

Statistical programming languages - *House Prices: Advanced Regression Techniques*

Christian Koopmann

Felix Skarke

Enno Tammena

Ladislaus von Bortkiewicz Chair of Statistics
Humboldt-Universität zu Berlin

<http://lvb.wiwi.hu-berlin.de>



Motivation

- demonstrate some of R's machine learning capabilities
- develop and test different models to predict the price of houses in Ames, Iowa
- dataset: Ames Housing data published in a *kaggle* competition with 79 explanatory variables

Outline

Exploratory Data Analysis

Exploratory Data Analysis Dependence

Data Preprocessing

Regression Models

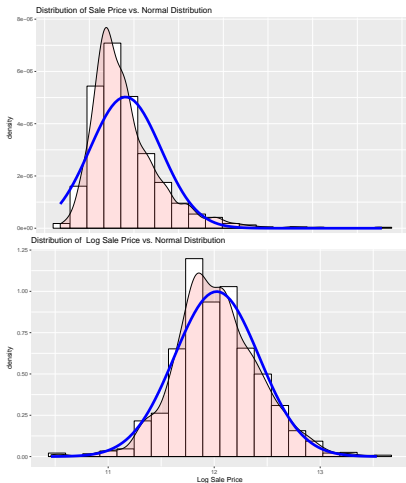
Random Forest

Model Comparison

Data Structure / Quality

- 2019 Observations separated into 1460 labelled training observations and 1459 unlabelled test observations to predict.
- Target Variable: Final Sale Price in USD
- 79 Explanatory Variables of which 36 are numeric (e.g. floor area) and 43 categoric (e.g. heating type).
- 19 variables contain NA's 16 of which are categoric variables
- NA usually means a variable is not applicable (Pool Quality if house has no pool) instead of missing data

Sale Price Distribution



- Original price distribution heavily skewed
- Log transformed House Prices seem to be normally distributed

Quantlet: Exploratory Data Analysis: Dependence

- aim of this quantlet: show relations between variables, especially the target variable SalePrice
- producing different graphical representations
 - ▶ correlation matrix
 - ▶ barplot
 - ▶ boxplots

Correlation Matrix of all numeric variables: Code

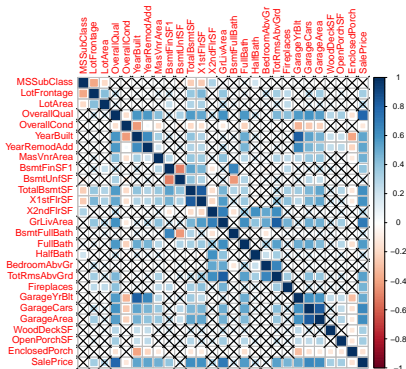
```
1 corr.func = function(data, cut.value, corr.mat = FALSE, corr.test =  
  FALSE, significance = 0.05) {  
2   corr.numeric = cor(na.omit(numeric.data))  
3   find.rows = apply(corr.numeric, 1, function(x) sum(abs(x) > abs(cut.  
    value)) > 1)  
4   corr.numeric.adjusted = corr.numeric[find.rows, find.rows]  
5   low.corr = colnames(corr.numeric) %in% colnames(corr.numeric.adjusted)  
6  
7   pdf("Corrplot.pdf")  
8   if (corr.test == FALSE) {  
9     corrplot(corr.numeric.adjusted, method = "square")  
10  } else {  
11    corrplot(corr.numeric.adjusted,  
12    p.mat = correlation.test(corr.numeric.adjusted),  
13    sig.level = significance, method = "square")  
14  }  
15  dev.off()  
16  
17  if (corr.mat == TRUE)  
18    return(corr.numeric.adjusted)  
19 }
```



