

# Sparse Learning with CART

Jason M. Klusowski

Princeton University, Operations Research and Financial Engineering

## Objectives

- When do decision trees adapt to the sparsity of a predictive model?

## Introduction

- Training data  
 $\mathcal{D} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}, \quad (\mathbf{X}_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$
- Predictor for decision tree  $T$   
 $\hat{Y} = \hat{Y}(T, \mathcal{D})$
- Prediction error  
 $\text{Err}(\hat{Y}(T)) = \mathbb{E}_{(\mathbf{X}', Y')}[(Y' - \hat{Y}(\mathbf{X}'))^2]$   
 for independent copy  $(\mathbf{X}', Y')$

## CART decision trees

- CART [1] methodology based on recursively minimizing impurity
- For regression, impurity in node  $t \in T$  is

$$\hat{\Delta}(t) = \frac{1}{N(t)} \sum_{\mathbf{X}_i \in t} (Y_i - \bar{Y}_t)^2,$$

where  $N(t) = \#\{\mathbf{X}_i \in t\}$  and  $\bar{Y}_t = \frac{1}{N(t)} \sum_{\mathbf{X}_i \in t} Y_i$

- Optimal direction  $\hat{j}$  and split point  $\hat{s}$  obtained by maximizing reduction in impurity

$$\hat{\Delta}(s, t) = \hat{\Delta}(t) - \frac{N(t_L)}{N(t)} \hat{\Delta}(t_L) - \frac{N(t_R)}{N(t)} \hat{\Delta}(t_R),$$

where

$$t_L = \{\mathbf{X} \in t : X_j \leq s\}, \quad t_R = \{\mathbf{X} \in t : X_j > s\}$$

are left and right child nodes

- Tree output  $\hat{Y}(\mathbf{x}) = \bar{Y}_t$  for  $\mathbf{x}$  in terminal node  $t$

## Main results

- Consider pruned tree

$$\hat{T} \in \arg \min_{T \preceq T_{\max}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}(\mathbf{X}_i))^2 + \alpha |T| \right\},$$

where  $T_{\max}$  is fully grown tree, temperature  $\alpha = \Theta((d/n) \log(n/d))$ , and  $|T|$  is # of terminal nodes

## Theorem

Suppose  $\mathbf{X}$  is uniformly distributed on  $[0, 1]^d$  and

$$Y = \sum_j g_j(X_j)$$

is a sparse additive model with  $d_0 \ll d$  smooth component functions  $g_j(\cdot)$ , where each function is not too locally 'flat'. Then,

$$\limsup_n \frac{\text{Err}(\hat{Y}(\hat{T}))}{((d/n) \log(n/d))^{\Omega(1/d_0)}} \stackrel{\text{a.s.}}{=} \mathcal{O}(1).$$

## Proof idea

- Reduction in impurity  $\hat{\Delta}(\hat{s}, t)$  can be written as

$$\hat{\Delta}(t) \times \hat{\rho}^2(\hat{Y}, Y | \mathbf{X} \in t),$$

where  $\hat{\rho} = \hat{\rho}(\hat{Y}, Y | \mathbf{X} \in t)$  is Pearson correlation between response data  $Y$  and optimal decision stump

