDePHICS lab, EPFL*

arXiv: 2105.15004

Why study kernels?

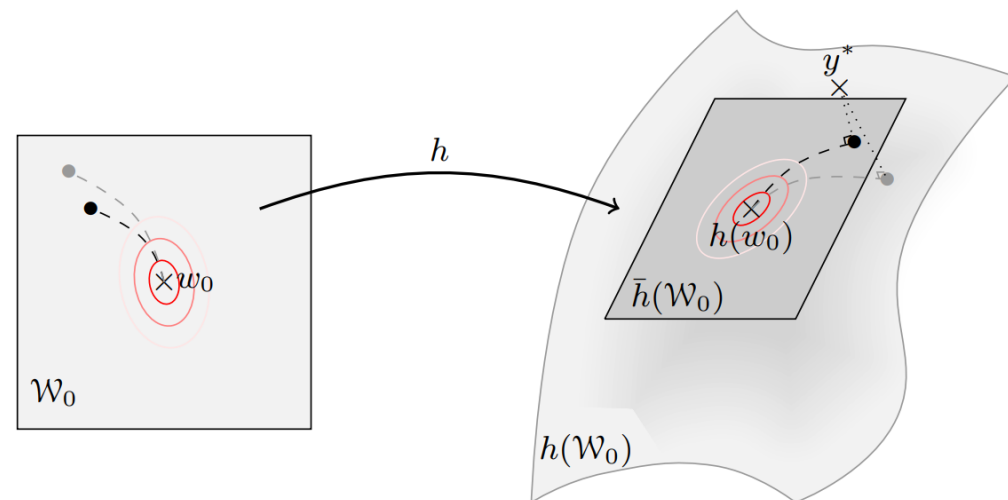


Kernels can outperform NN...



Nigam, Pozdnyakiv, Ceriotti, J. of Chem. Phys, 2020

... and be an interesting limit of NNs



Chizat, Oyallon, Bach, NeurIPS 2018
Jacot, Gabriel, Hongler, NeurIPS 2018

Quick appetizer

For a dataset characterized by

α : effective dimension of the dataset

r : complexity of the label distribution

*At which rate does the excess error decay with the number of samples n for **kernel ridge**?*

Quick appetizer

For a dataset characterized by

α : effective dimension of the dataset

r : complexity of the label distribution

*At which rate does the excess error decay with the number of samples n for **kernel ridge**?*

$$n^{-2\alpha\min(1,r)}$$

$$n^{-\frac{2\alpha\min(1,r)}{1+2\alpha\min(1,r)}}$$

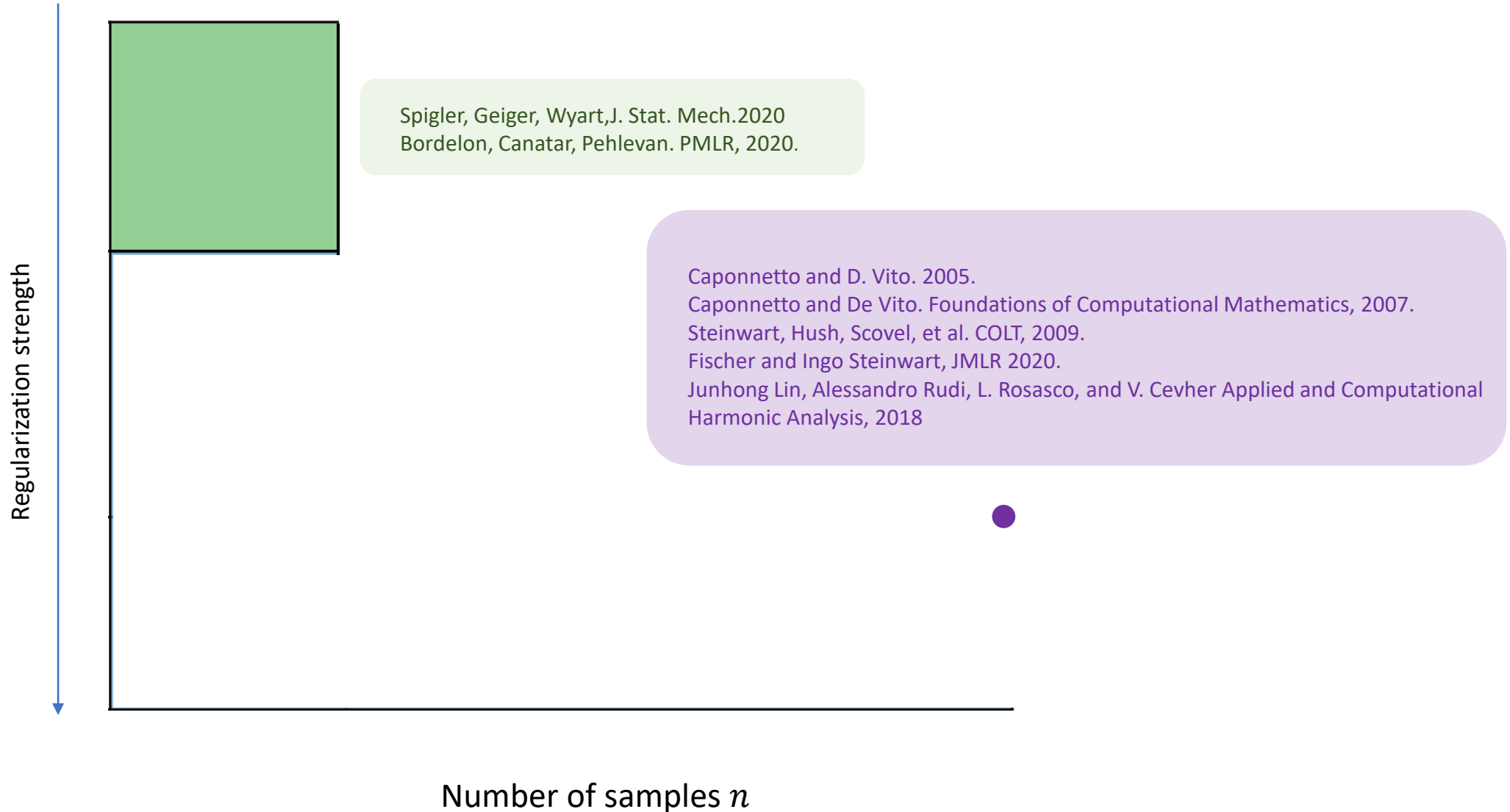
Spigler, Geiger, Wyart, J. Stat. Mech. 2020
Bordelon, Canatar, Pehlevan. PMLR, 2020.

