



# Early stopping for kernel boosting algorithms

Yuting Wei\*, Fanny Yang\*, Martin Wainwright

Department of Statistics and EECS, University of California, Berkeley



## PROBLEM SETTING

- Given **arbitrary regular loss function**  $\phi : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$ ,  $n$  fixed covariates  $x_i$  and corresponding random  $Y_i \sim \mathbb{P}_{x_i}$
- Object of interest: minimizer of population loss function over some **function class  $\mathcal{F}$**

$$\mathcal{L}(f) := \mathbb{E}_{Y_1^n} \left[ \frac{1}{n} \sum_{i=1}^n \phi(Y_i, f(x_i)) \right] \text{ and } f^* := \arg \min_{f \in \mathcal{F}} \mathcal{L}(f)$$

- In practice: minimizer of empirical loss function based on observed  $\{x_i, Y_i\}_{i=1}^n$

$$\mathcal{L}_n(f) := \frac{1}{n} \sum_{i=1}^n \phi(Y_i, f(x_i)) \text{ and } \hat{f} := \arg \min_{f \in \mathcal{F}} \mathcal{L}_n(f)$$

- If function class large  $\mathcal{F} \rightarrow$  risk of **overfitting** to noise!
- Standard way to prevent overfitting: additive penalty function

## BOOSTING ALGORITHMS

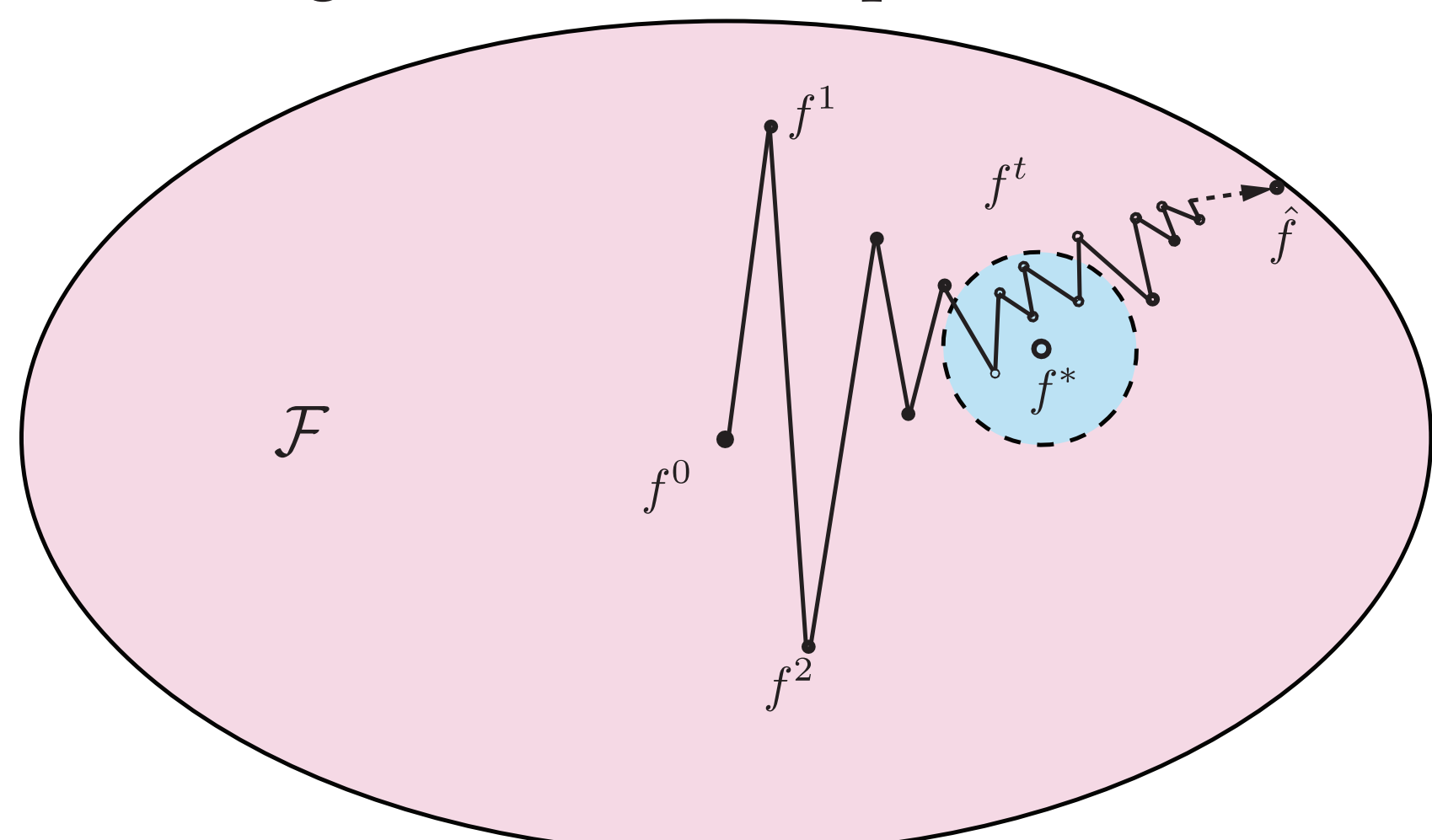
- Based on a sequence of additive updates (weak learners) to improve the fit of a function, see e.g. [1]
- Can be viewed as gradient descent on some  $\mathcal{F}$  with updates

