

# Model Training and Validation with caret

Alexander C. Mueller, PhD

September 14, 2017

# What is caret? Why should I use it?

## What?

- It's an R library!
- **classification and regression training**
- caret contains tools for creating predictive models.
- Created by Max Kuhn at Pfizer.

## Why?

- Access many algorithms from one interface.
- Handle preliminary work uniformly.
- Focus on training and testing.

## Links:

- <http://topepo.github.io/caret/index.html>
- [https://github.com/BirdMueller/rug/blob/master/caret\\_walkthrough.R](https://github.com/BirdMueller/rug/blob/master/caret_walkthrough.R)

# What is covered in this talk?

## Part 1: Motivating caret

- Subsample data for testing and training.
- How much model complexity is desirable?
- Fear noise and avoid learning about it.

## Part 2: Working with caret

- Play with the Titanic dataset.
- Spotlight some key functions in caret.
- Create and compare some binary predictive models using different algorithms.
- Philosophy of working with caret.

## An Illustrative Example

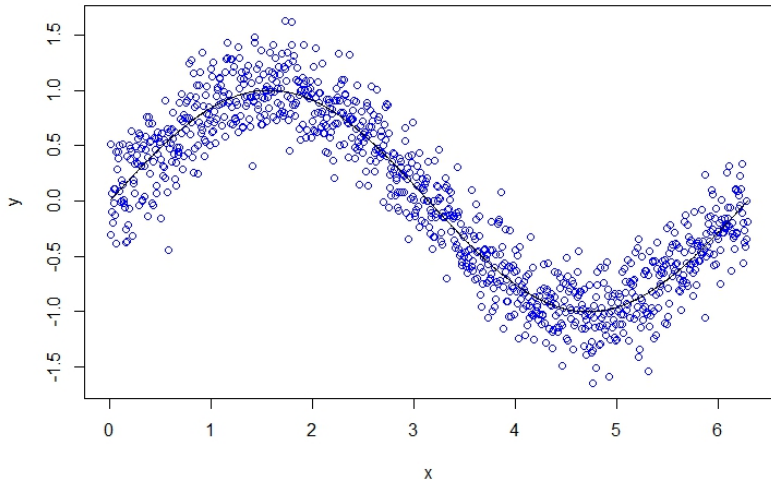
Train a model to compute the sin function given some sample points.

- Noise has been added to represent real-world measurement errors.
- We know  $\sin(x)$  has a Taylor series  $x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$
- The powers  $x$ ,  $x^2$ ,  $x^3$  of  $x$  should work as features.

We already know the perfect-world, zero-noise answer. The point of this exercise is to examine the effect of the noise.

# Learning sine

$$y = \sin(x)$$



## Learning sine

```
13 #our number of points
14 N<-1000
15 #spread x coords out evenly in (0,2*pi]
16 xRaw<-(1:N)*2*pi/(N)
17 xNoise<-0
18 #compute y = sin(x) and "measurement errors"
19 yRaw<-sin(xRaw)
20 yNoise<-0.25*rnorm(N)
21 #compute final x and y coordinates
22 xCoord<-xRaw+xNoise
23 yCoord<-yRaw+yNoise
```

## How much complexity?

Model complexity, and how to get the right amount of it, is key.

- Too little complexity and you're not really trying.
- Too much complexity is idle hands turning to evil deeds.
- Complex models memorize noise, relate new data to dumb memories of noise.

Our models of sine will provide a window into these issues.

- For our polynomial model, complexity is the degree (one less than the number of features).
- Degree 0 is the constant polynomial  $f(x) = 0$ , not very sine-like.
- We'll next explore just what goes wrong in high degree.