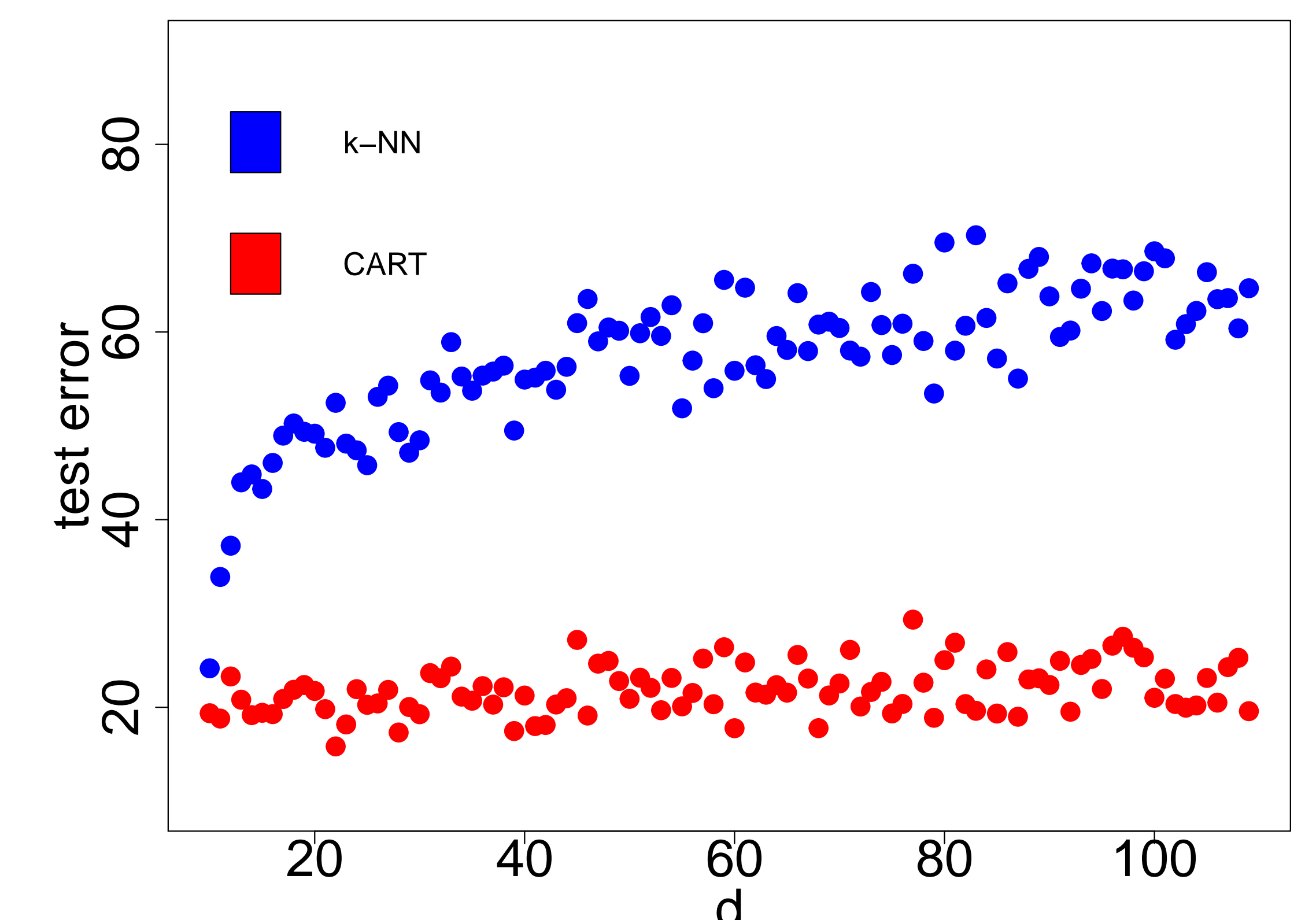
$$\hat{Y} = \bar{Y}_L \mathbf{1}(X_j \leq \hat{s}) + \bar{Y}_R \mathbf{1}(X_j > \hat{s})$$

- Training error bound

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}(\mathbf{X}_i))^2 \leq \widehat{\text{Var}}(Y) \exp(-K \times \min_t \hat{\rho}^2),$$

where  $K = \Theta(\log_2(n))$  is tree depth and  $\liminf_n \min_t \hat{\rho}^2 = \Omega(1/d_0)$  a.s.

## Experiments



**Figure:** Boston housing dataset [1] ( $d_0 = 10$  and  $n = 506$ ) with  $d - d_0$  noisy features added. Plot shows prediction error of pruned CART vs. cross-validated  $k$ -NN as  $d$  varies.

## Conclusion

- CART adapts to underlying sparsity, whereas kernel methods with nonadaptive weights (like  $k$ -NN) suffer from curse of dimensionality

## References

- [1] Leo Breiman, Jerome Friedman, RA Olshen, and Charles J Stone.  
*Classification and regression trees.*  
 Chapman and Hall/CRC, 1984.