Caponnetto and D. Vito. 2005.
Caponnetto and De Vito. Foundations of Computational Mathematics, 2007.
Steinwart, Hush, Scovel, et al. COLT, 2009.
Fischer and Ingo Steinwart, JMLR 2020.
Junhong Lin, Alessandro Rudi, L. Rosasco, and V. Cevher Applied and Computational Harmonic Analysis, 2018

Why the discrepancy?

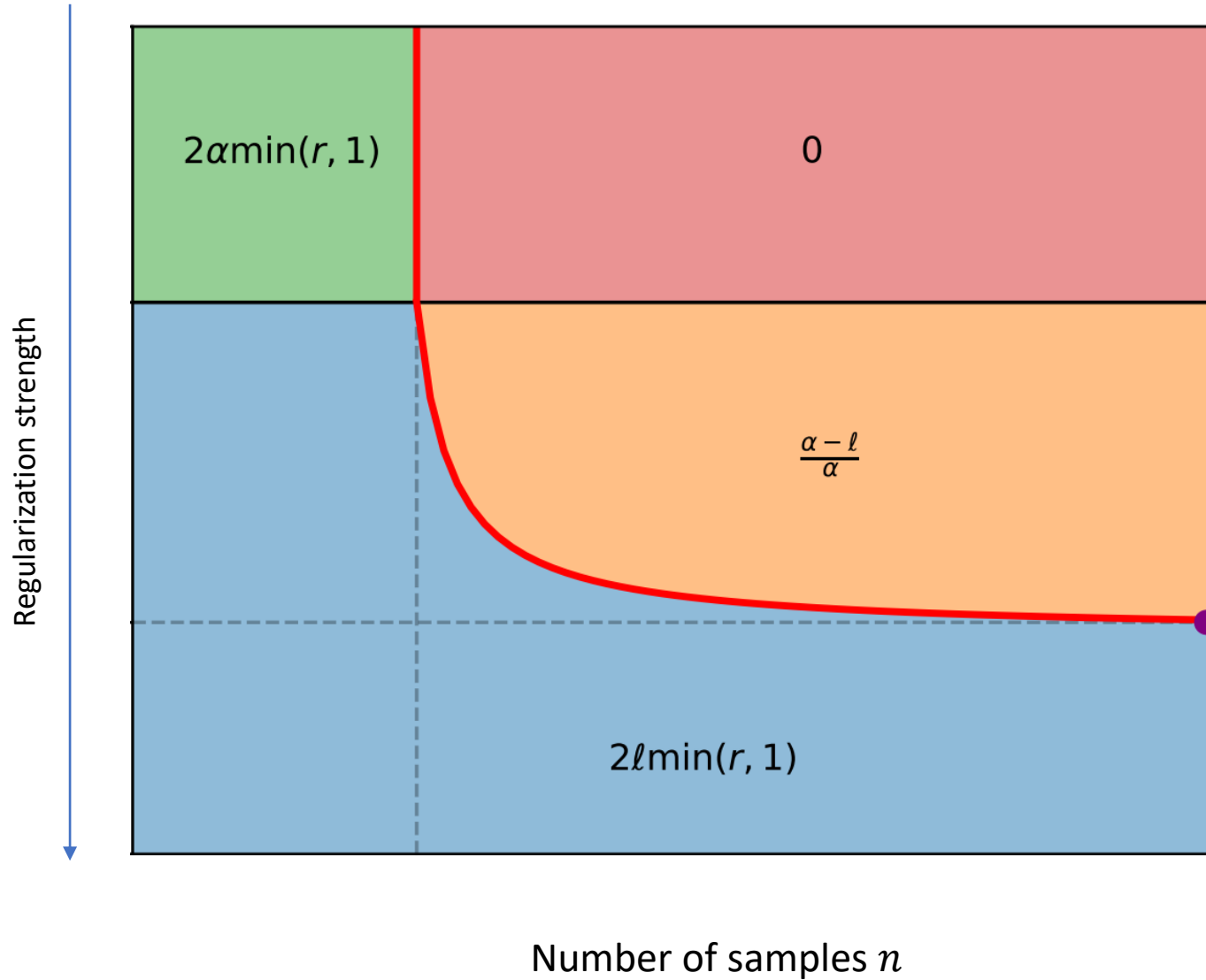
Quick appetizer

Actually these are two different regimes. We locate them on the (regularization, sample complexity plane)....



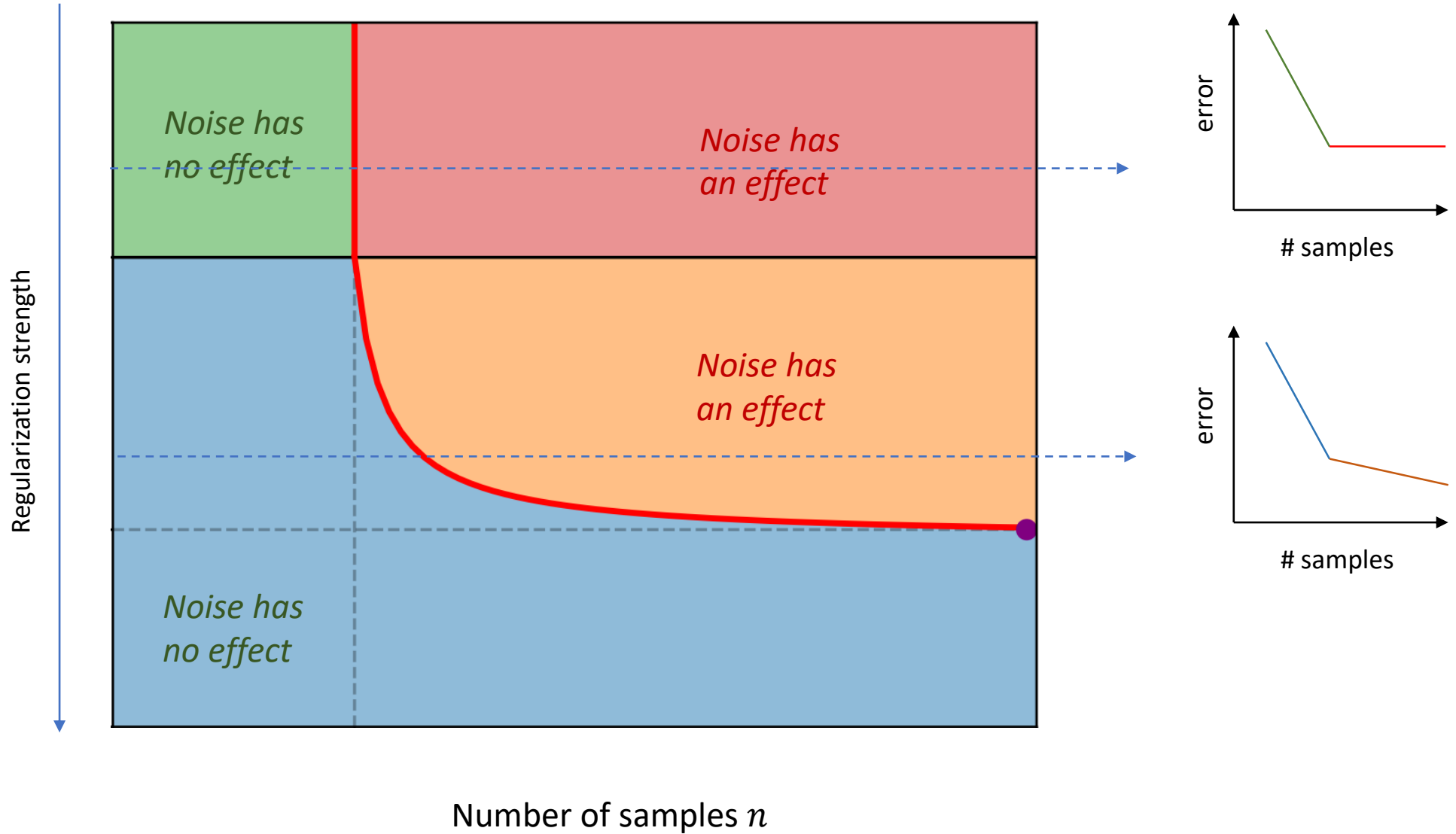
Quick appetizer

... and provide an unifying picture of the four regimes of KRR...