# Training Error

```
43 #plot training error with increasing degree 1 to 100
44 plot(
45     1:100,
46     sapply(1:100,function(i){
47         regrModel<-lm(regr_fm1a(i),higherPowers)
48         return((1/N)*sum(abs( yCoord-predict(regrModel,higherPowers) )))
49     }),
50     xlab = "Polynomial Degree (Number of Independent Variables)",
51     ylab = "Sum of Absolute Values of Errors",
52     main = "Training Error",
53     col='blue'
54 )
```

For each degree  $i$

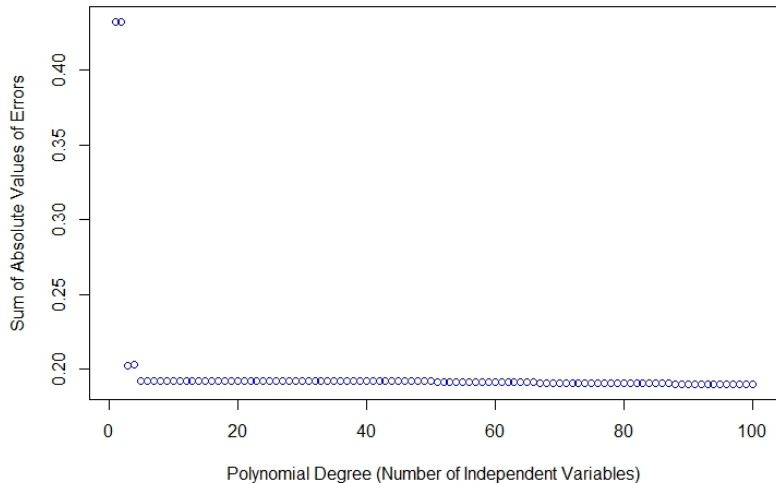
- train a regression model,
- measure its accuracy (the  $L^1$  error is non-standard),
- and plot!

What do you expect this plot to look like?

