

00-getdata

2022-04-24

Get Data

```
train_data <- read.csv("train.csv")
```

Check Data

```
str(train_data)
```

```
## 'data.frame':    1460 obs. of  81 variables:
## $ Id             : int  1 2 3 4 5 6 7 8 9 10 ...
## $ MSSubClass     : int  60 20 60 70 60 50 20 60 50 190 ...
## $ MSZoning       : chr  "RL" "RL" "RL" "RL" ...
## $ LotFrontage    : int  65 80 68 60 84 85 75 NA 51 50 ...
## $ LotArea        : int  8450 9600 11250 9550 14260 14115 10084 10382 6120 7420 ...
## $ Street         : chr  "Pave" "Pave" "Pave" "Pave" ...
## $ Alley          : chr  NA NA NA NA ...
## $ LotShape       : chr  "Reg" "Reg" "IR1" "IR1" ...
## $ LandContour    : chr  "Lvl" "Lvl" "Lvl" "Lvl" ...
## $ Utilities      : chr  "AllPub" "AllPub" "AllPub" "AllPub" ...
## $ LotConfig      : chr  "Inside" "FR2" "Inside" "Corner" ...
## $ LandSlope      : chr  "Gtl" "Gtl" "Gtl" "Gtl" ...
## $ Neighborhood   : chr  "CollgCr" "Veenker" "CollgCr" "Crawfor" ...
## $ Condition1     : chr  "Norm" "Feedr" "Norm" "Norm" ...
## $ Condition2     : chr  "Norm" "Norm" "Norm" "Norm" ...
## $ BldgType       : chr  "1Fam" "1Fam" "1Fam" "1Fam" ...
## $ HouseStyle     : chr  "2Story" "1Story" "2Story" "2Story" ...
## $ OverallQual    : int  7 6 7 7 8 5 8 7 7 5 ...
## $ OverallCond    : int  5 8 5 5 5 5 5 6 5 6 ...
## $ YearBuilt      : int  2003 1976 2001 1915 2000 1993 2004 1973 1931 1939 ...
## $ YearRemodAdd   : int  2003 1976 2002 1970 2000 1995 2005 1973 1950 1950 ...
## $ RoofStyle      : chr  "Gable" "Gable" "Gable" "Gable" ...
## $ RoofMatl       : chr  "CompShg" "CompShg" "CompShg" "CompShg" ...
## $ Exterior1st    : chr  "VinylSd" "MetalSd" "VinylSd" "Wd Sdng" ...
## $ Exterior2nd    : chr  "VinylSd" "MetalSd" "VinylSd" "Wd Shng" ...
## $ MasVnrType     : chr  "BrkFace" "None" "BrkFace" "None" ...
## $ MasVnrArea     : int  196 0 162 0 350 0 186 240 0 0 ...
## $ ExterQual      : chr  "Gd" "TA" "Gd" "TA" ...
## $ ExterCond      : chr  "TA" "TA" "TA" "TA" ...
## $ Foundation     : chr  "PConc" "CBlock" "PConc" "BrkTil" ...
## $ BsmtQual       : chr  "Gd" "Gd" "Gd" "TA" ...
```

```

## $ BsmtCond      : chr  "TA" "TA" "TA" "Gd" ...
## $ BsmtExposure  : chr  "No" "Gd" "Mn" "No" ...
## $ BsmtFinType1  : chr  "GLQ" "ALQ" "GLQ" "ALQ" ...
## $ BsmtFinSF1    : int   706 978 486 216 655 732 1369 859 0 851 ...
## $ BsmtFinType2  : chr  "Unf" "Unf" "Unf" "Unf" ...
## $ BsmtFinSF2    : int   0 0 0 0 0 0 0 32 0 0 ...
## $ BsmtUnfSF     : int   150 284 434 540 490 64 317 216 952 140 ...
## $ TotalBsmtSF   : int   856 1262 920 756 1145 796 1686 1107 952 991 ...
## $ Heating       : chr  "GasA" "GasA" "GasA" "GasA" ...
## $ HeatingQC     : chr  "Ex" "Ex" "Ex" "Gd" ...
## $ CentralAir    : chr  "Y" "Y" "Y" "Y" ...
## $ Electrical    : chr  "SBrkr" "SBrkr" "SBrkr" "SBrkr" ...
## $ X1stFlrSF     : int   856 1262 920 961 1145 796 1694 1107 1022 1077 ...
## $ X2ndFlrSF     : int   854 0 866 756 1053 566 0 983 752 0 ...
## $ LowQualFinSF  : int   0 0 0 0 0 0 0 0 0 0 ...
## $ GrLivArea     : int   1710 1262 1786 1717 2198 1362 1694 2090 1774 1077 ...
## $ BsmtFullBath  : int   1 0 1 1 1 1 1 1 0 1 ...
## $ BsmtHalfBath  : int   0 1 0 0 0 0 0 0 0 0 ...
## $ FullBath      : int   2 2 2 1 2 1 2 2 2 1 ...
## $ HalfBath      : int   1 0 1 0 1 1 0 1 0 0 ...
## $ BedroomAbvGr : int   3 3 3 3 4 1 3 3 2 2 ...
## $ KitchenAbvGr  : int   1 1 1 1 1 1 1 1 2 2 ...
## $ KitchenQual   : chr  "Gd" "TA" "Gd" "Gd" ...
## $ TotRmsAbvGrd : int   8 6 6 7 9 5 7 7 8 5 ...
## $ Functional    : chr  "Typ" "Typ" "Typ" "Typ" ...
## $ Fireplaces    : int   0 1 1 1 1 0 1 2 2 2 ...
## $ FireplaceQu   : chr  NA "TA" "TA" "Gd" ...
## $ GarageType    : chr  "Attchd" "Attchd" "Attchd" "Detchd" ...
## $ GarageYrBlt   : int   2003 1976 2001 1998 2000 1993 2004 1973 1931 1939 ...
## $ GarageFinish  : chr  "RFn" "RFn" "RFn" "Unf" ...
## $ GarageCars    : int   2 2 2 3 3 2 2 2 2 1 ...
## $ GarageArea    : int   548 460 608 642 836 480 636 484 468 205 ...
## $ GarageQual    : chr  "TA" "TA" "TA" "TA" ...
## $ GarageCond    : chr  "TA" "TA" "TA" "TA" ...
## $ PavedDrive    : chr  "Y" "Y" "Y" "Y" ...
## $ WoodDeckSF    : int   0 298 0 0 192 40 255 235 90 0 ...
## $ OpenPorchSF   : int   61 0 42 35 84 30 57 204 0 4 ...
## $ EnclosedPorch : int   0 0 0 272 0 0 0 228 205 0 ...
## $ X3SsnPorch    : int   0 0 0 0 0 320 0 0 0 0 ...
## $ ScreenPorch   : int   0 0 0 0 0 0 0 0 0 0 ...
## $ PoolArea      : int   0 0 0 0 0 0 0 0 0 0 ...
## $ PoolQC        : chr  NA NA NA NA ...
## $ Fence         : chr  NA NA NA NA ...
## $ MiscFeature    : chr  NA NA NA NA ...
## $ MiscVal       : int   0 0 0 0 0 700 0 350 0 0 ...
## $ MoSold        : int   2 5 9 2 12 10 8 11 4 1 ...
## $ YrSold        : int   2008 2007 2008 2006 2008 2009 2007 2009 2008 2008 ...
## $ SaleType      : chr  "WD" "WD" "WD" "WD" ...
## $ SaleCondition : chr  "Normal" "Normal" "Normal" "Abnorml" ...
## $ SalePrice     : int   208500 181500 223500 140000 250000 143000 307000 200000 129900 118000 ...

