



# Application of Generalized Boosted Regression in the XM Tool Project

**Ken Lau**

Co-op Student

Environment Canada

Air Quality Science Unit

***August 25, 2011***



Environment  
Canada

Environnement  
Canada

Canada

# Outline

---

- List of Predictors
- Leaps BIC Predictor Selection
- Generalized Boosted Models (GBM)
  - Introduction
  - Parameters
  - Optimization
  - Results
- Conclusion and Future Work



# List of Predictors

---

## Predictors

84 AQ Model Predictors

3 **Persistence Predictors**

27 **Antecedent Predictors**

## **Persistence Predictors**

OBS at Hr 00

## **Antecedent Predictors**

*Lag 24/48/72 hrs*

Max

Min

OBS



# Leaps BIC Predictor Selection

---

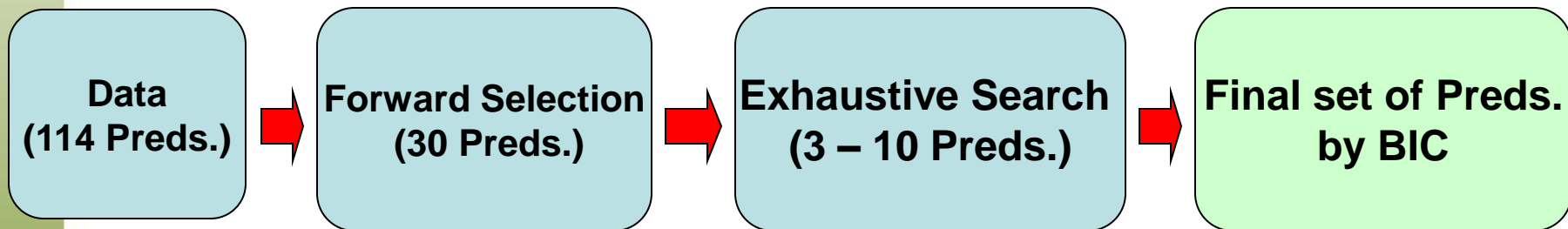
- Leaps Package in R Provides Automated Predictor Selection Routines
- Combination of Forward, Exhaustive, and minimum BIC Selection
  - Fast (Forward Selection)
  - Accurate (Exhaustive)
  - Avoids Overfitting (minimum BIC)
- $BIC = n \cdot \log(RSS/n) + (\log(n)) \cdot (p+1)$

$n = \# \text{ Obs}$  |  $p = \# \text{ Predictors}$  |  $RSS = \text{Residual Sum of Squares}$