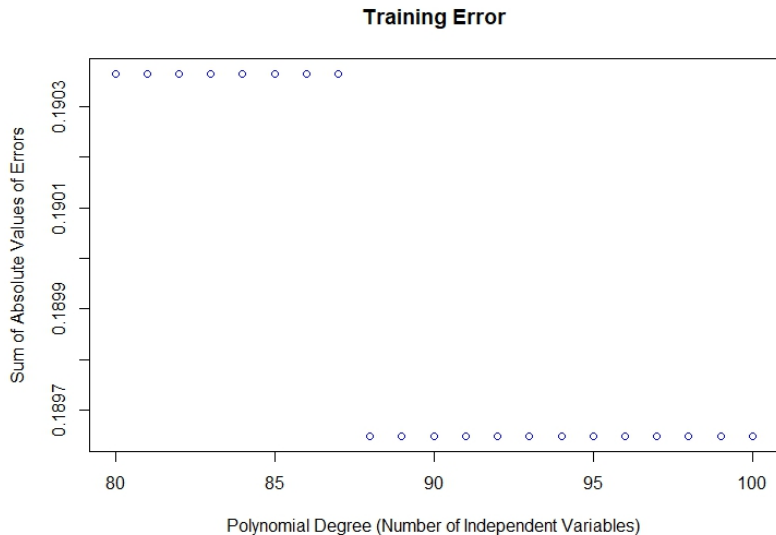


# Training Error

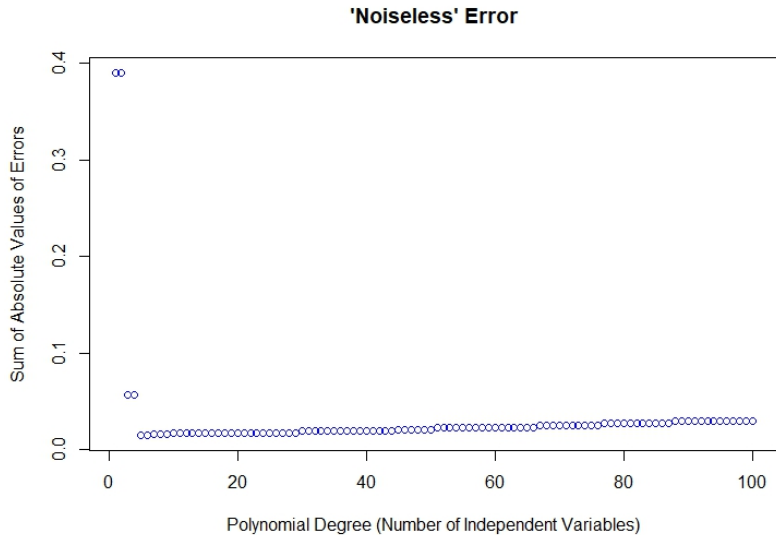
**Training Error**



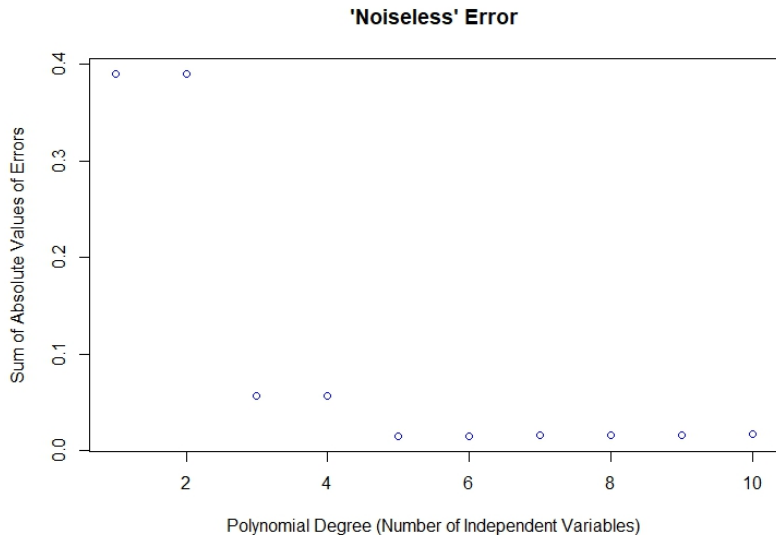
## Training Error, High Complexity (Degree)