```

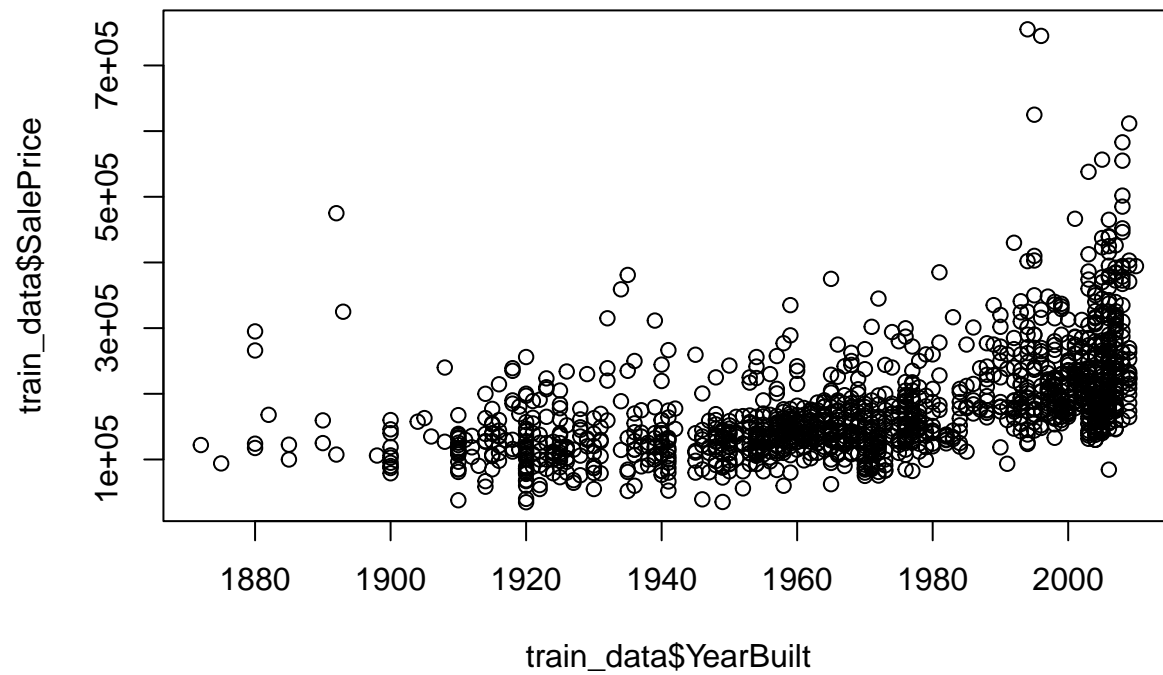
```
dim(train_data)
```

```
## [1] 1460 81
```