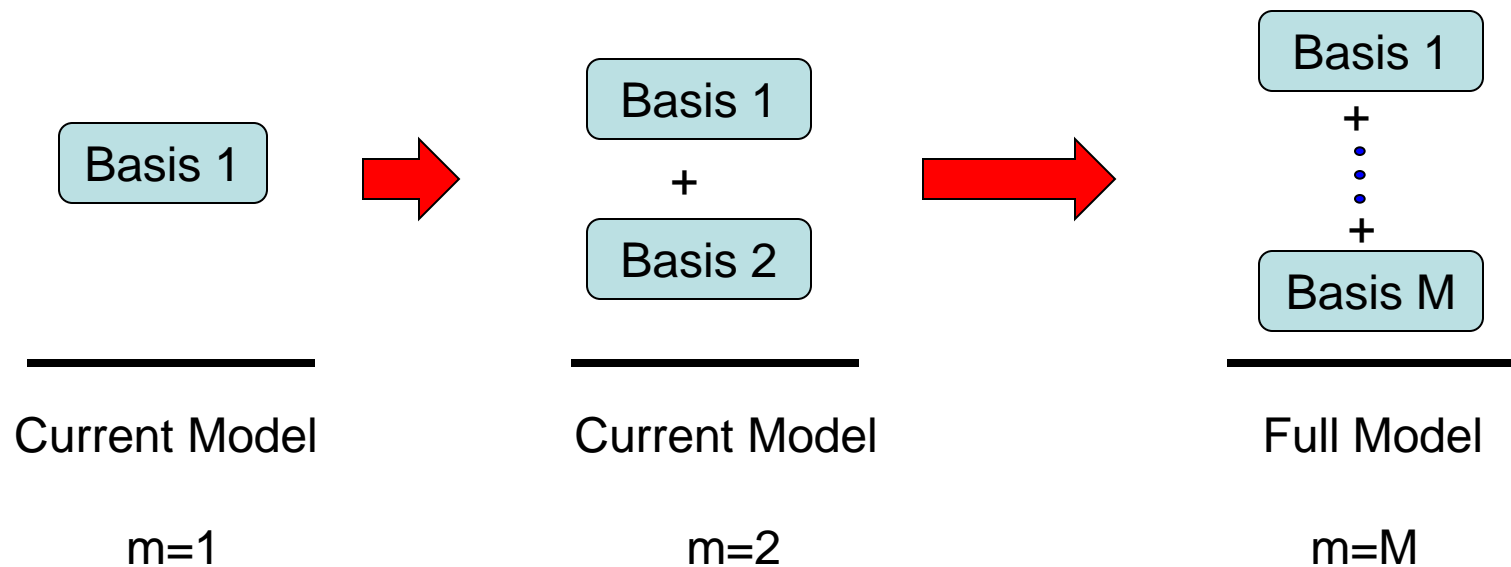
# Flow Diagram for Leaps BIC Selection

---



# Generalized Boosted Regression Model (GBM)

- Regression method that iteratively combines weak prediction models into one strong model
- Fit regression using current model, search basis function based on minimizing loss function, then update current model by adding basis function



# Basis and Loss Function

Full Model	$f(x) = \sum_{m=1}^M \beta_m b(x)$
Loss Function	$\min \sum_{i=1}^N L(y_i, \beta b(x))$ $L(y, f(x)) = (y - f(x))^2$

$b(x)$  represents basis function

$\beta_m$  represents corresponding coefficient



# Advantages of GBM

---

- Easy to compare between the 2 techniques (GBM and MLR)
  - Uses same predictors
  - Uses the same loss function
- Applicable on many techniques
  - Regression, classification, exponential, decision tree, support vector machine





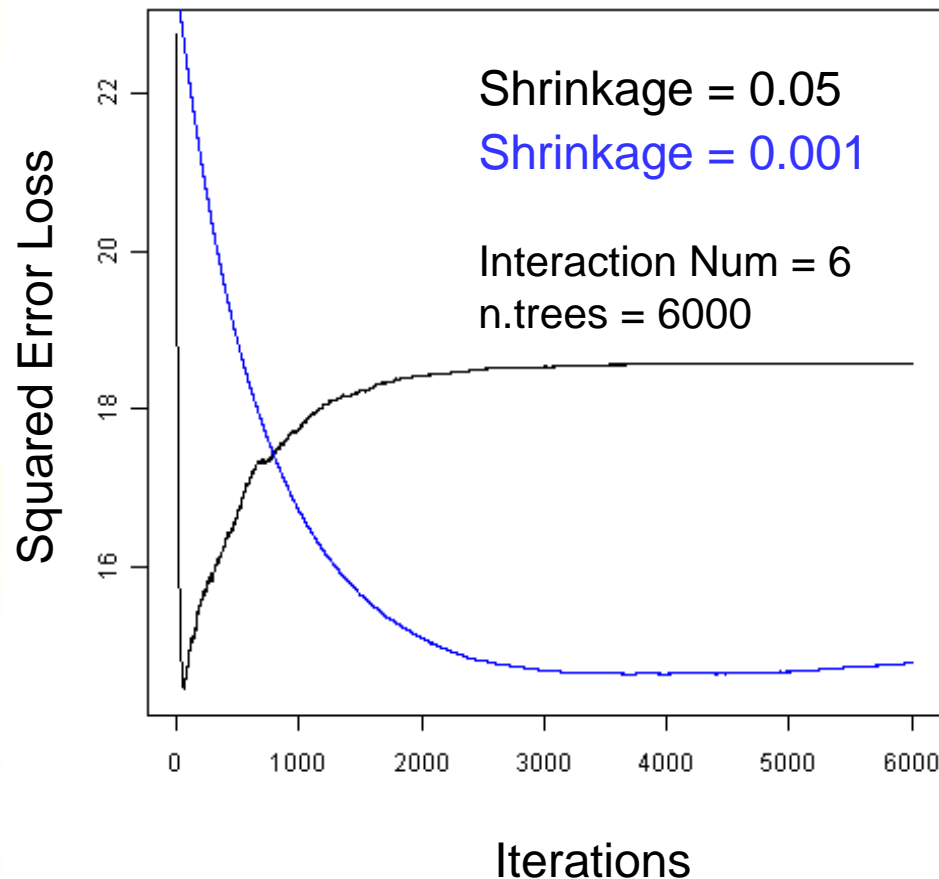
# GBM Parameters

	Distribution	Shrinkage	Interaction .depth	N.trees
Description	Specifying Loss Function	Learning Rate / Weight for Basis Function	Num. of Variable Interactions	Iterations
Recommended Specification	Gaussian	0.01 – 0.001	4 - 8	1000 - 10000