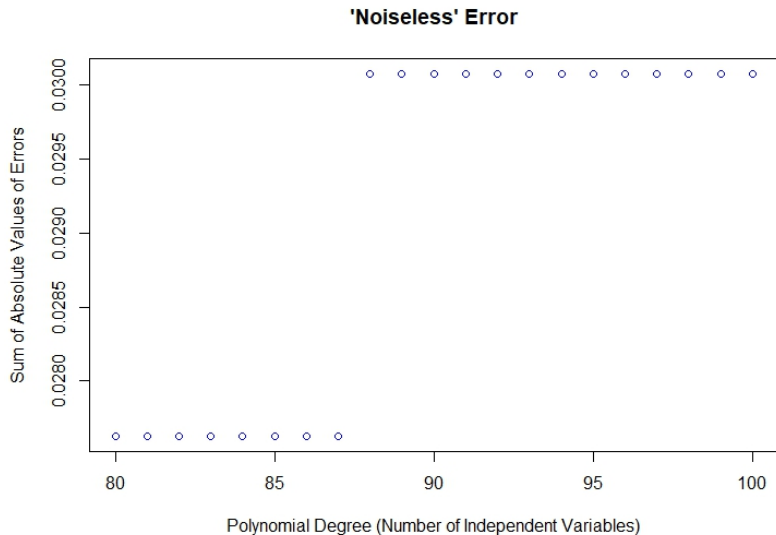
# “Noiseless” Error



# “Noiseless” Error, Optimal Complexity (Degree)



## “Noiseless” Error, High Complexity (Degree)



## Testing on “New” Data

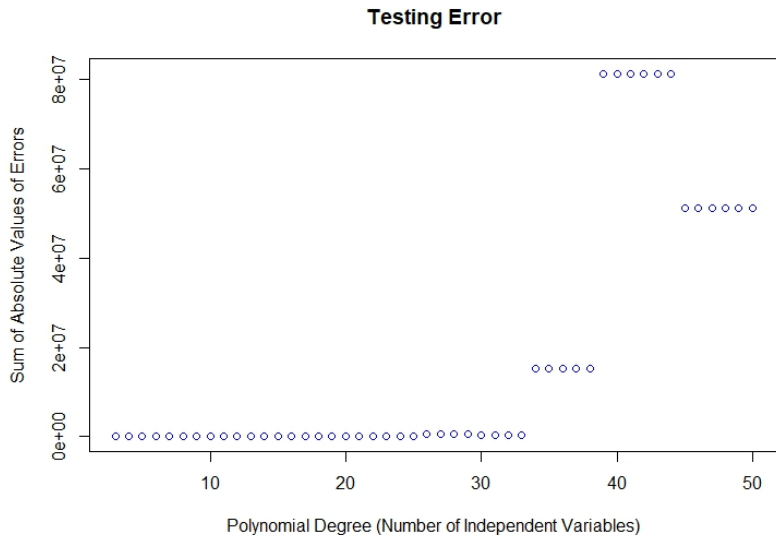
```
113 #take the first three quarters of my N points
114 inTrain<-1:floor(0.75*N)
115 #create data frames for training and then testing afterwards
116 training<-higherPowers[inTrain,]
117 testing<-higherPowers[-inTrain,]
```

How will our models generalize to new data?

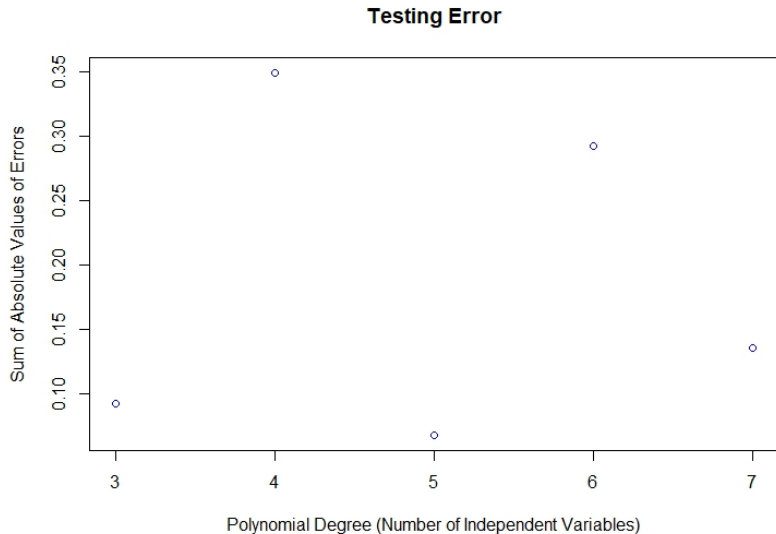
- Separate the points with  $x$  in a smaller interval.
- Repeat our training loop for these points.
- Compute  $y$  values for the remaining “test” points.

What will be the degree of the most accurate polynomial?

# Testing Error, High Complexity (Degree)



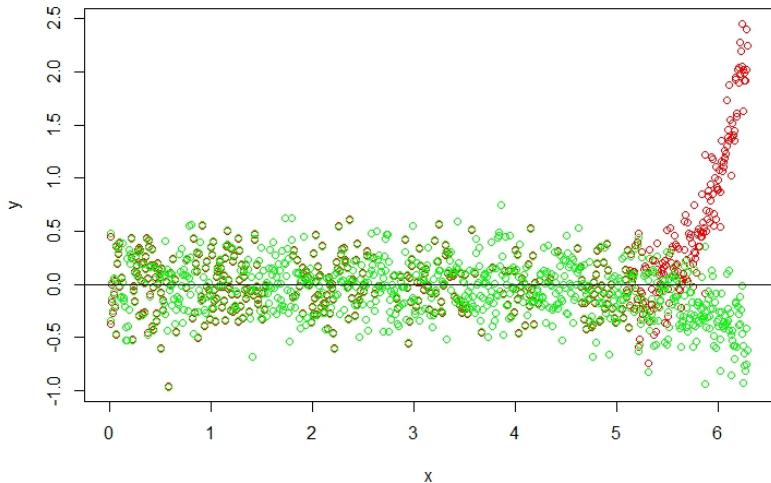
# Testing Error, Optimal Complexity (Degree)





# Residues for the Optimal Model vs. More Complex (Higher Degree)

**Residues, Degree 5 vs Degree 7**



# Training Error vs. “Noiseless” Error

Some important, vague, truthy ideas:

- **Overtraining** is when model complexity memorizes noise.
- Complexity without meaning inevitably does complex, meaningless things to new data.
- A good model is smart enough to learn the signal, too dumb to learn the noise.