Quick appetizer

... and discuss when label noise affects the learning speed



A refresher on Kernels

Take a kernel K with Reproducing Kernel Hilbert Space \mathcal{H}

Kernel Ridge Regression (KRR)

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{\mu=1}^n (f(x^\mu) - y^\mu)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

samples

label

data

regularization

A refresher on Kernels

Take a kernel K with Reproducing Kernel Hilbert Space \mathcal{H}

Using a feature map $\psi(x^\mu) \in \mathbb{R}^p$

Kernel Ridge Regression (KRR)

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{\mu=1}^n (f(x^\mu) - y^\mu)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

samples
label
data
regularization

$$\hat{\mathcal{R}}_n(w) = \frac{1}{n} \sum_{\mu=1}^n (w^\top \psi(x^\mu) - y^\mu)^2 + \lambda w^\top w.$$

A refresher on Kernels

Take a kernel K with Reproducing Kernel Hilbert Space \mathcal{H}

Using a feature map $\psi(x^\mu) \in \mathbb{R}^p$

Chosen so that the covariance is diagonal

Kernel Ridge Regression (KRR)

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{\mu=1}^n (f(x^\mu) - y^\mu)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

samples
data
label
regularization

$$\hat{\mathcal{R}}_n(w) = \frac{1}{n} \sum_{\mu=1}^n (w^\top \psi(x^\mu) - y^\mu)^2 + \lambda w^\top w.$$

$$\Sigma \equiv \mathbb{E}_{x \sim \rho_x} [\psi(x)\psi(x)^\top] = \text{diag}(\eta_1, \eta_2, \dots, \eta_p)$$

Working assumptions

Gaussian design

$$\psi(x) \stackrel{d}{=} \mathcal{N}(0, \Sigma)$$

Gaussian features

$$y^\mu = \theta^* \psi(x^\mu) + \sigma \mathcal{N}(0, 1)$$

teacher + additive gaussian noise

Regularization

$$\lambda = n^{-\ell}$$

Dicker et al. Bernoulli, 2016.

Hsu, Kakade, and Zhang. PMLR 2012.

Dobriban and Wager. The Annals of Statistics, 2018.

Ledoit and P ech e. Probability Theory and Related Fields, 2011

Working assumptions

$$\exists \alpha > 1, r \geq 0$$

$$\text{tr } \Sigma^{\frac{1}{\alpha}} < \infty$$

$$\|\Sigma^{\frac{1}{2}-r} \theta^*\| < \infty$$

Working assumptions

$$\exists \alpha > 1, r \geq 0$$

$$\text{tr } \Sigma^{\frac{1}{\alpha}} < \infty$$

Capacity condition

Eigenvalues of Σ

$$\eta_k = k^{-\alpha}$$

Spigler, Geiger, Wyart, J. Stat. Mech. 2020

Bordelon, Canatar, Pehlevan. PMLR, 2020.

Dobriban and Wager. Ann. of Stat., 2018.

$$\|\Sigma^{\frac{1}{2}-r} \theta^*\| < \infty$$

Source condition

$$\theta_k^* = k^{-\frac{1+\alpha(2r-1)}{2}}$$

Working assumptions

$$\exists \alpha > 1, r \geq 0$$

$$\text{tr } \Sigma^{\frac{1}{\alpha}} < \infty$$

$$\|\Sigma^{\frac{1}{2}-r} \theta^*\| < \infty$$

Capacity condition

Source condition

Eigenvalues of Σ

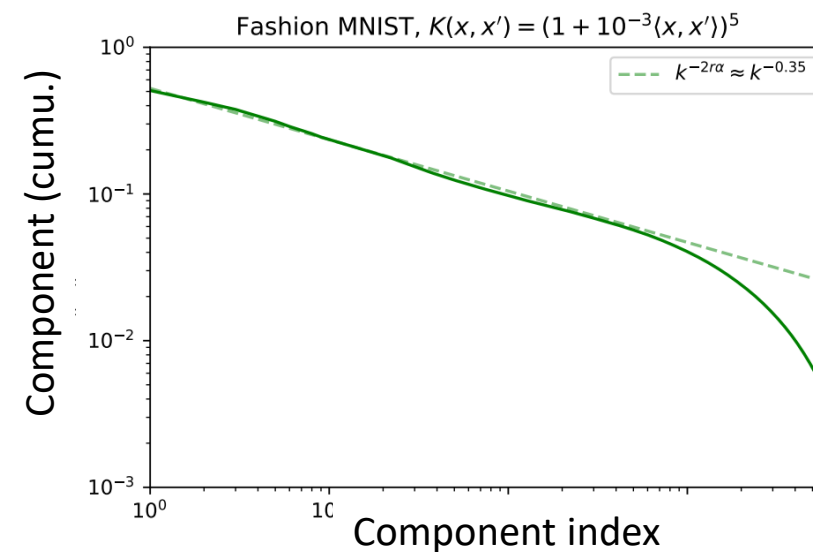
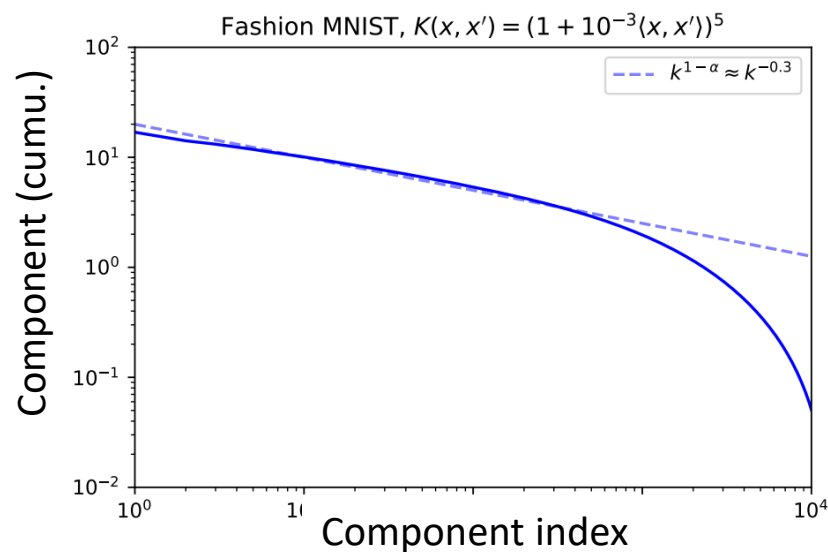
$$\eta_k = k^{-\alpha}$$

