

# MACHINE LEARNING: REGRESSION IN R

Jan-Philipp Kolb

13 Januar, 2020

# WHY A PART ON SIMPLE REGRESSION

- ▶ OLS can be seen as a simple machine learning technique
- ▶ Some other machine learning concepts are based on regression (e.g. regularization).
- ▶ We would like to remind you how simple regression works in R.
- ▶ We also want to show the constraints
- ▶ In a next step we will learn, how to coop with these constraints

# VARIABLES OF THE MTCARS DATASET

Help for the mtcars dataset:

```
?mtcars
```

- ▶ mpg - Miles/(US) gallon
- ▶ cyl - Number of cylinders
- ▶ disp - Displacement (cu.in.)
- ▶ hp - Gross horsepower
- ▶ drat - Rear axle ratio
- ▶ wt - Weight (1000 lbs)
- ▶ qsec - 1/4 mile time
- ▶ vs - Engine (0 = V-shaped, 1 = straight)
- ▶ am - Transmission (0 = automatic, 1 = manual)
- ▶ gear - Number of forward gears
- ▶ carb - Number of carburetors

## DATASET MTCARS

	mpg	cyl	displacement	hp	drat	wt	qsec	vs	am
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	0
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0
Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	1
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0
Duster 360	14.3	8	360.0	245	3.21	3.570	15.84	0	1
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0
Merc 450SE	16.4	8	275.8	180	3.07	4.070	17.40	0	1
Merc 450SL	17.3	8	275.8	180	3.07	3.730	17.60	0	1
Merc 450SLC	15.2	8	275.8	180	3.07	3.780	18.00	0	1
Cadillac Fleetwood	10.4	8	472.0	205	2.93	5.250	17.98	0	1
Lincoln Continental	10.4	8	460.0	215	3.00	5.424	17.82	0	1

# A SIMPLE REGRESSION MODEL

DEPENDENT VARIABLE - MILES PER GALLON (MPG)

INDEPENDENT VARIABLE - WEIGHT (WT)

```
m1 <- lm(mpg ~ wt,data=mtcars)
m1
##
## Call:
## lm(formula = mpg ~ wt, data = mtcars)
##
## Coefficients:
## (Intercept)          wt
##      37.285      -5.344
```

## GET THE MODEL SUMMARY

```
summary(m1)

##
## Call:
## lm(formula = mpg ~ wt, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.5432 -2.3647 -0.1252  1.4096  6.8727
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  37.2851     1.8776   19.858 < 2e-16 ***
## wt          -5.3445     0.5591   -9.559 1.29e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.046 on 30 degrees of freedom
```

# THE MODEL FORMULA

## MODEL WITHOUT INTERCEPT

```
m2 <- lm(mpg ~ - 1 + wt,data=mtcars)
summary(m2)$coefficients
##      Estimate Std. Error  t value    Pr(>|t|)
## wt  5.291624   0.5931801  8.920771 4.55314e-10
```

## ADDING FURTHER VARIABLES

```
m3 <- lm(mpg ~ wt + cyl,data=mtcars)
summary(m3)$coefficients
##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) 39.686261   1.7149840 23.140893 3.043182e-20
## wt          -3.190972   0.7569065 -4.215808 2.220200e-04
## cyl          -1.507795   0.4146883 -3.635972 1.064282e-03
```

# THE COMMAND `as.formula`

```
?as.formula
```

```
class(fo <- mpg ~ wt + cyl)
```

```
## [1] "formula"
```

```
# The formula object can be used in the regression:
```

```
m3 <- lm(fo,data=mtcars)
```



# FURTHER POSSIBILITIES TO SPECIFY THE FORMULA

## TAKE ALL AVAILABLE PREDICTORS

```
m3_a<-lm(mpg~.,data=mtcars)
```

## INTERACTION EFFECT

```
# effect of cyl and interaction effect:
```

```
m3a<-lm(mpg~wt*cyl,data=mtcars)
```

```
# only interaction effect:
```

```
m3b<-lm(mpg~wt:cyl,data=mtcars)
```

## TAKE THE LOGARITHM

```
m3d<-lm(mpg~log(wt),data=mtcars)
```

# THE COMMAND SETDIFF

- We can use the command to create a dataset with only the features, without the dependent variable

```
names(mtcars)

## [1] "mpg" "cyl" "disp" "hp" "drat" "wt" "qsec" "vs"
## [11] "carb"

features <- setdiff(names(mtcars), "mpg")
features

## [1] "cyl" "disp" "hp" "drat" "wt" "qsec" "vs" "am"

featdat <- mtcars[,features]
```