Why do the “noiseless” and testing errors increase while training error continues to decrease?

# Training vs. Testing Data

In practice, model building typically involves training many models of varying complexity, evaluating their accuracy, and finding the sweet spot.

- **Testing data** is excluded from the training process so the model can't learn its noise.
- “Noiseless” data is not available, but testing data can serve the same purpose.
- Evaluate a candidate model according to its behavior on the testing data.

## Enter caret

caret handles shared model-building concerns in a uniform way over many types of model. We'll use

- **createDataPartition(...)** for splitting datasets,
- **train(...)** for training models,
- the **grid** and **trControl** options, and
- **confusionMatrix(...)** for evaluation

as we build some models on the Titanic dataset. This dataset contains information about passengers on the famously doomed ocean liner Titanic including their fate. Life and death is the ultimate binary outcome.

```
> titanic_train[1,]  
 PassengerId Survived Pclass      Name Sex Age SibSp Parch  Ticket Fare Cabin Embarked  
1           1         0       3 Braund, Mr. Owen Harris male   22     1     0 A/5  21171  7.25      S
```

## The Titanic Dataset

```
170 titanic_train$outFactor<-as.factor(  
171   sapply(titanic_train$Survived,function(row){  
172     if(row){  
173       return('winslet')  
174     }else{  
175       return('DiCaprio')  
176     }  
177   })
```

Make sure your outcome variable is a factor or caret may get confused about what sort of model you are building. This dataset comes with “titanic\_train” and “titanic\_test” data frames but we will work only with the former.

## Splitting data

```
179 #find a random three quarters subset of all indices
180 inTrainTitanic<-createDataPartition(
181   y=titanic_train$outFactor,
182   p=0.75,
183   list=FALSE
184 )
185 #create data frames of testing and training data
186 trainTitanic<-titanic_train[inTrainTitanic,]
187 testTitanic<-titanic_train[-inTrainTitanic,]
```

createDataPartition splits our data into testing and training sets according to a proportion we input. It is necessary to tell the function which is your outcome variable so it can ensure outcomes are distributed evenly between the two groups. In this case, three quarters of the data goes to the training set.

## Train a Random Forest

```
189 rf <- train(  
190   outFactor ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked,  
191   data = trainTitanic,  
192   method = "rf"  
193 )  
194 confusionMatrix(  
195   data = predict(rf,newdata = testTitanic),  
196   na.omit(testTitanic)$outFactor  
197 )
```

The string "rf" is the only part of this code in any way unique to the random forest algorithm. At time of writing, 238 algorithms can be called (from other R libraries) using `train(...)` and each has a unique code.

# Random Forest Under the Hood

```
> rf
Random Forest

669 samples
12 predictor
2 classes: 'DiCaprio', 'Winslet'

No pre-processing
Resampling: Bootstrapped (25 reps)
summary of sample sizes: 542, 542, 542, 542, 542, 542, ...
Resampling results across tuning parameters:
```

mtry	Accuracy	Kappa	Accuracy SD	Kappa SD
2	0.8107367	0.5994673	0.02389990	0.05013360
5	0.7984620	0.5779757	0.03371800	0.07040022
9	0.7846143	0.5501639	0.03258991	0.06775746

Accuracy was used to select the optimal model using the largest value.  
The final value used for the model was mtry = 2.

In accordance with our theme, random forest works by training models of varying complexity on subsampled data.



## Confusion Matrix and Statistics

	Reference	
Prediction	DiCaprio	winslet
DiCaprio	90	22
winslet	9	51

Accuracy : 0.8198

95% CI : (0.754, 0.8741)

No Information Rate : 0.5756

P-Value [Acc > NIR] : 8.361e-12

Kappa : 0.6223

McNemar's Test P-Value : 0.03114

Sensitivity : 0.9091

Specificity : 0.6986

Pos Pred Value : 0.8036

Neg Pred Value : 0.8500

Prevalence : 0.5756

Detection Rate : 0.5233

Detection Prevalence : 0.6512

Balanced Accuracy : 0.8039

'Positive' Class : DiCaprio

## Some Nice Features

caret will automatically import libraries needed for models called via `train(...)`. I have thus far only explicitly referenced three libraries.

```
2 library(caret)
3 library(mlbench)
4 library(titanic)
```

caret will even install libraries automatically (with your permission) if necessary.