$$\theta_k^* = k^{-\frac{1+\alpha(2r-1)}{2}}$$

Spigler, Geiger, Wyart, J. Stat. Mech. 2020

Bordelon, Canatar, Pehlevan. PMLR, 2020.

Dobriban and Wager. Ann. of Stat., 2018.



Decay rates for the prediction error

$$\epsilon_g - \sigma^2 = \mathbb{E}_{u \sim \mathcal{N}(0, \Sigma)} (u^T \underbrace{\theta^*}_{\text{Teacher}} - u^T \underbrace{\hat{w}}_{\text{KRR Estimator}})^2$$

At which rate does the excess error decay with the number of samples n ?

Decay rates for the prediction error

$$\epsilon_g - \sigma^2 = \mathbb{E}_{u \sim \mathcal{N}(0, \Sigma)} (u^T \theta^* - u^T \hat{w})^2$$

Teacher *KRR Estimator*

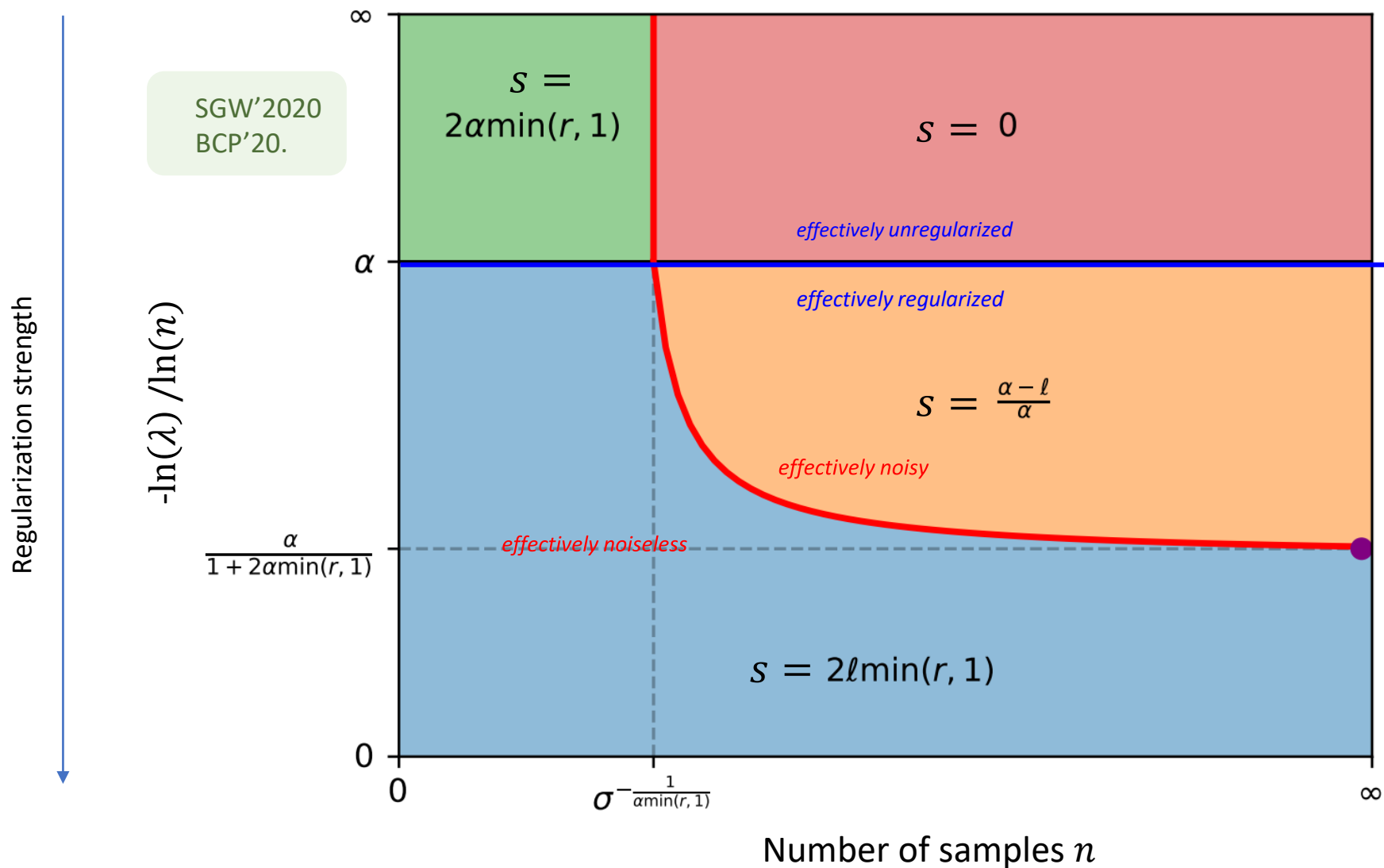
Spigler, Geiger, Wyart, J. Stat. Mech. 2020
Bordelon, Canatar, Pehlevan. PMLR, 2020.