Correlation Matrix of all numeric variables: Code ctd.

```
21 correlation.test = function(corr.data) {  
22   corr.data      = as.matrix(corr.data)  
23   n              = ncol(corr.data)  
24   p.value.matrix = matrix(NA, n, n)  
25   diag(p.value.matrix) = 0  
26  
27   for (i in 1:(n - 1)) {  
28     for (j in (i + 1):n) {  
29       tmp = cor.test(corr.data[, i], corr.data[, j])  
30       p.value.matrix[i, j] = p.value.matrix[j, i] = tmp$p.value  
31     }  
32     colnames(p.value.matrix) = rownames(p.value.matrix) =  
33                               colnames(corr.numeric.adjusted)  
34   }  
35   return(p.value.matrix)  
36 }
```


Plot of Correlations (Corrplot)



- A cut-off value of 0.3 was used: Variables with no correlation over 0.3 do not show up in the plot
- A significance level of 0.5 was used to test the correlations (crosses indicate no significance)

Barplot: Code

```
36 corr.barplot = function(numb.corr = 36) {  
37   correlation.vars = names(numeric.data) %in% c("SalePrice")  
38   correlation.data = numeric.data[!correlation.vars]  
39   correlations = vector(length = length(names(correlation  
40     .data)))  
41   names(correlations) = names(correlation.data)  
42   for (i in names(correlation.data)) {  
43     correlations[i] = cor(numeric.data$SalePrice, correlation.  
44       data[i], use = "pairwise.complete.obs")  
45  
46     y.plotting = correlations[order(abs(correlations),  
47       decreasing = TRUE)][1:numb.corr]  
48     x.plotting = names(y.plotting)  
49     names(y.plotting) = NULL  
50     df = data.frame(x.plotting, y.plotting)  
51     df$x.plotting = factor(df$x.plotting, levels =  
52       df[order(df$y.plotting, decreasing = TRUE), "x.plotting"])  
53  
54     ggplot(data = df, aes(x.plotting, y.plotting), fill = as.factor  
55       (x.plotting)) + geom_bar(stat = "identity") +  
56       theme(axis.title.x = element_blank(), axis.text.x = element  
57         _text(angle = 90, vjust = 0.5, size = 10)) +  
58       ylab("Correlation") + ggtitle(paste("Barplot of the", numb.  
59         corr, "highest bivariate correlations with SalePrice",  
60         sep = " "))  
61   }
```



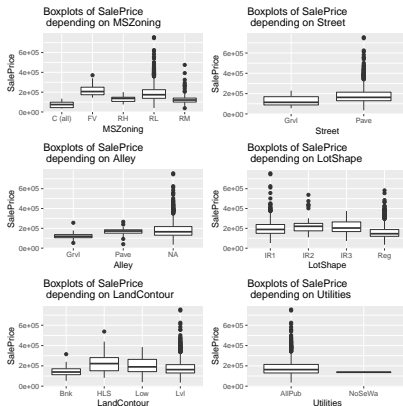
Ordered barplot for the correlations with SalePrice

- The plot shows the 20 highest correlated numeric variables
- The few negative correlations are not plotted, since they are very low

Boxplot: Code

```
58 boxplot.target = function(categoric) {  
59   categoric.x = data[, categoric]  
60   plot.data = as.data.frame(cbind(data$SalePrice, categoric.x))  
61   plot.data[[2]] = as.factor(plot.data[[2]])  
62   levels(plot.data[[2]]) = levels(categoric.x)  
63  
64   ggplot(plot.data, aes(x = categoric.x, y = V1)) + geom_boxplot  
65     () +  
66     labs(title = paste("Boxplots of SalePrice",  
67       "\n", "depending on", categoric, sep = " "),  
67       x = categoric, y = "SalePrice")  
68 }
```

Boxplots of SalePrice

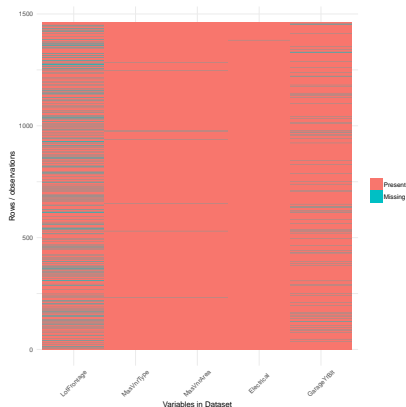


- Example of boxplots for SalePrice based on the levels of categorical variables
- Differences in median and spread across levels indicate a possibly good predictor variable