$$f^{t+1} = f^t - \alpha^t g^t \quad \text{with} \quad g^t \propto \arg \max_{\|d\|_{\mathcal{F}} \leq 1} \langle \nabla \mathcal{L}_n(f^t), d(x_1^n) \rangle \quad (1)$$

- We consider a **Reproducing Kernel Hilbert Space** (RKHS)  $\mathcal{F} = \mathcal{H}$ , with functions  $f(\cdot) = \sum_{i=1}^n \omega_i \mathbb{K}(\cdot, x_i)$  with  $\mathbb{K}$  kernel functions, e.g. Gaussian  $\mathbb{K}(x, z) = e^{-\frac{(x-z)^2}{2\sigma^2}}$ , Sobolev  $\mathbb{K}(x, z) = 1 + \min\{x, z\}$

## EARLY STOPPING PREVENTS OVERFITTING

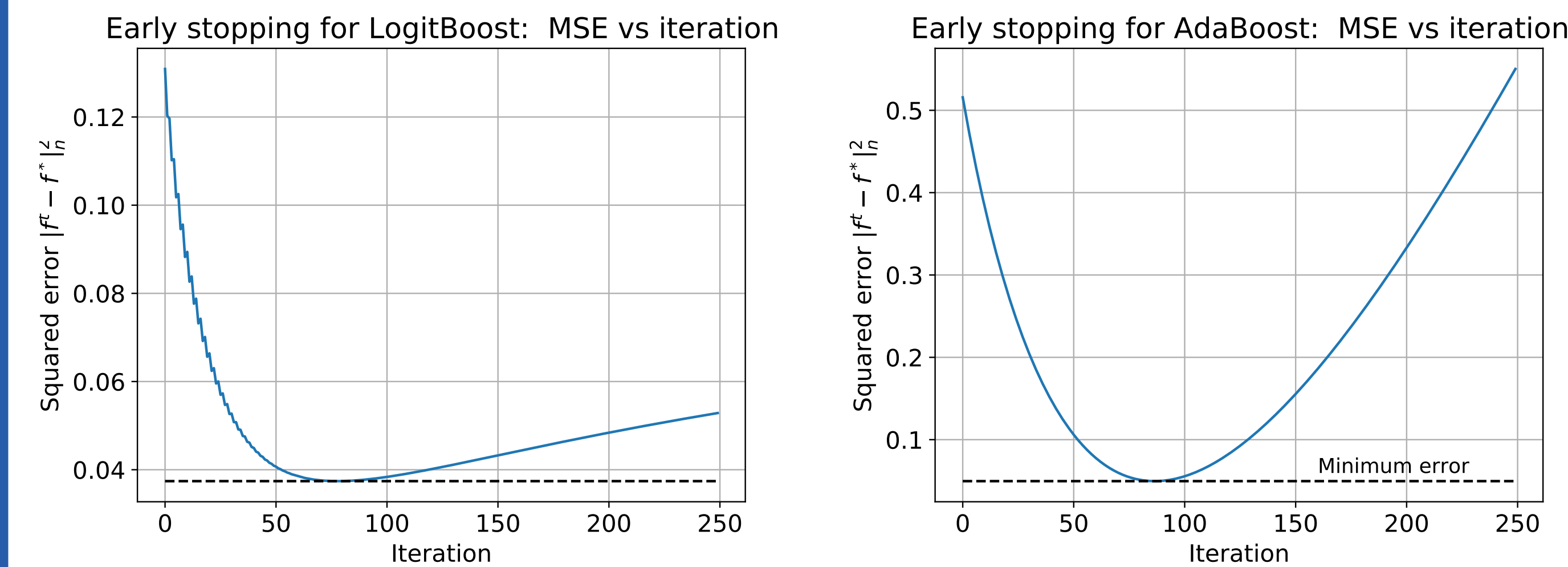
Running until convergence to  $\hat{f}$  again risks overfitting!  
Illustration of boosting in functional space:



$\rightarrow$  better to stop at some **earlier time**  $T(\mathcal{F}, n)$  when  $f^t$  closest to  $f^*$ .  
Early stopped boosting estimators are **consistent** (see e.g. [2,3])!

## EARLY STOPPED VS. PENALTY REGULARIZED

Distance to population optimum  $f^*$  over iterates:



For **least-squares loss**, early stopped boosting (*algorithmic regular.*) and penalized estimators behave similarly statistically [4], i.e.

$$\|f_{\text{pen}} - f^*\|_n^2 \sim \|f^T - f^*\|_n^2$$

## MAIN CONTRIBUTIONS

- Is there a common principle behind statistical behavior of algorithmic and penalized regularization?

**YES!** We prove statistical rates for early stopping using **the same key quantities as in penalized regularization**

- Can we extend to other loss functions?

**YES!** Our new proof technique allows to extend to a **broad class of loss functions (e.g. AdaBoost, LogitBoost ...)**

## KEY QUANTITIES AND MAIN THEOREM

- Key quantity I: **Localized Gaussian complexity**

$$\mathcal{G}_n(\mathcal{E}(\delta, 1)) := \mathbb{E} \left[ \sup_{g \in \mathcal{E}(\delta, 1)} \frac{1}{n} \sum_{i=1}^n w_i g(x_i) \right], \quad w_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$$

where:  $\mathcal{E}(\delta, 1) := \{f - g \mid f, g \in \mathcal{H}, \|f - g\|_{\mathcal{H}} \leq 1, \|f - g\|_n \leq \delta\}$

- Key quantity II: **Critical radius**  $\delta_n$  smallest  $\delta$  s.t.  $\frac{\mathcal{G}_n(\mathcal{E}(\delta, 1))}{\delta} \leq \frac{\delta}{\sigma}$

**Theorem 1.** Given some regular loss function  $\phi$  and the iterates  $\{f^t\}_{t=0}^{\infty}$  in (1), for all iterations  $T = 0, 1, \dots, \lfloor 1/(8\delta_n^2) \rfloor$ , the averaged estimate  $\bar{f}^T$  satisfies w.h.p.

$$\|\bar{f}^T - f^*\|_n^2 \leq O\left(\frac{1}{T} + \delta_n^2\right).$$

## REFERENCES

[1] R. E. Schapire, "The boosting approach to machine learning: An overview," in *Nonlinear estimation and classification*. Springer, 2003, pp. 149–171.  
 [2] T. Zhang and B. Yu, "Boosting with early stopping: Convergence and consistency," *Annals of Statistics*, vol. 33, no. 4, pp. 1538–1579, 2005.  
 [3] P. L. Bartlett and M. Traskin, "Adaboost is consistent," *Journal of Machine Learning Research*, vol. 8, no. Oct, pp. 2347–2368, 2007.  
 [4] G. Raskutti, M. Wainwright, and B. Yu, "Early stopping and non-parametric regression: An optimal data-dependent stopping rule," *Journal of Machine Learning Research*.