Typical case, fast decay

Caponnetto and D. Vito. 2005.
Caponnetto and De Vito. Foundations of Computational Mathematics, 2007.
Steinwart, Hush, Scovel, et al. COLT, 2009.
Fischer and Ingo Steinwart, JMLR 2020.
Junhong Lin, Alessandro Rudi, L. Rosasco, and V. Cevher Applied and Computational Harmonic Analysis, 2018

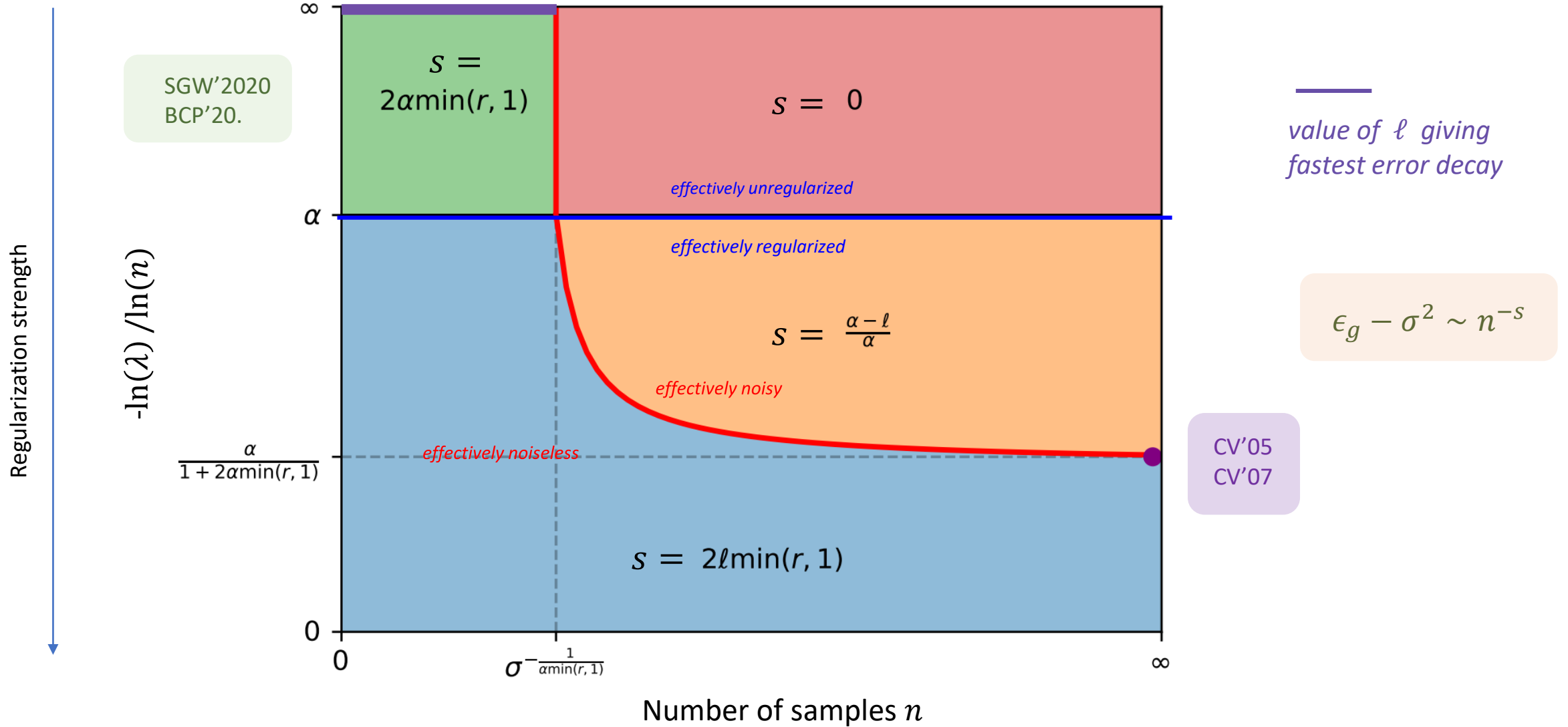
Worst case, mild decay