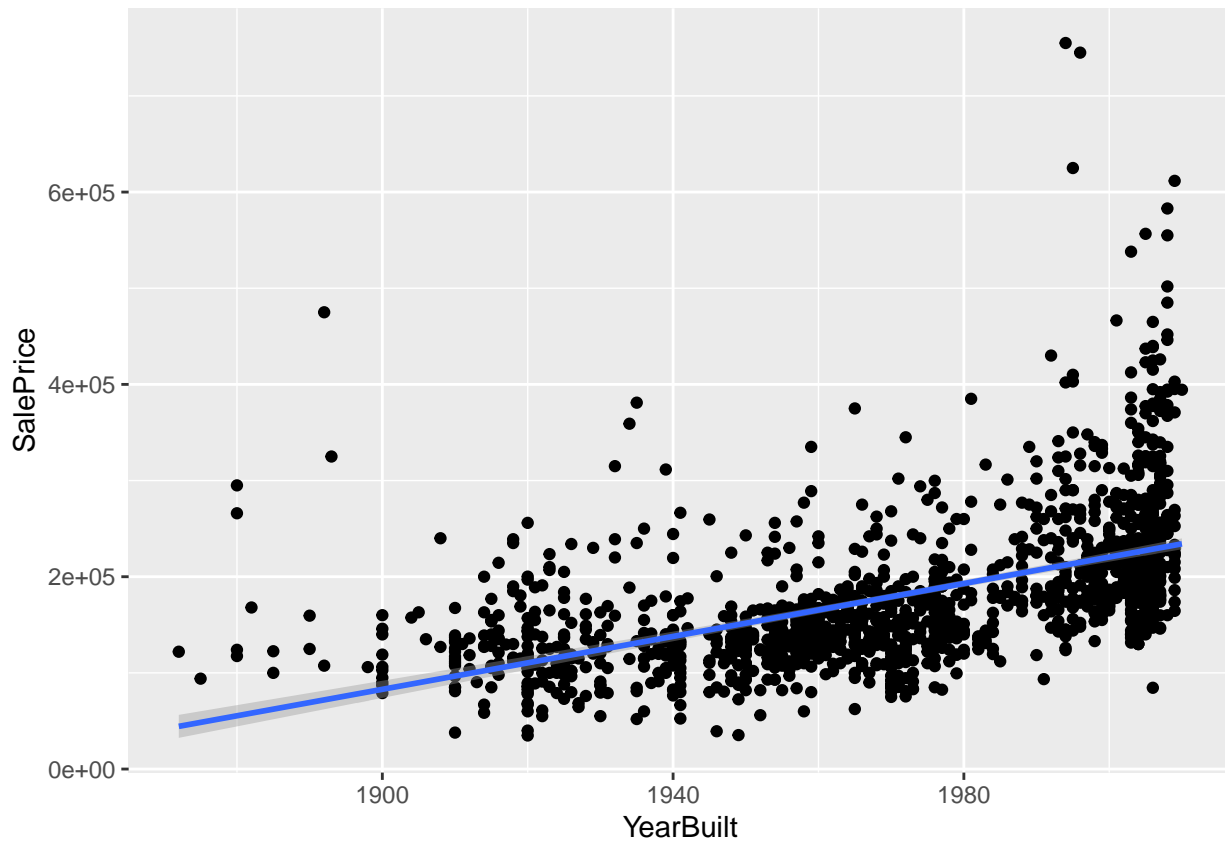
Plot Data (skip)

```
plot(train_data$YearBuilt, train_data$SalePrice)
# plot(x,y)

library(ggplot2)
```



```
ggplot(data=train_data, aes(x=YearBuilt, y=SalePrice))+
  geom_point()+
  geom_smooth(method="lm", formula=y~x)
```



Fit Simple Linear Regression

```
model_0 <- lm( formula = SalePrice ~ YearBuilt, data = train_data)
```

```
# print summary statistics
summary(model_0)
```

```
##
## Call:
## lm(formula = SalePrice ~ YearBuilt, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -144191  -40999  -15464   22685  542814
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.530e+06  1.158e+05  -21.86  <2e-16 ***
## YearBuilt    1.375e+03   5.872e+01   23.42  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 67740 on 1458 degrees of freedom
## Multiple R-squared:  0.2734, Adjusted R-squared:  0.2729
## F-statistic: 548.7 on 1 and 1458 DF,  p-value: < 2.2e-16
```

$$\mu(Y|X) = \beta_0 + \beta_1 X$$

$$\hat{\mu}(Y|X) = \hat{\beta}_0 + \hat{\beta}_1 X$$