Quantlet: Data Preprocessing

- aim of this quantlet: ensure, that the data quality is sufficient for all types of models that will be used
- handling of missing data, including imputation
- reduction of dimensionality
 - ▶ merging of factor levels
 - ▶ principal component analysis
- detection and handling of outliers

Missing Data



- based on data description, most NA's have a meaning (None/Other)
- imputation of remaining NA's
 - ▶ numeric variables: median
 - ▶ factor variables: mode

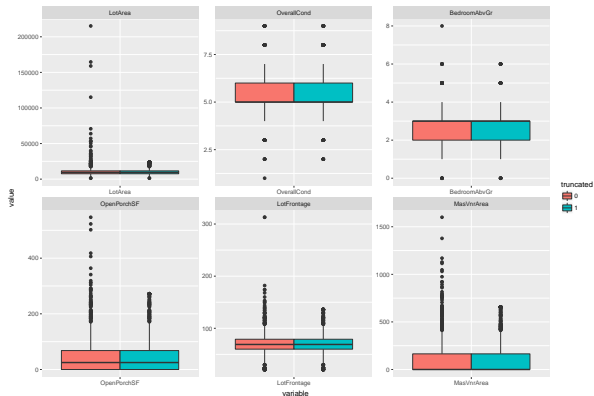
Missing Data: Code

```
1 Mode = function(x) {  
2     ux = unique(x)  
3     ux[which.max(tabulate(match(x, ux)))]  
4 }  
5 impute.mode = function(x){  
6     nas      = is.na(x)  
7     x[nas] = Mode(x[!nas])  
8     as.factor(x)  
9 }  
10 categorical.imputed = as.data.frame(sapply(categorical.  
    data, impute.mode))
```


Outlier Handling: Code

```
1 outlier.count = function(x){
2     low =as.numeric(quantile(x)[2] - IQR(x)*3)
3     high=as.numeric(IQR(x)*3 + quantile(x)[4])
4     sum(x >= high | x <= low)
5 }
6 outlier.truncate = function(x){
7     low =as.numeric(quantile(x)[2] - IQR(x)*3)
8     high=as.numeric(IQR(x)*3 + quantile(x)[4])
9     x[x < low] = low
10    x[x > high] = high
11    return(x)
12 }
13 df_outlier_trunc = as.data.frame(sapply(df.temp.
    numeric, outlier.truncate))
```

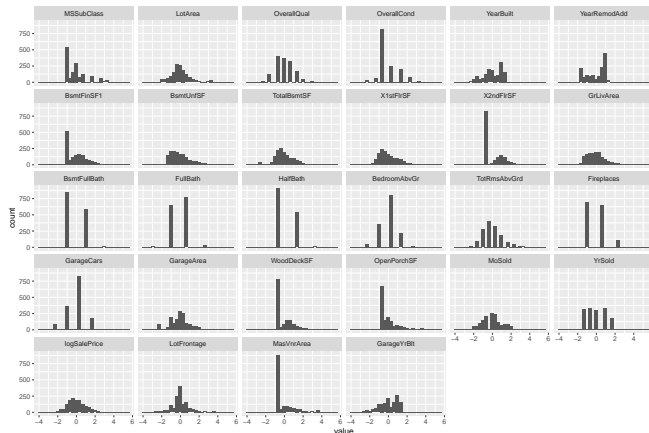
Outlier: Before and after truncation



Dimensionality Reduction: Code

```
1 single.factors = function(data) {  
2   for(var in names(data)) {  
3     if(is.factor(data[[var]])){  
4       tbl = table(data[[var]])  
5       ren = names(tbl)[tbl <= 20]  
6       levels(data[[var]])[levels(data[[var]])  
7         %in% ren] = "Other"  
8       tbl      = table(data[[var]])  
9       tbl_sum = sum(tbl < 20)  
10      if(nlevels(data[[var]]) < 3 & tbl_sum >=  
11        1 ) data[[var]] = NA  
12    }  
13  }  
14  return(data)  
15 }
```