The four regimes



$$\epsilon_g - \sigma^2 \sim n^{-s}$$

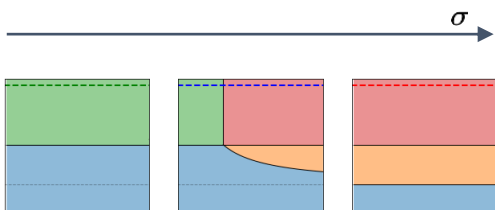
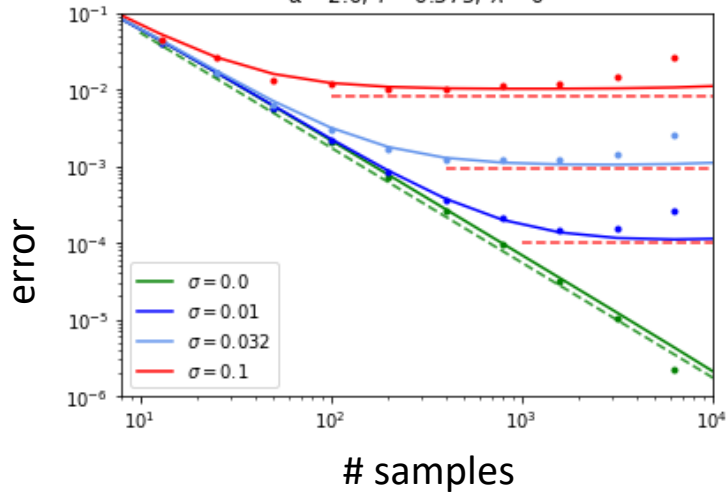
The four regimes



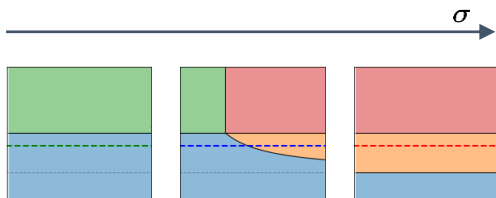
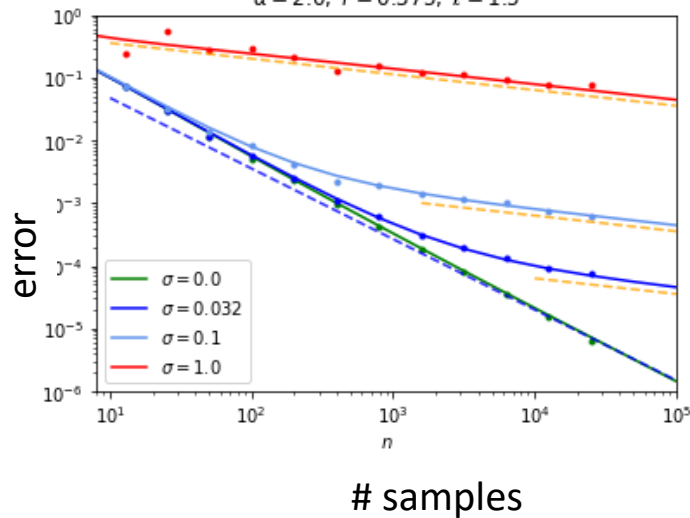
Crossovers



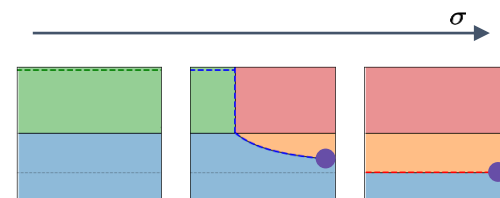
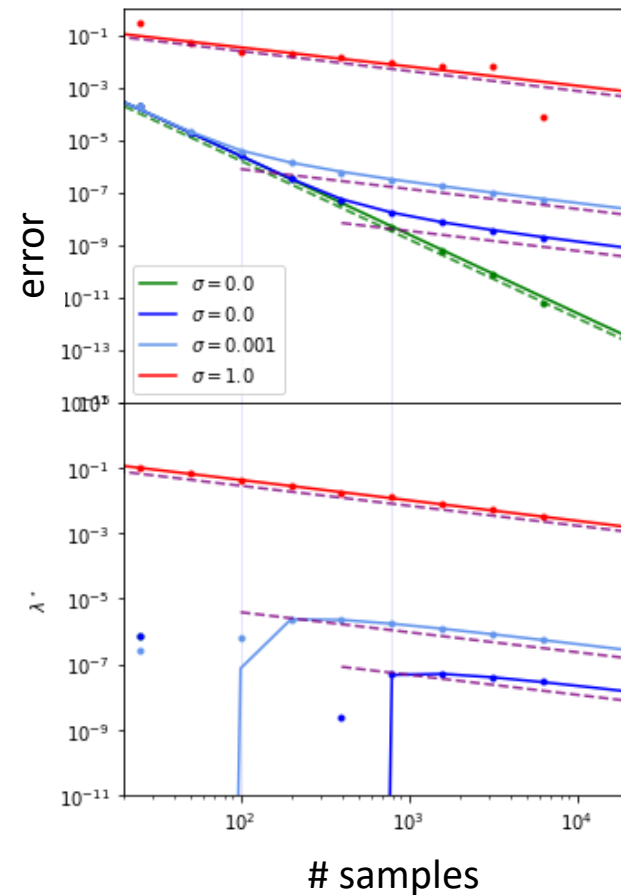
$\alpha=2.0, r=0.375, \lambda=0$