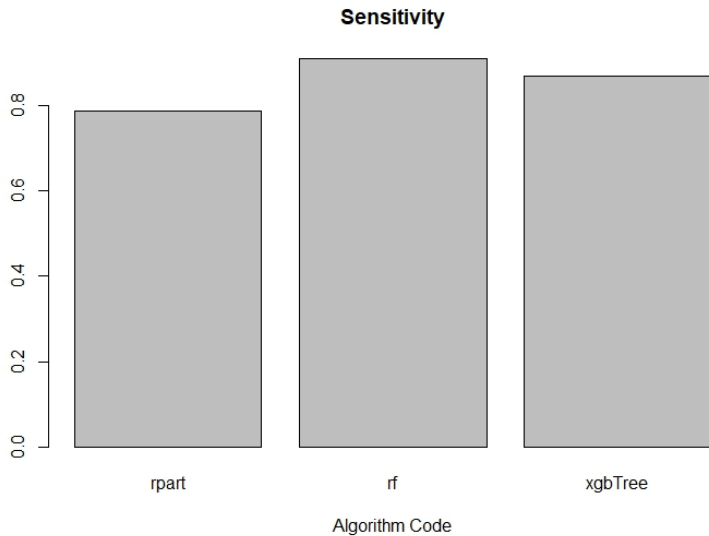
```
1 package is needed for this model and is not installed. (rFerts). Would you like to try to install it now?
1: yes
2: no
```

# Model Selection

```
200 models<-c('rpart','rf','xgbTree')
201 modelList<-lapply(models,function(name){
202   return(
203     train(
204       outFactor ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked,
205       data = trainTitanic,
206       method = name
207     )
208   )
209 })
```

One benefit of a uniform interface is that allows easy iteration over a group of candidate algorithms. In this case, this is as simple as iterating a nearly uniform call to `train(...)` over a vector of the relevant code strings.

# Model Selection Results



# Model Tuning

```
225 grid <- expand.grid(  
226   mtry=c(2,3)  
227 )  
228 ctrl <- trainControl(  
229   method = "repeatedcv",  
230   repeats = 3,  
231   classProbs = TRUE,  
232   summaryFunction = twoClassSummary  
233 )  
234 rfTwo <- train(  
235   outFactor ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked,  
236   data = titanic_train,  
237   method = "rf",  
238   trControl = ctrl,  
239   tuneGrid = grid  
240 )
```

The **tuneGrid** and **trControl** options provide additional control over the training process. Once again, the interface is uniform although not every parameter is relevant to every model.

# Model Tuning

Random Forest

891 samples  
12 predictor  
2 classes: 'Dicaprio', 'winslet'

No pre-processing

Resampling: Cross-validated (10 fold, repeated 3 times)

summary of sample sizes: 643, 643, 643, 642, 642, 642, ...

Resampling results across tuning parameters:

mtry	ROC	Sens	Spec	ROC SD	Sens SD	Spec SD
2	0.8679272	0.9103359	0.6655172	0.04565705	0.04116911	0.08157597
3	0.8679513	0.8889996	0.7068966	0.04381717	0.05336908	0.08953161

ROC was used to select the optimal model using the largest value.  
The final value used for the model was mtry = 3.

Note the choices of mtry and the new-cross validation procedure.

# Practical and Philosophical Advice

No tool is an end all be all. caret is great when

- doing model selection.
- sizing up the predictive power of a dataset.
- cross-validation concerns are key.

caret might get in the way when

- a decision has been made about what algorithm to use.
- when significant and specific parameter tuning is required.
- algorithm specific issues are key.

Good software, like good mathematical notation, should by its structure what is important and how to appropriately focus on it.