## THE COMMAND `MODEL.MATRIX`

- ▶ With `model.matrix` the qualitative variables are automatically dummy encoded

```
?model.matrix
```

```
model.matrix(m3d)
```

##	(Intercept)	log(wt)
## Mazda RX4	1	0.9631743
## Mazda RX4 Wag	1	1.0560527
## Datsun 710	1	0.8415672
## Hornet 4 Drive	1	1.1678274
## Hornet Sportabout	1	1.2354715
## Valiant	1	1.2412686
## Duster 360	1	1.2725656
## Merc 240D	1	1.1600209
## Merc 230	1	1.1474025
## Merc 280	1	1.2354715
## Merc 280C	1	1.2354715

## MODEL MATRIX (II)

- ▶ We can also create a model matrix directly from the formula and data arguments
- ▶ See `Matrix::sparse.model.matrix` for increased efficiency on large dimension data.

```
ff <- mpg ~ log(wt):cyl  
m <- model.frame(ff, mtcars)  
  
(mat <- model.matrix(ff, m))
```

##	(Intercept)	log(wt):cyl
## Mazda RX4	1	5.779046
## Mazda RX4 Wag	1	6.336316
## Datsun 710	1	3.366269
## Hornet 4 Drive	1	7.006964
## Hornet Sportabout	1	9.883772
## Valiant	1	7.447612
## Duster 360	1	10.180525
## Merc 240D	1	4.640084

# A MODEL WITH INTERACTION EFFECT

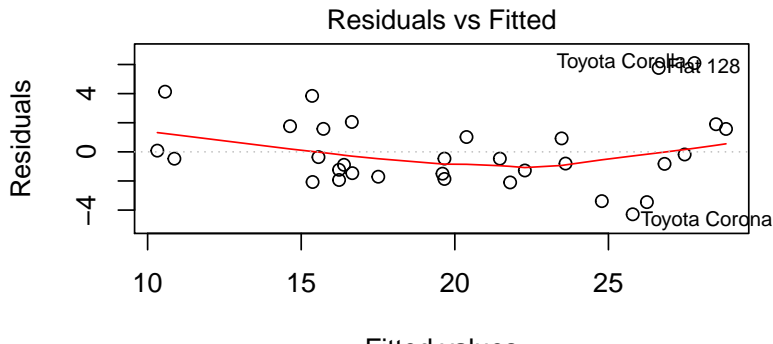
```
# disp - Displacement (cu.in.)  
m3d<-lm(mpg~wt*disp,data=mtcars)  
m3dsum <- summary(m3d)  
m3dsum$coefficients
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	44.08199770	3.123062627	14.114990	2.955567e-14
## wt	-6.49567966	1.313382622	-4.945763	3.216705e-05
## disp	-0.05635816	0.013238696	-4.257078	2.101721e-04
## wt:disp	0.01170542	0.003255102	3.596022	1.226988e-03

## RESIDUAL PLOT - MODEL ASSUMPTIONS VIOLATED?

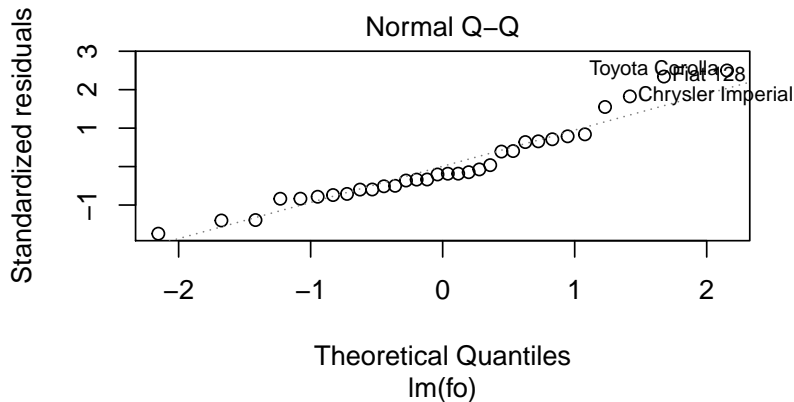
- ▶ We have model assumptions violated if points deviate with a pattern from the line