Preprocessing: Summary



Quantlet: Regression Models

- aim of this quantlet: select appropriate variables
- four selection procedures
 - ▶ Backwards Selection based on significance
 - ▶ Forward Selection based on AIC
 - ▶ LASSO
 - ▶ Ridge

Significance Selection: Code

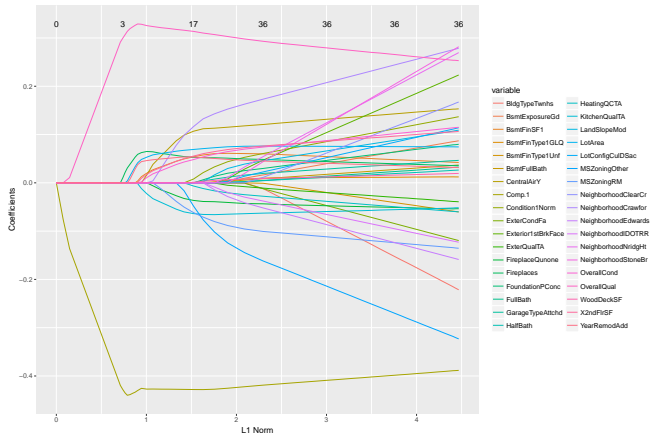
```
1 sign.select = function(dframe, y) {  
2   pvals = 1, z = 1, i = 1  
3   vars = names(dframe)  
4   vars = vars[!vars %in% y]  
5   while(z>0){  
6     df.lm = cbind(dframe[vars], dframe[y])  
7     lm1 = lm(formula(paste(y, "~ . ")), data=df.lm)  
8     pvals = summary(lm1)$coefficients[,4]  
9     pvals = pvals[!names(pvals) %in% "(Intercept)"]  
10    vars = names(pvals[pvals<0.05])  
11    #z = sum(pvals>0.05), i = i+1  
12    if(i==300){  
13      warning("... No signif. Vars in Data Set?")  
14      break  
15    }  
16  }  
17  return(vars)  
18 }
```



LASSO: Code

```
1 lm.penal = function(type, x, y) {  
2   if (type == "lasso") {           alpha = 1  
3   } else if ( type == "ridge") {alpha = 0  
4   } else stop("type must be ridge or lasso")  
5   cvfit      = cv.glmnet(x, y, alpha = alpha,  
6     nolds = 10)  
7   fit        = predict(cvfit,newx=x,s="lambda.1se")  
8   rsq        = 1 - sum(y^2)/sum((fit - y)^2)  
9   c          = coef(cvfit, s = "lambda.1se")  
10  inds       = which(c != 0)  
11  var = row.names(c)[inds]  
12  vars.sele = var[!variables %in% "(Intercept)"]  
13  coeftable = data.frame(var = var, coeff = c[inds  
14    ], stringsAsFactors = F)  
15  output    = list(vars.sele, coeftable, c, cvfit,  
16    fit, rsq)  
17 }
```

LASSO Results



Regression Models Results (Excerpt)

.variable	.stat	Sign. Selec.	AIC Selec.	Lasso	Ridge
(Intercept)	Estimate	-0.849***	-0.914***	-0.116	-0.229
LotArea	Estimate	0.112***	0.062***		-0.019
OverallQual	Estimate	0.245***	0.201***	0.076	0.057
OverallCond	Estimate	0.089***	0.104***	0.293	0.118
YearRemodAdd	Estimate	0.049***	0.044***	0.074	0.047
BsmtUnfSF	Estimate	-0.066***	-0.165***	0.045	0.052
X2ndFlrSF	Estimate	0.208***	0.215***	0.06	0.054
LotShapeOther	Estimate	-0.415***	-0.345***		0.034
LotConfigCulDSac	Estimate	0.141***	0.136***	0.048	0.043
LandSlopeMod	Estimate	0.119**	0.137***	0.009	0.029
NeighborhoodClearCr	Estimate	0.184**	0.225***		0.026
NeighborhoodCrawfor	Estimate	0.323***	0.291***		0.003
NeighborhoodEdwards	Estimate	-0.153***	-0.199***		-0.004
NeighborhoodIDOTRR	Estimate	-0.209***	-0.199***		0.033
ElectricalOther	Estimate				0.051
ElectricalSBrkr	Estimate				0.035
	N	1168	1168	1168	1168
	R2	0.895	0.914	0.88	0.89
	Number Vars	33	77	36	177

