

Intelligible Models for Classification and Regression

Yin Lou¹ Rich Caruana² Johannes Gehrke¹

Department of Computer Science¹
Cornell University

Microsoft Research²
Microsoft Corporation

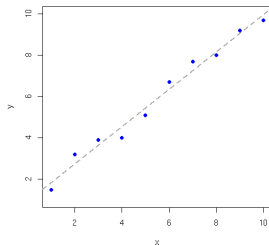
Aug. 13, 2012

Motivation

Simple Model

- Linear regression, logistic regression
- Regression: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$
- Classification: $\text{logit}(y) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$

Linear Regression

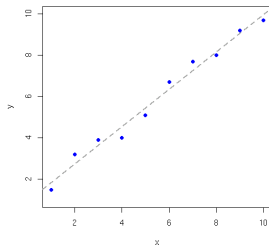


Motivation

Simple Model

- Linear regression, logistic regression
- Regression: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$
- Classification: $\text{logit}(y) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$

Linear Regression

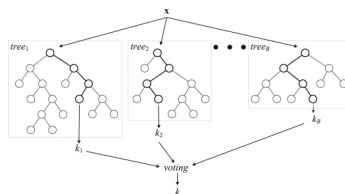


Intelligible but usually **less accurate**

Complex Model

- Random forest, SVMs with RBF kernel, etc.
- $y = f(x_1, \dots, x_n)$

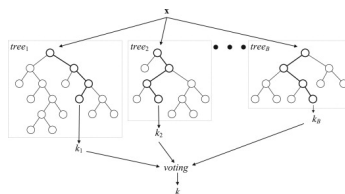
Random Forest



Complex Model

- Random forest, SVMs with RBF kernel, etc.
- $y = f(x_1, \dots, x_n)$

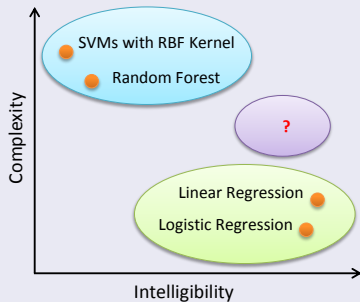
Random Forest



Unintelligible but usually more accurate

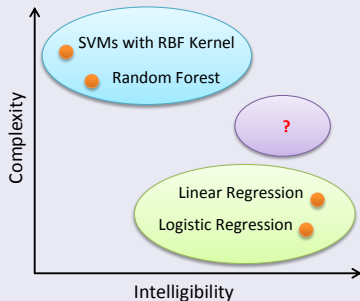
Motivation

The tradeoff



Motivation

The tradeoff



Intelligibility is important

- Medical applications
- Domains where we want scientific understanding
- Efficient model engineering
 - Impact of features in a ranker

Outline

- 1 Motivation
- 2 Towards More Accurate Models
- 3 Algorithms
- 4 Experiments
- 5 Discussion
- 6 Conclusion

Outline

- 1 Motivation
- 2 Towards More Accurate Models
- 3 Algorithms
- 4 Experiments
- 5 Discussion
- 6 Conclusion

Generalized Additive Models

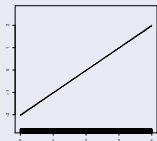
- Developed by Hastie and Tibshirani
- Regression: $y = f_1(x_1) + \dots + f_n(x_n)$
- Classification: $\text{logit}(y) = f_1(x_1) + \dots + f_n(x_n)$
- Each feature is “shaped” by shape function f_i
- **Intelligible** and **accurate**



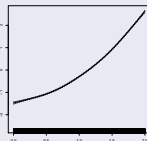
T. Hastie and R. Tibshirani.
Generalized additive models.
Chapman & Hall/CRC, 1990.