```
plot(m3,1)
```



# RESIDUAL PLOT

`plot(m3,2)`



## ANOTHER EXAMPLE FOR OBJECT ORIENTATION

- ▶ m3 is now a special regression object
- ▶ Various functions can be applied to this object

```
predict(m3) # Prediction  
resid(m3) # Residuals
```

##	Mazda RX4	Mazda RX4 Wag	Datsun 710	Horn
##	22.27914	21.46545	26.25203	
##	Hornet Sportabout	Valiant		
##	16.64696	19.59873		
##	Mazda RX4	Mazda RX4 Wag	Datsun 710	Horn
##	-1.2791447	-0.4654468	-3.4520262	
##	Hornet Sportabout	Valiant		
##	2.0530424	-1.4987281		



# MAKE MODEL PREDICTION

```
pre <- predict(m1)
head(mtcars$mpg)
```

```
## [1] 21.0 21.0 22.8 21.4 18.7 18.1
```

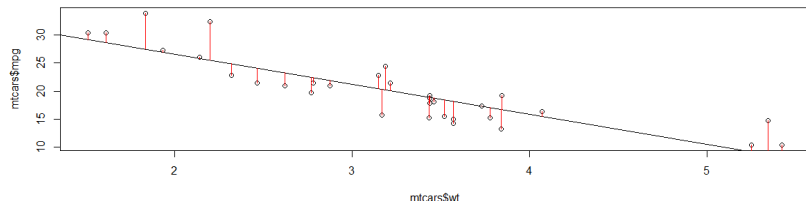
```
head(pre)
```

```
##           Mazda RX4      Mazda RX4 Wag      Datsun 710      Horn
##           23.28261           21.91977           24.88595
## Hornet Sportabout           Valiant
##           18.90014           18.79325
```

# REGRESSION DIAGNOSTIC WITH BASE-R

## VISUALIZING RESIDUALS

```
plot(mtcars$wt,mtcars$mpg)  
abline(m1)  
segments(mtcars$wt, mtcars$mpg, mtcars$wt, pre, col="red")
```



# THE MEAN SQUARED ERROR (MSE)

- ▶ The **MSE** measures the average of the squares of the errors
- ▶ **The lower the better**

```
(mse5 <- mean((mtcars$mpg - pre)^2)) # model 5  
## [1] 8.697561
```

```
(mse3 <- mean((mtcars$mpg - predict(m3))^2))  
## [1] 5.974124
```

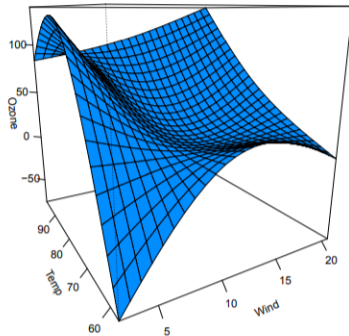
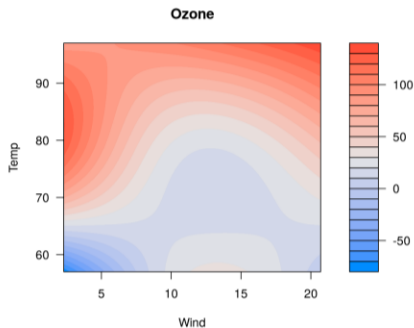
## PACKAGE METRICS TO COMPUTE MSE

```
library(Metrics)  
mse(mtcars$mpg, predict(m3))  
## [1] 5.974124
```

# THE VISREG-PACKAGE

```
install.packages("visreg")
```

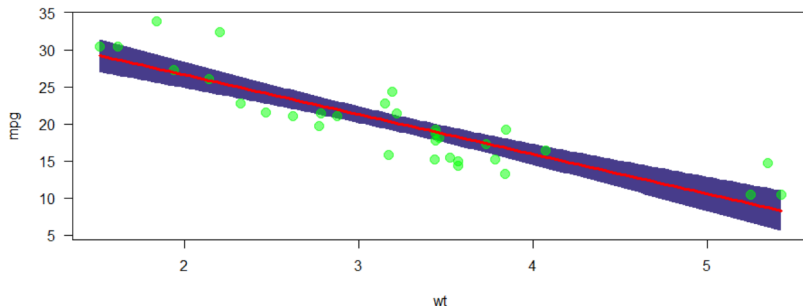
```
library(visreg)
```



# THE VISREG-PACKAGE

- ▶ The default-argument for type is conditional.
- ▶ Scatterplot of mpg and wt plus regression line and confidence bands