Rule of Thumb: # Trees x Shrinkage = [1:100]



# Optimizing Shrinkage Parameter



Advantages  
of Low  
Shrinkage

More  
consistent  
error loss

improves  
Predictive  
Performance  
for new cases

Avoids  
Overfitting

Disadvantages  
of Low  
Shrinkage

Computation  
Cost



Environment  
Canada

Environnement  
Canada

Canada

# Optimizing Interaction Depth and Number of Iterations

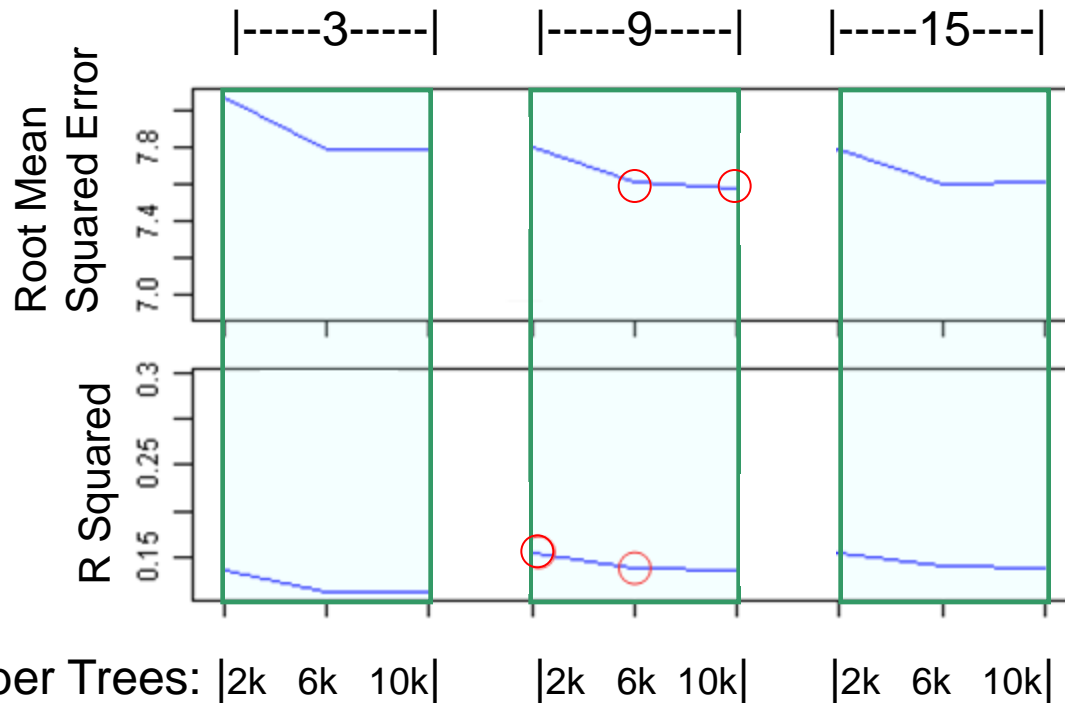
Station: 30120

Hour: 29

Pollutant: O3

Shrinkage: 0.001

## Interactions



Above 75<sup>th</sup> percentile



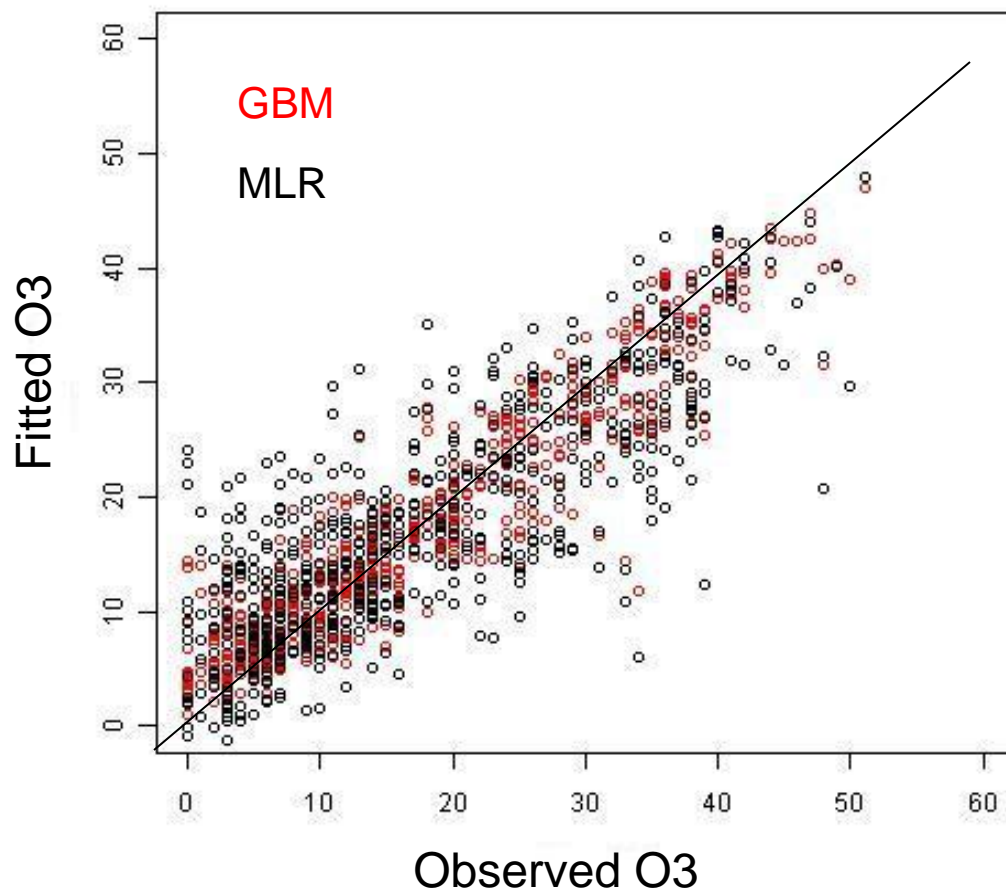
Environment  
Canada

Environnement  
Canada

Canada