Example

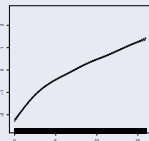
$$y = x_1 + x_2^2 + \sqrt{x_3} + \log x_4 + e^{x_5} + 2 \sin x_6 + \epsilon$$



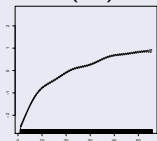
$f_1(x_1)$



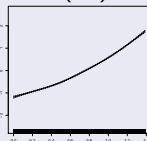
$f_2(x_2)$



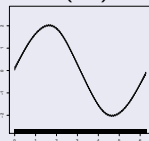
$f_3(x_3)$



$f_4(x_4)$



$f_5(x_5)$



$f_6(x_6)$

Figure: Shape Functions for Synthetic Dataset.

Model	Form	Intelligibility	Accuracy
Linear Model	$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$	+++	+
Generalized Linear Model	$g(y) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$	+++	+
Additive Model	$y = f_1(x_1) + \dots + f_n(x_n)$	++	++
Generalized Additive Model	$g(y) = f_1(x_1) + \dots + f_n(x_n)$	++	++
Full Complexity Model	$y = f(x_1, \dots, x_n)$	+	+++

Table: From Linear to Additive Models.

Outline

- 1 Motivation
- 2 Towards More Accurate Models
- 3 Algorithms**
- 4 Experiments
- 5 Discussion
- 6 Conclusion

Fitting GAMs

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Shape Functions

- Splines (SP)
- Single Tree (TR)
- Bagged Trees (bagTR)
- Boosted Trees (bstTR)
- Boosted Bagged Trees (bbTR)

Fitting GAMs

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Shape Functions

- Splines (SP)
- Single Tree (TR)
- Bagged Trees (bagTR)
- Boosted Trees (bstTR)
- Boosted Bagged Trees (bbTR)

Learning Methods

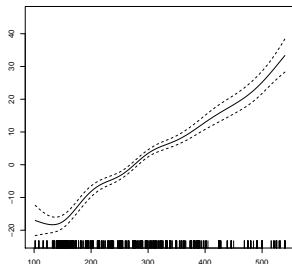
- Penalized Least Squares (P-LS/P-IRLS)
- Backfitting (BF)
- Gradient Boosting (BST)

Fitting GAMs

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Shape Function: Splines (SP)

- $f_i(x_i) = \sum_{k=1}^d \beta_k b_k(x_i)$

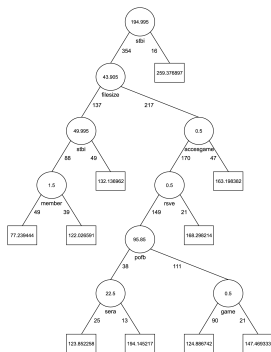


Fitting GAMs

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Shape Function: Single Tree (TR)

- $f_i(x_i) = \text{RegressionTree}(x_i, \text{response})$

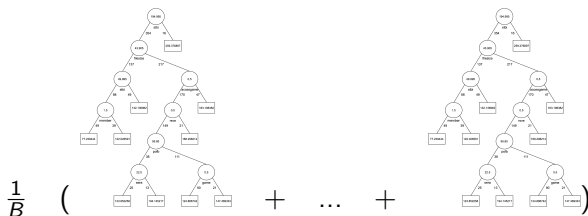


Fitting GAMs

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Shape Function: Bagged Trees (bagTR)

- $f_i(x_i) = \frac{1}{B} \sum_{j=1}^B \text{Regression Tree}(x_i, \text{bootstrap sample } j)$

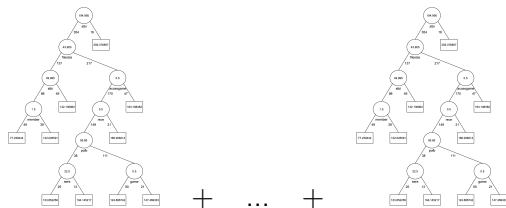


Fitting GAMs

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Shape Function: Boosted Trees (bstTR)

- $f_i(x_i) = \sum_{j=1}^B \text{RegressionTree}(x_i, \text{residual}_j)$



$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Shape Function: Boosted Bagged Trees (bbTR)

- $f_i(x_i) = \sum_{j=1}^B \text{BaggedRegressionTree}(x_i, \text{residual}_j)$

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Learning Method: Penalized Least Squares (P-LS/P-IRLS)

- Works only on Splines ($f_i(x_i) = \sum_{k=1}^d \beta_k b_k(x_i)$)
- Converts the optimization problem to fitting linear regression/logistic regression with different basis



S. Wood.

Generalized additive models: an introduction with R.
CRC Press, 2006.