```
visreg(m1, "wt", type = "conditional")
```



# REGRESSION WITH FACTORS

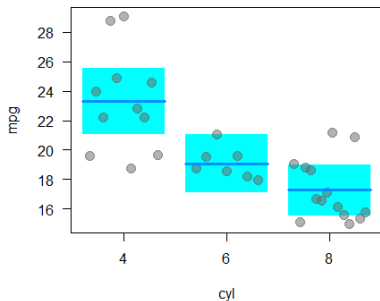
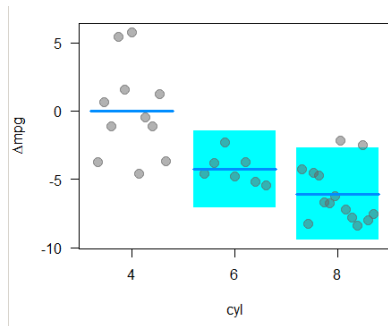
- The effects of factors can also be visualized with visreg:

```
mtcars$cyl <- as.factor(mtcars$cyl)
m4 <- lm(mpg ~ cyl + wt, data = mtcars)
# summary(m4)
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	33.990794	1.8877934	18.005569	6.257246e-17
## cyl6	-4.255582	1.3860728	-3.070244	4.717834e-03
## cyl8	-6.070860	1.6522878	-3.674214	9.991893e-04
## wt	-3.205613	0.7538957	-4.252065	2.130435e-04

# EFFECTS OF FACTORS

```
par(mfrow=c(1,2))  
visreg(m4, "cyl", type = "contrast")  
visreg(m4, "cyl", type = "conditional")
```



# THE PACKAGE VISREG - INTERACTIONS

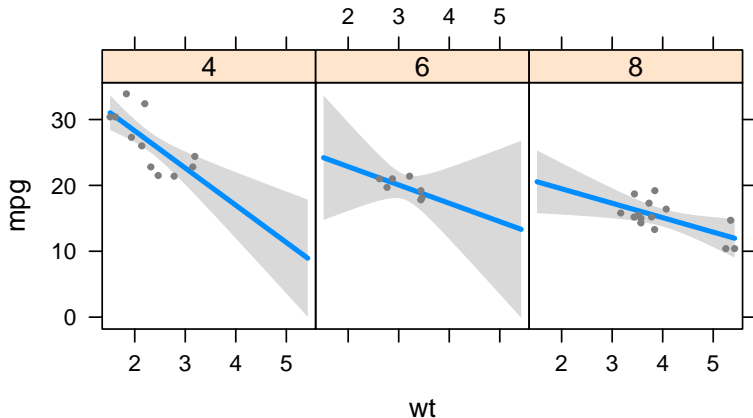
```
m5 <- lm(mpg ~ cyl*wt, data = mtcars)
# summary(m5)
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	39.571196	3.193940	12.3894599	2.058359e-12
## cyl6	-11.162351	9.355346	-1.1931522	2.435843e-01
## cyl8	-15.703167	4.839464	-3.2448150	3.223216e-03
## wt	-5.647025	1.359498	-4.1537586	3.127578e-04
## cyl6:wt	2.866919	3.117330	0.9196716	3.661987e-01
## cyl8:wt	3.454587	1.627261	2.1229458	4.344037e-02



## CONTROL OF THE GRAPHIC OUTPUT WITH LAYOUT.

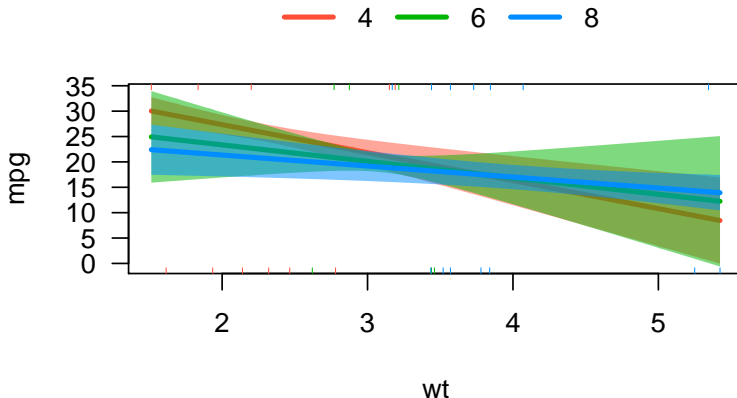
```
visreg(m5, "wt", by = "cyl", layout=c(3,1))
```



# THE PACKAGE VISREG - INTERACTIONS OVERLAY

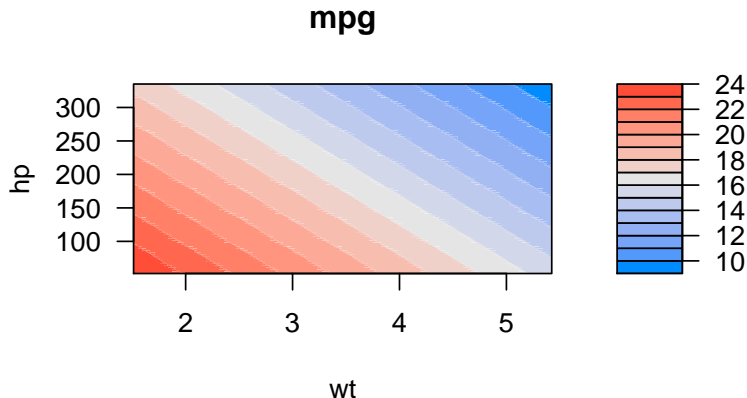
```
m6 <- lm(mpg ~ hp + wt * cyl, data = mtcars)
```