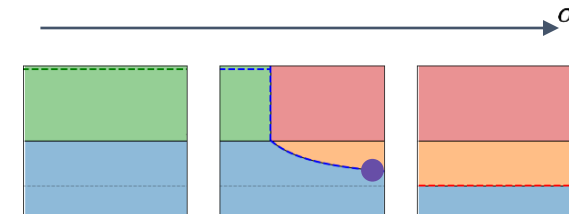
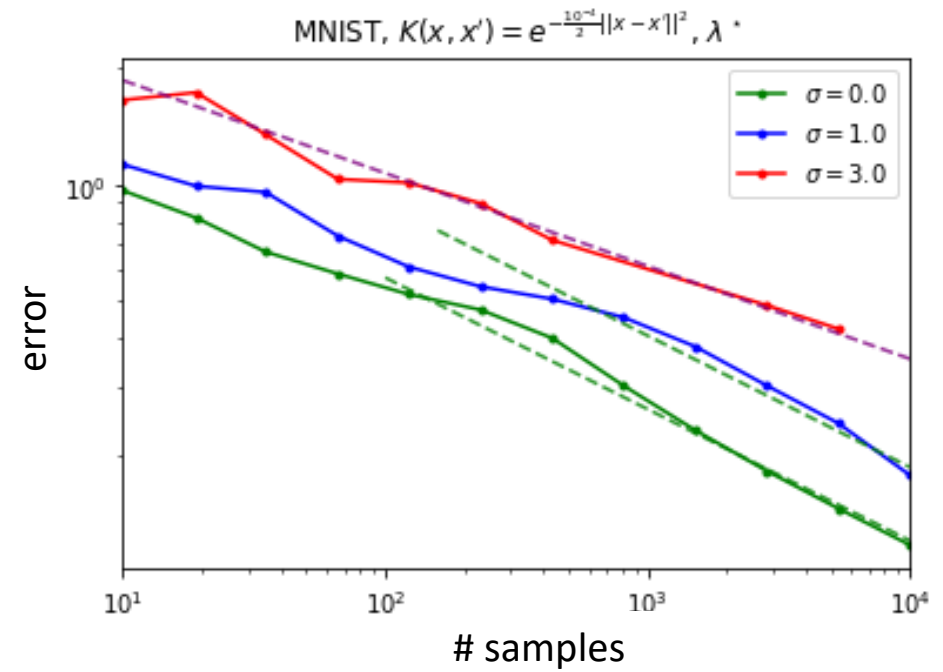
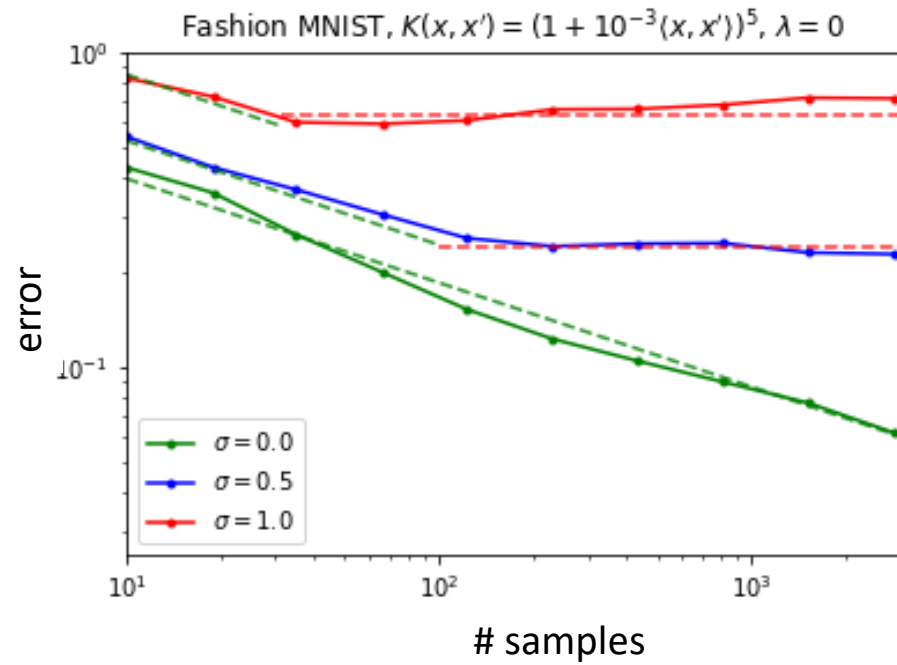
$\alpha=2.0, r=0.375, l=1.5$



$\alpha=2.5, r=0.6, l^*$



Crossovers



Thank you for your attention!

For questions, see you at the virtual poster session