Fitting GAMs

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Learning Method: Backfitting (BF)

- 1: $f_j \leftarrow 0$
- 2: **for** $m = 1$ to M **do**
- 3: **for** $j = 1$ to n **do**
- 4: $\mathcal{R} \leftarrow \{x_{ij}, y_i - \sum_{k \neq j} f_k\}_1^N$
- 5: Learn shaping function $S : x_j \rightarrow y$ using \mathcal{R} as training dataset
- 6: $f_j \leftarrow S$
- 7: **end for**
- 8: **end for**



T. Hastie and R. Tibshirani.
Generalized additive models.
Chapman & Hall/CRC, 1990.

Fitting GAMs

$$g(y) = f_1(x_1) + \dots + f_n(x_n)$$

Learning Method: Gradient Boosting (BST)

```
1:  $f_j \leftarrow 0$ 
2: for  $m = 1$  to  $M$  do
3:   for  $j = 1$  to  $n$  do
4:      $\mathcal{R} \leftarrow \{x_{ij}, y_i - \sum_k f_k\}_1^N$ 
5:     Learn shaping function  $S : x_j \rightarrow y$  using  $\mathcal{R}$  as training dataset
6:      $f_j \leftarrow f_j + S$ 
7:   end for
8: end for
```



J. Friedman.

Greedy function approximation: a gradient boosting machine.

Annals of Statistics, 29:1189–1232, 2001.

- First large-scale study that uses trees as shape function for GAMs
- Novel methods for using trees as shape functions
- Largest empirical study of fitting GAMs

Outline

- 1 Motivation
- 2 Towards More Accurate Models
- 3 Algorithms
- 4 Experiments**
- 5 Discussion
- 6 Conclusion

	Dataset	Size	Attributes	%Pos
Regression	Concrete	1030	9	-
	Wine	4898	12	-
	Delta	7192	6	-
	CompAct	8192	22	-
	Music	50000	90	-
	Synthetic	10000	6	-
Classification	Spambase	4601	58	39.40
	Insurance	9823	86	5.97
	Magic	19020	11	64.84
	Letter	20000	17	49.70
	Adult	46033	9/43	16.62
	Physics	50000	79	49.72

Shape Function	Least Squares	Gradient Boosting	Backfitting
Splines	P-LS/P-IRLS	BST-SP	BF-SP
Single Tree	N/A	BST-TR _x	BF-TR
Bagged Trees	N/A	BST-bagTR _x	BF-bagTR
Boosted Trees	N/A	BST-TR _x	BF-bstTR _x
Boosted Bagged Trees	N/A	BST-bagTR _x	BF-bbTR _x

Table: Notation for learning methods and shape functions.

- 9 different methods
- 5-fold cross validation for each method

Model	Regression	Classification	Mean
Linear/Logistic			
P-LS/P-IRLS BST-SP BF-SP			
BST-bagTR2 BST-bagTR3 BST-bagTR4 BST-bagTRX			
Random Forest			

Results

Model	Regression	Classification	Mean
Linear/Logistic	1.68	1.22	1.45
P-LS/P-IRLS BST-SP BF-SP			
BST-bagTR2 BST-bagTR3 BST-bagTR4 BST-bagTRX			
Random Forest			

Results

Model	Regression	Classification	Mean
Linear/Logistic	1.68	1.22	1.45
P-LS/P-IRLS BST-SP BF-SP			
BST-bagTR2 BST-bagTR3 BST-bagTR4 BST-bagTRX			
Random Forest	0.88	0.80	0.84