```
visreg(m6, "wt", by="cyl", overlay=TRUE, partial=FALSE)
```



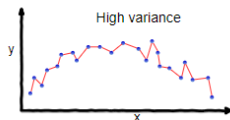
# THE PACKAGE VISREG - VISREG2D

```
visreg2d(m6, "wt", "hp", plot.type = "image")
```

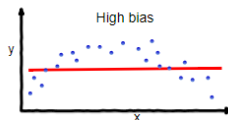


# THE BIAS-VARIANCE TRADEOFF (I)

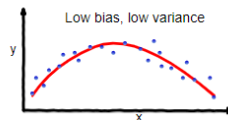
- ▶ The bias–variance tradeoff is the property of a set of predictive models whereby models with a lower bias in parameter estimation have a higher variance of the parameter estimates across samples, and vice versa.



**overfitting**



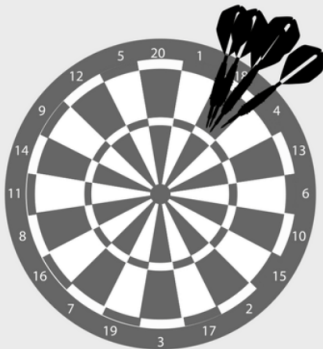
**underfitting**



**Good balance**

## THE BIAS-VARIANCE TRADEOFF (II)

**High Bias**  
Low Variance



**High bias**, low variance  
algorithms train models that

**High Variance**  
Low Bias



**High variance**, low bias  
algorithms train models that

## EXERCISE: REGRESSION AMES HOUSING DATA

- 1) Install the package `AmesHousing` and create a **processed version** of the Ames housing data with (at least) the variables `Sale_Price`, `Gr_Liv_Area` and `TotRms_AbvGrd`
- 2) Create a regression model with `Sale_Price` as dependent and `Gr_Liv_Area` and `TotRms_AbvGrd` as independent variables. Then create separated models for the two independent variables. Compare the results. What do you think?

# THE AMES IOWA HOUSING DATA

```
ames_data <- AmesHousing::make_ames()
```

## SOME VARIABLES

- ▶ Gr\_Liv\_Area: Above grade (ground) living area square feet
- ▶ TotRms\_AbvGrd: Total rooms above grade (does not include bathrooms)
- ▶ MS\_SubClass: Identifies the type of dwelling involved in the sale.
- ▶ MS\_Zoning: Identifies the general zoning classification of the sale.
- ▶ Lot\_Frontage: Linear feet of street connected to property
- ▶ Lot\_Area: Lot size in square feet
- ▶ Street: Type of road access to property
- ▶ Alley: Type of alley access to property
- ▶ Lot\_Shape: General shape of property
- ▶ Land\_Contour: Flatness of the propert

# MULTICOLLINEARITY

- ▶ As  $p$  increases we are more likely to capture multiple features that have some multicollinearity.
- ▶ When multicollinearity exists, we often see high variability in our coefficient terms.
- ▶ E.g. we have a correlation of 0.801 between Gr\_Liv\_Area and TotRms\_AbvGrd
- ▶ Both variables are strongly correlated to the response variable (Sale\_Price).