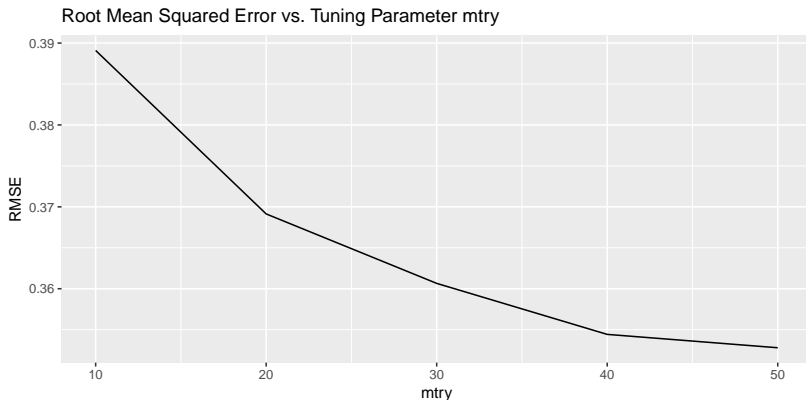
Random Forest

- ▣ Applied Random Forest Regression Model to the data
- ▣ Collection of Regression Trees, where the predicted value is the average value in the leaf node.
- ▣ Results are averaged across all Regression Trees in the Forest
- ▣ Tuning Parameter: *mtry* - Number of randomly selected variables per decision tree
- ▣ Model is tuned using the *caret* package and 5 Fold Cross Validation in a parallelized fashion using *doParallel*

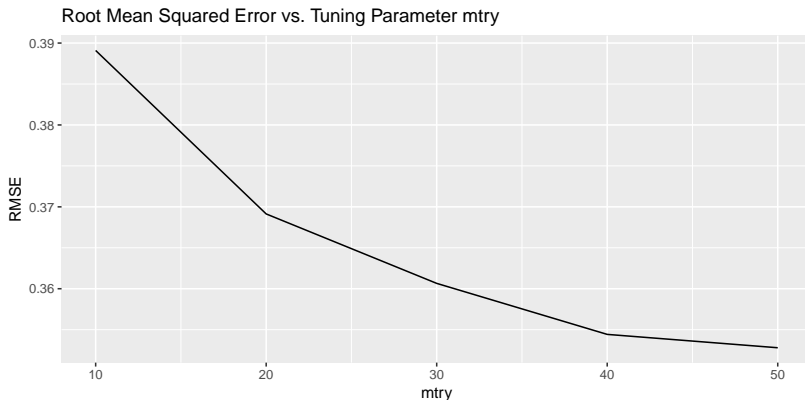
Random Forest: Code

```
1 CrossValidation = trainControl(method = "repeatedcv"  
  , number = nfolds, p = ptrain, repeats = nrepeats)  
2 rfGrid          = expand.grid(mtry = mtry.tuning)  
3 cores           = max(c(detectCores() - 1, 0))  
4 cl              = makeCluster(cores)  
5 registerDoParallel(cl)  
6 set.seed(123)  
7 rftuned = train(log(SalePrice) ~ ., data = df,  
  method = "rf", importance = TRUE, trControl =  
  CrossValidation, verbose = TRUE,  
8   tuneGrid = rfGrid)  
9 stopCluster(cl)
```