## OPTIMALITY

Early stopped boosted estimators are usually statistically optimal:

**Corollary 1.** For the class of regular kernels and any function  $f^*$  with  $\|f^*\|_{\mathcal{H}} \leq 1$ , running  $T := \lfloor \frac{1}{8\delta_n^2} \rfloor$  iterations

$$\mathbb{E} \|\bar{f}^T - f^*\|_n^2 \asymp \inf_{\hat{f} \|f^*\|_{\mathcal{H}} \leq 1} \sup \mathbb{E}_{Y_1^n} \|\hat{f} - f^*\|_n^2.$$

Examples for specific kernel spaces:

- $\gamma$ -exponential decay:** the kernel eigenvalues  $\mu_j \leq c_1 \exp(-c_2 j^\gamma)$ , when stopped after  $T \asymp \frac{n}{\log^{1/\gamma} n}$  steps:

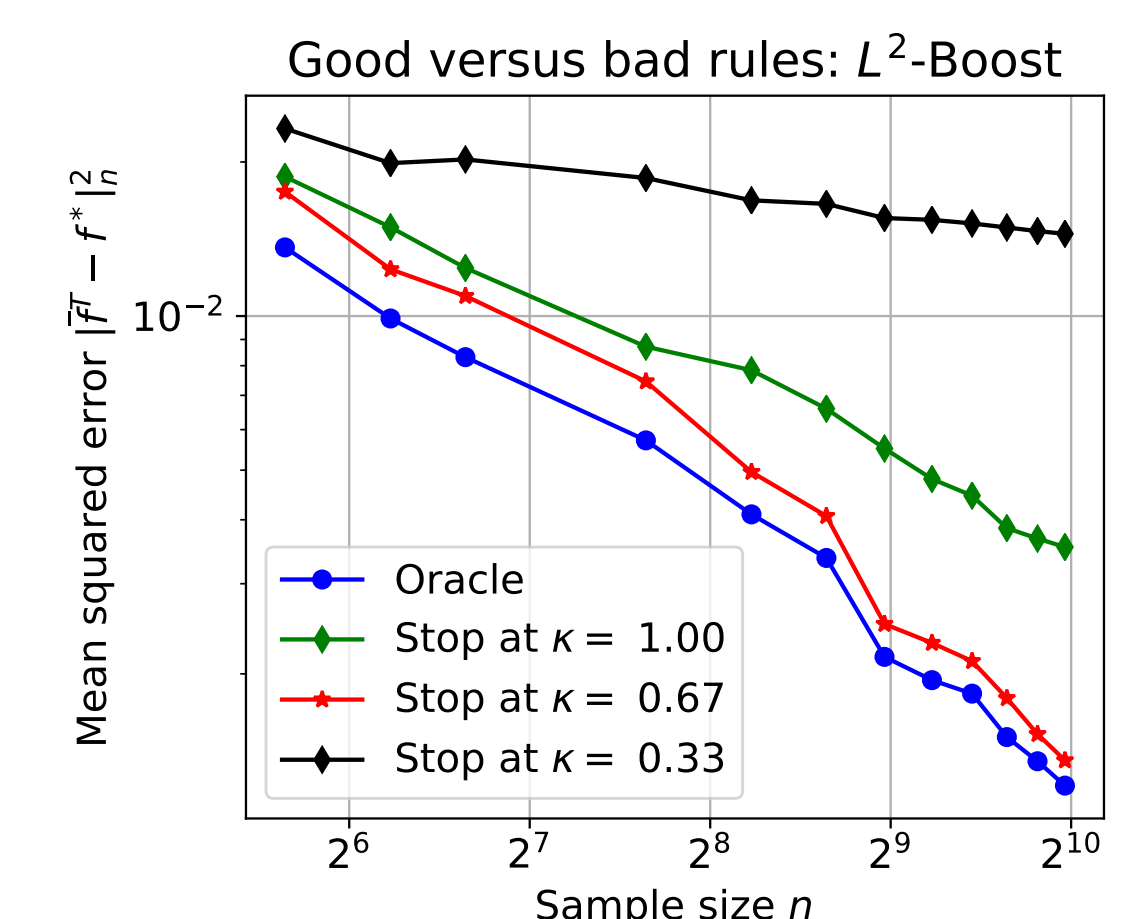
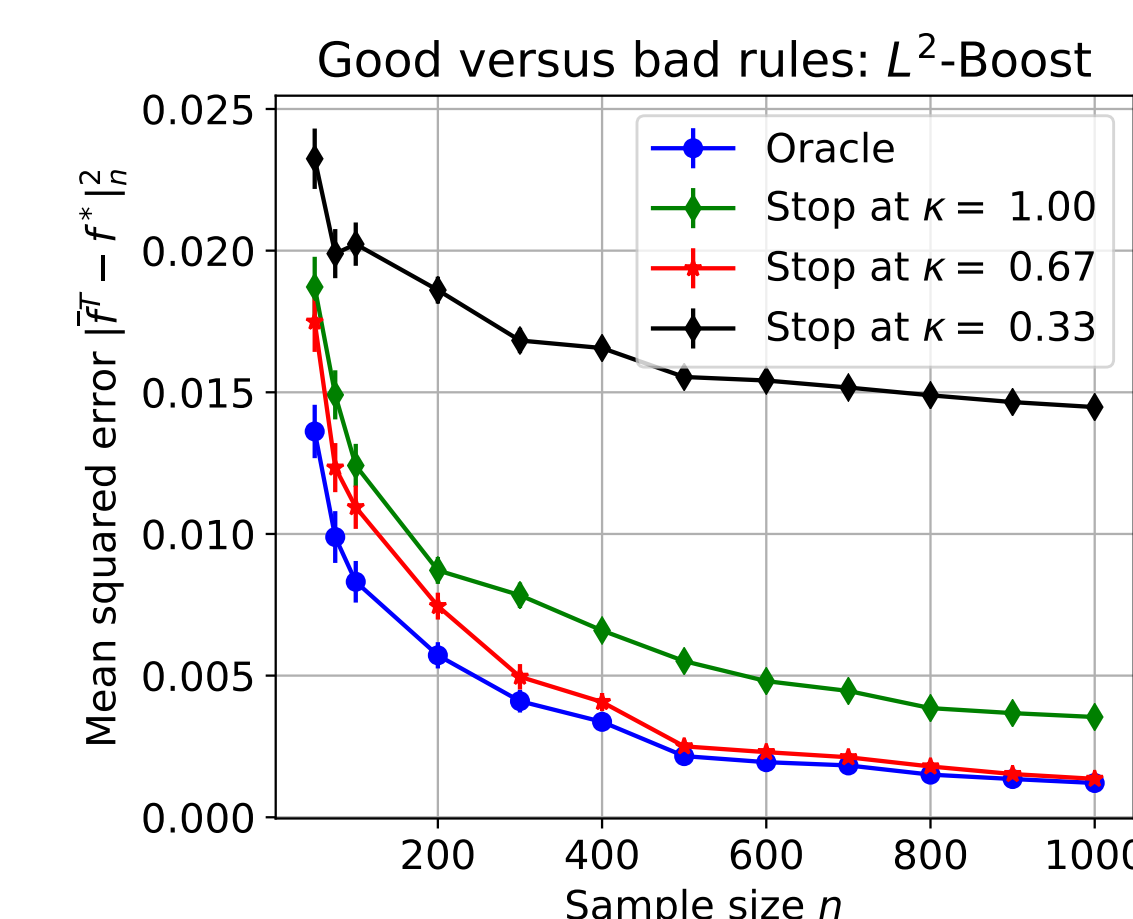
$$\|\bar{f}^T - f^*\|_n^2 \lesssim \frac{\log^{1/\gamma} n}{n}$$

- $\beta$ -polynomial decay:** the kernel eigenvalues  $\mu_j \leq c_1 j^{-2\beta}$ , when stopped after  $T \asymp n^{2\beta/(2\beta+1)}$  steps:

$$\|\bar{f}^T - f^*\|_n^2 \lesssim n^{-\frac{2\beta}{2\beta+1}}$$

## NUMERICAL RESULTS

$\mathcal{H}$ : first order Sobolev space, stop iterates after  $T = (7n)^\kappa$



Log-scale: Error scales worse when stopped too early or too late