```
ames_data <- AmesHousing::make_ames()  
cor(ames_data[,c("Sale_Price", "Gr_Liv_Area", "TotRms_AbvGrd")])
```

```
##           Sale_Price Gr_Liv_Area TotRms_AbvGrd  
## Sale_Price      1.0000000    0.7067799    0.4954744  
## Gr_Liv_Area     0.7067799    1.0000000    0.8077721  
## TotRms_AbvGrd  0.4954744    0.8077721    1.0000000
```



# EFFECTS OF MULTICOLLINEARITY

```
lm(Sale_Price ~ Gr_Liv_Area + TotRms_AbvGrd, data = ames_data)

##
## Call:
## lm(formula = Sale_Price ~ Gr_Liv_Area + TotRms_AbvGrd, data =
##
## Coefficients:
##      (Intercept)      Gr_Liv_Area  TotRms_AbvGrd
##          42767.6           139.4         -11025.9
```

- ▶ When we fit a model with both these variables we get a positive coefficient for Gr\_Liv\_Area but a negative coefficient for TotRms\_AbvGrd, suggesting one has a positive impact to Sale\_Price and the other a negative impact.

## SEPERATED MODELS

- ▶ If we refit the model with each variable independently, they both show a positive impact.
- ▶ The Gr\_Liv\_Area effect is now smaller and the TotRms\_AbvGrd is positive with a much larger magnitude.

```
lm(Sale_Price ~ Gr_Liv_Area, data = ames_data)$coefficients
```

```
## (Intercept) Gr_Liv_Area
```

```
##    13289.634    111.694
```

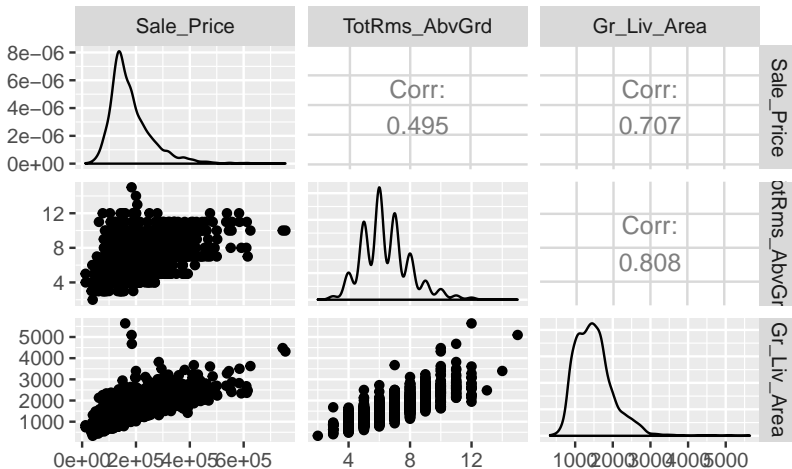
```
lm(Sale_Price ~ TotRms_AbvGrd, data = ames_data)$coefficients
```

```
## (Intercept) TotRms_AbvGrd
```

```
##    18665.40    25163.83
```

- ▶ This is a common result when collinearity exists.
- ▶ Coefficients for correlated features become over-inflated and can fluctuate significantly.