# Improvements in Fitted Model

Training Data for Station: 30120 Hour: 29



Shrinkage: 0.001

Iterations: 6000

Interactions: 9



Environment  
Canada

Environnement  
Canada

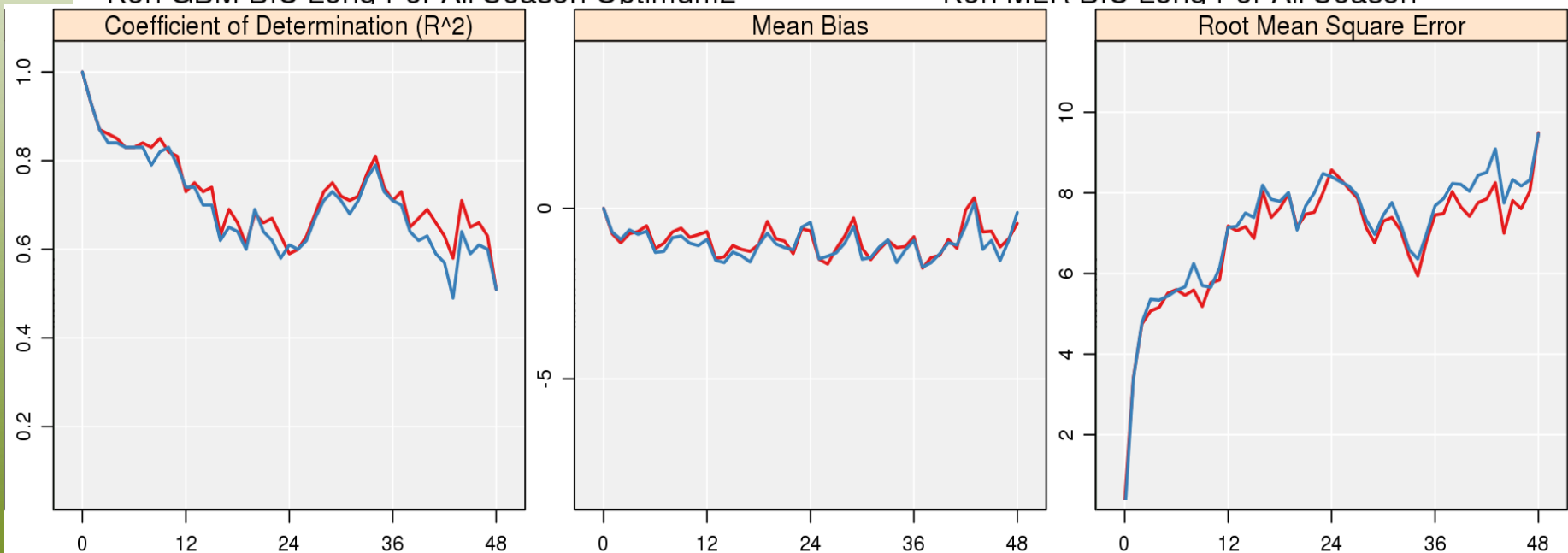
Canada

# Verification Results for all cases

## Statistics TEST Set All Season GEM15 Hourly Forecast of O3 for 75% Quantile Station 00030120

Ken GBM BIC Long Per All Season Optimum2

Ken MLR BIC Long Per All Season



Shrinkage: 0.001 Iterations: 6000 Interactions: 9



Environment  
Canada

Environnement  
Canada

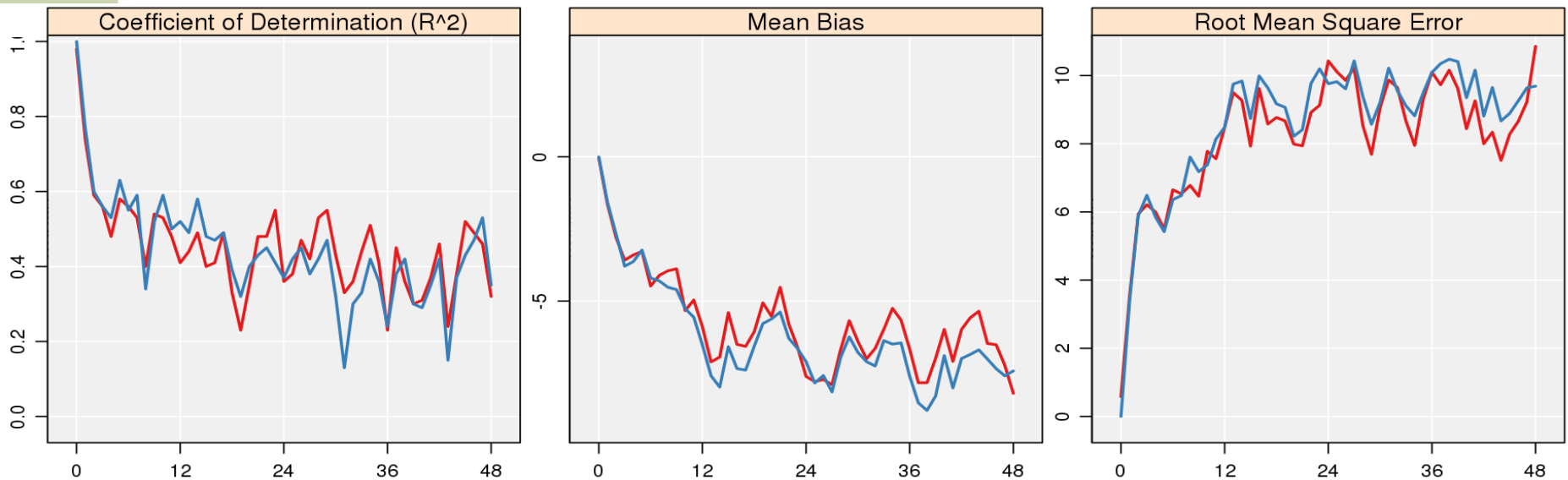
Canada

# Verification Results for Above 75<sup>th</sup> Percentile cases

## Statistics TEST Set All Season GEM15 Hourly Forecast of O3 for 75% Quantile Station 00030120

Ken GBM BIC Long Per All Season Optimum2

Ken MLR BIC Long Per All Season



Shrinkage: 0.001   Iterations: 6000   Interactions: 9



Environment  
Canada

Environnement  
Canada

Canada

# Conclusion and Future Work

---

## Conclusions:

- GBM is an iterative regression method that searches for basis functions to minimize a selected loss function
- With proper optimization of parameters, GBM is capable of producing marginal improvements over MLR

## Future Work:

- `gbm.more` function to find the optimal number of iterations
- Boosted Regression with other loss functions