Results

Model	Regression	Classification	Mean
Linear/Logistic	1.68	1.22	1.45
P-LS/P-IRLS	1.00	1.00	1.00
BST-SP	1.04	1.00	1.02
BF-SP	1.00	1.00	1.00
BST-bagTR2			
BST-bagTR3			
BST-bagTR4			
BST-bagTRX			
Random Forest	0.88	0.80	0.84

Results

Model	Regression	Classification	Mean
Linear/Logistic	1.68	1.22	1.45
P-LS/P-IRLS	1.00	1.00	1.00
BST-SP	1.04	1.00	1.02
BF-SP	1.00	1.00	1.00
BST-bagTR2	0.96	0.96	0.96
BST-bagTR3	0.97	0.95	0.96
BST-bagTR4	0.99	0.95	0.97
BST-bagTRX	0.95	0.94	0.95
Random Forest	0.88	0.80	0.84

Model	Regression	Classification	Mean
Linear/Logistic	1.68	1.22	1.45
P-LS/P-IRLS	1.00	1.00	1.00
BST-SP	1.04	1.00	1.02
BF-SP	1.00	1.00	1.00
BST-bagTR2	0.96	0.96	0.96
BST-bagTR3	0.97	0.95	0.96
BST-bagTR4	0.99	0.95	0.97
BST-bagTRX	0.95	0.94	0.95
Random Forest	0.88	0.80	0.84

Observations

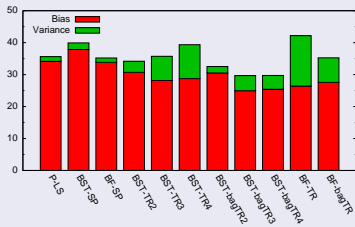
- Two accuracy gaps: shaping and interactions
- Tree-base shaping methods are more accurate than spline-based methods

Outline

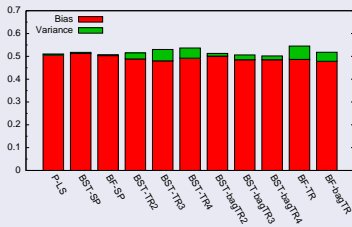
- 1 Motivation
- 2 Towards More Accurate Models
- 3 Algorithms
- 4 Experiments
- 5 Discussion**
- 6 Conclusion

Bias Variance Decomposition

$$\text{Expected Loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$



(a) Concrete



(b) Wine

Figure: Bias-variance analysis.

Learned Shaped Function: Splines vs. Trees

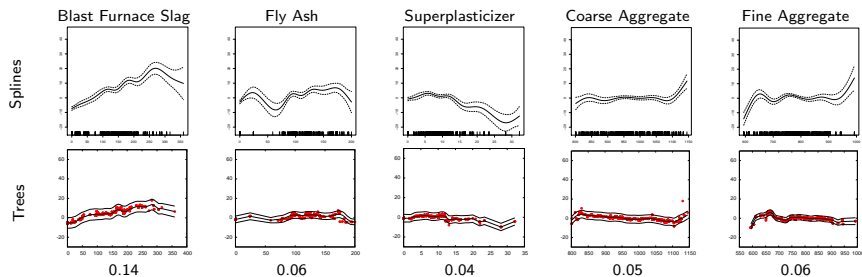


Figure: Shapes of features for the "Concrete" dataset produced by P-LS (top) and BST-bagTR3 (bottom).

Outline

- 1 Motivation
- 2 Towards More Accurate Models
- 3 Algorithms
- 4 Experiments
- 5 Discussion
- 6 Conclusion**

Conclusion

- Generalized additive models are accurate and intelligible
- Tree has low bias but high variance
- Bagging reduces variance and makes tree-based method stand out
- Bagged shallow trees with gradient boosting are most accurate

Future Work

- Feature selection
- Scalability
- Statistical interaction detection

Thank You

Questions?