```
library(GGally)
ggpairs(ames_data[,c("Sale_Price", "TotRms_AbvGrd", "Gr_Liv_Area")])
```

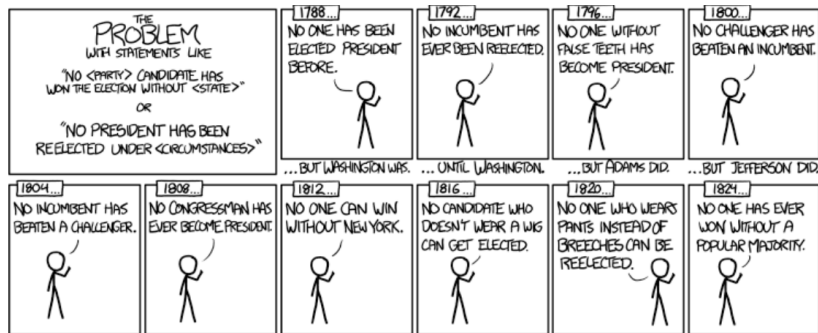


# CONSEQUENCES

- ▶ One consequence of these large fluctuations in the coefficient terms is **overfitting**, which means we have high variance in the bias-variance tradeoff space.
- ▶ We can use tools such as **variance inflation factors** (Myers, 1994) to identify and remove those strongly correlated variables, but it is not always clear which variable(s) to remove.
- ▶ Nor do we always wish to remove variables as this may be removing signal in our data.

# THE PROBLEM - OVERFITTING

- Our model doesn't generalize well from our training data to unseen data.



# WHAT CAN BE DONE AGAINST OVERFITTING

- ▶ **Cross Validation**
- ▶ Train with more data
- ▶ Remove features
- ▶ **Regularization** - e.g. ridge and lasso regression
- ▶ Ensembling - e.g. bagging and boosting

# CROSS VALIDATION

- ▶ Cross-validation is a powerful preventative measure against overfitting.
- ▶ Use your initial training data to generate multiple mini train-test splits. Use these splits to tune your model.

## NECESSARY PACKAGES

```
library(tidyverse)  
library(caret)
```

# CROSS VALIDATION IN R

## SPLIT DATA INTO TRAINING AND TESTING DATASET

```
training.samples <- ames_data$Sale_Price %>%  
createDataPartition(p = 0.8, list = FALSE)  
train.data <- ames_data[training.samples, ]  
test.data <- ames_data[-training.samples, ]
```

## BUILD THE MODEL AND MAKE PREDICTIONS

```
model <- lm(Sale_Price ~ Gr_Liv_Area + TotRms_AbvGrd,  
            data = train.data)
```

# Make predictions and compute the R2, RMSE and MAE

```
(predictions <- model %>% predict(test.data))
```

##	1	2	3	4	5	6
##	161758.45	204342.39	163035.96	165545.85	207800.76	117238.90
##	9	10	11	12	13	14
##	281748.37	299059.34	313066.89	263101.73	223408.37	167487.96
##	17	18	19	20	21	22
##	209923.45	128826.86	111606.19	131478.70	164932.92	207136.17
##	25	26	27	28	29	30



# MODEL WITH CROSS VALIDATION

## ► Loocv: **leave one out cross validation**

```
train.control <- caret::trainControl(method = "LOOCV")  
  
# Train the model  
model2 <- train(Sale_Price ~ Gr_Liv_Area + TotRms_AbvGrd,  
                data = train.data, method = "lm",  
                trControl = train.control)  
model2 %>% predict(test.data)
```

# LINKS - LINEAR REGRESSION

- ▶ Regression - **r-bloggers**
- ▶ The complete book of **Faraway**- very intuitive
- ▶ Good introduction on **Quick-R**
- ▶ **Multiple regression**
- ▶ **15 Types of Regression you should know**
- ▶ **ggeffects - Create Tidy Data Frames of Marginal Effects for 'ggplot' from Model Outputs**
- ▶ **Machine learning iteration**

## NICE TABLE OUTPUT WITH STARGAZER

```
library(stargazer)
stargazer(m3, type="html")
```

EXAMPLE HTML OUTPUT:

	<i>Dependent variable:</i>
	mpg
wt	-3.125*** (0.911)
cyl	-1.510*** (0.422)
am	0.176 (1.304)
Constant	39.418*** (2.641)
Observations	32
R <sup>2</sup>	0.830

# SHINY APP - DIAGNOSTICS FOR LINEAR REGRESSION

- ▶ Shiny App - **Simple Linear Regression**
- ▶ Shiny App - **Multicollinearity in multiple regression**

## Diagnostics for simple linear regression

Select a trend:

- ☐ Linear up
- ☐ Linear down
- ☐ Curved up
- ☐ Curved down
- ☒ Fan-shaped

☒ Show residuals

This applet uses ordinary least squares (OLS) to fit a regression line to the data with the selected trend. The applet is designed to help you practice evaluating whether or not the linear model is an appropriate fit to the data. The three diagnostic plots on the lower half of the page are provided to help you identify undesirable patterns in the residuals that may arise from non-linear trends in the data.

[Rate this app!](#)

[View code](#)

[Check out other apps](#)

[Want to learn more for free?](#)