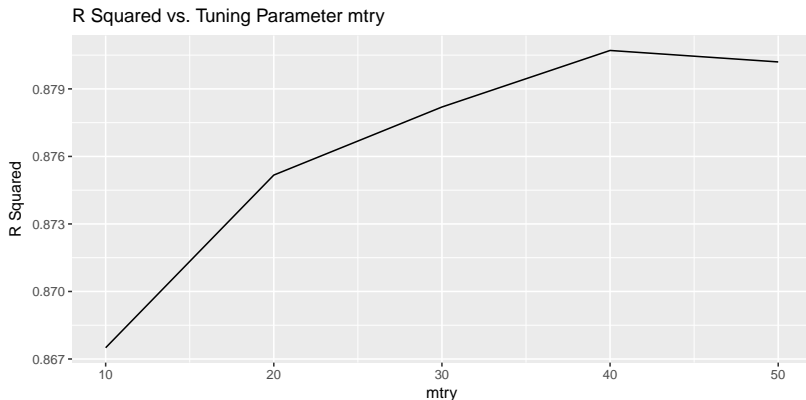
Random Forest Tuning Results - RMSE



Random Forest Tuning Results - RMSE



Random Forest Tuning Results - R Squared



Quantlet: Model Comparison

- aim of this quantlet: after building and training all the models on the training data the goal is to measure how they perform on new data (=test data)
- computing different measures and show results graphically
 - ▶ Mean Squared Error
 - ▶ Bias
 - ▶ real Vs. estimated plots

Performance measure example: MSE Code

```
1 model.mse = function(model, test.data = test) {  
2   if (class(model)[1] %in% c("train", "lm")) {  
3     pred = predict(model, newdata = test.data)  
4     mse = (1/ncol(test.data)) * sum((pred -  
5       test.data$logSalePrice)^2)  
6   } else {  
7     pred = predict(model, newx = as.matrix(test.  
8       data[!names(test.data) %in% "logSalePrice"  
9       ]), s = "lambda.1se")  
10    mse = (1/ncol(test.data)) * sum((pred -  
11      test.data$logSalePrice)^2)  
12  }  
13  return(mse)  
14 }
```


Result comparing all models

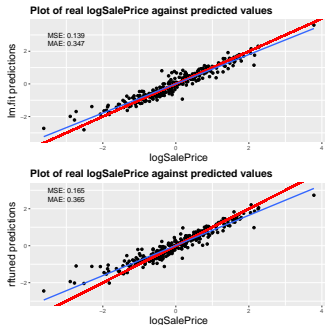
	lm	fwd	lasso	ridge	gbm	rf
MSE	0.14	0.12	0.19	0.19	0.14	0.16
MSEtrain	0.70	0.57	0.94	0.88	0.48	0.14
MAE	0.35	0.33	0.40	0.40	0.37	0.36
BIAS	0.02	0.03	0.03	0.03	0.04	0.03
RSQ	0.91	0.92	0.88	0.88	0.91	0.90

- In comparison the self built linear model and the forward algorithm have the best results amongst all implemented models
- The random forest and the gradient boosting method show a far better Mean Squared Error on the training data. This might indicate, that these models are overfitted to the data.

Performance measure example: MSE Code

```
70 predictions.lm = predict(lm.fit, newdata = test)
71 df.lm.fit      = data.frame(cbind(test$logSalePrice, predictions.lm))
72
73 lm.plot = ggplot(df.lm.fit, aes(test$logSalePrice, predictions.lm)) +
  geom_point() + geom_segment(x = -4, y = -4, xend = 4, yend = 4,
  color = "red", size = 1.3) +
74   stat_smooth(method = "lm", se = FALSE) +
75   labs(title = "Plot of real logSalePrice against predicted values",
76        x = "logSalePrice", y = "lm.fit predictions") + theme(axis.title =
  element_text(size = 16), plot.title = element_text(size = 16,
  face = "bold")) +
77   annotate("text", label = paste("MSE:", comparison.result["MSE", "lm
  "], sep = " "), x = -3, y = 3) +
78   annotate("text", label = paste("MAE:", comparison.result["MAE", "lm
  "], sep = " "), x = -3, y = 2.5)
```

Real Vs. Estimated Plots



- If the house prices would be predicted perfectly, all points would lie on the red line
- The blue line comes from an OLS regression: The better the predictions the more blue and red line should be the same