



Regression analysis is the most widely used statistical tool for understanding relationships among variables

It provides a conceptually simple method for investigating functional relationships between one or more factors and an outcome of interest

The relationship is expressed in the form of an equation or a model connecting the response or dependent variable and one or more explanatory or predictor variable

# Why?



Straight-up **prediction**:

- How much will I sell my house for?

**Explanation** and understanding:

- What is the impact of economic freedom on growth?

# Predicting house prices



To keep things super simple, let's focus only on size. The value

that we seek to predict is called the  
**dependent (or output)** variable, and we denote this:

- $Y$  = price of house (e.g. thousands of dollars)

The variable that we use to guide prediction is the  
**explanatory (or input)** variable, and this is labeled

- $X$  = size of house (e.g. thousands of square feet)

# Predicting house prices



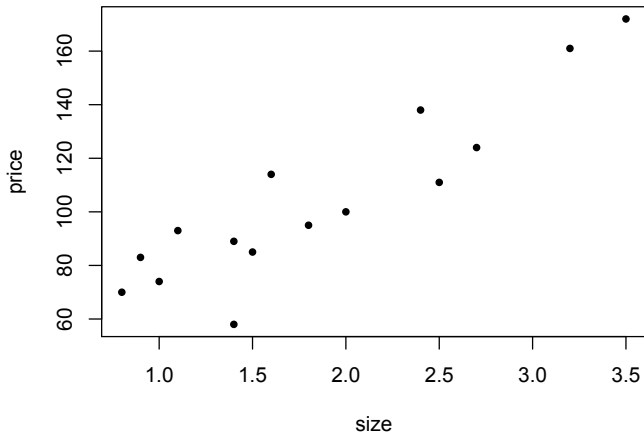
What does this data look like?

Size	Price
0.80	70
0.90	83
1.00	74
1.10	93
1.40	89
1.40	58
1.50	85
1.60	114
1.80	95
2.00	100
2.40	138
2.50	111
2.70	124
3.20	161
3.50	172

# Predicting house prices



It is much more useful to look at a scatterplot



In other words, view the data as points in the  $X \times Y$  plane.

# Regression model



$Y$  = response or outcome variable

$X$  = explanatory or input variables

A linear relationship is written

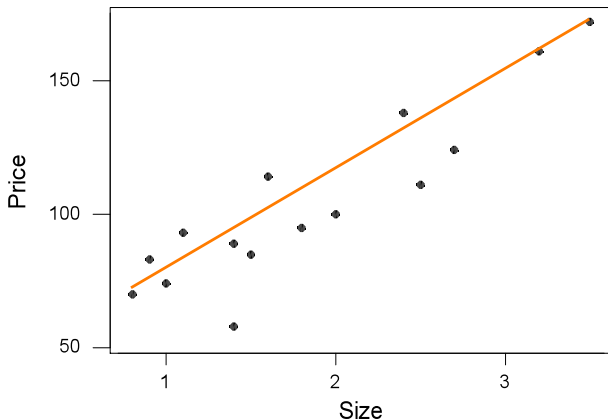
$$Y = b_0 + b_1X + e$$

# Linear prediction



There seems to be a linear relationship between price and size:

As size goes up, price goes up.





Recall that the equation of a line is:

$$Y = b_0 + b_1X$$

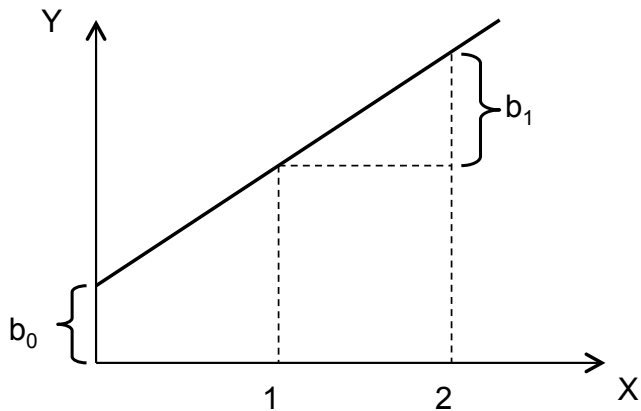
Where  $b_0$  is the **intercept** and  $b_1$  is the **slope**.

→ The **intercept** value is in units of  $Y$  (\$1,000)

→ The **slope** is in units of  $Y$  *per* units of  $X$  (\$1,000/1,000 sq ft)



# Linear prediction



$$Y = b_0 + b_1 X$$



## Q: How to find the “best line”?

We desire a strategy for estimating the slope and intercept parameters in the model  $\hat{Y} = b_0 + b_1X$

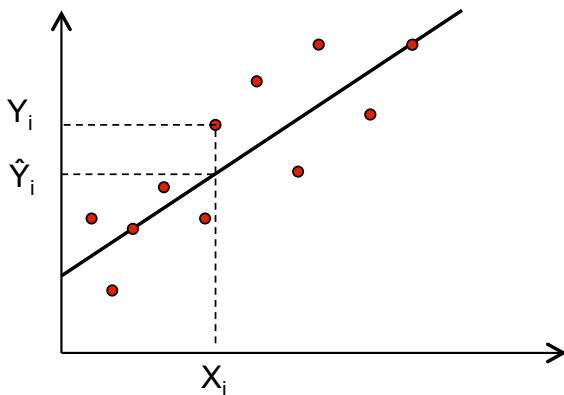
A reasonable way to fit a line is to minimize the amount by which the **fitted value** differs from the actual value.

This amount is called the **residual**.

# Linear prediction



What is the “fitted value”?

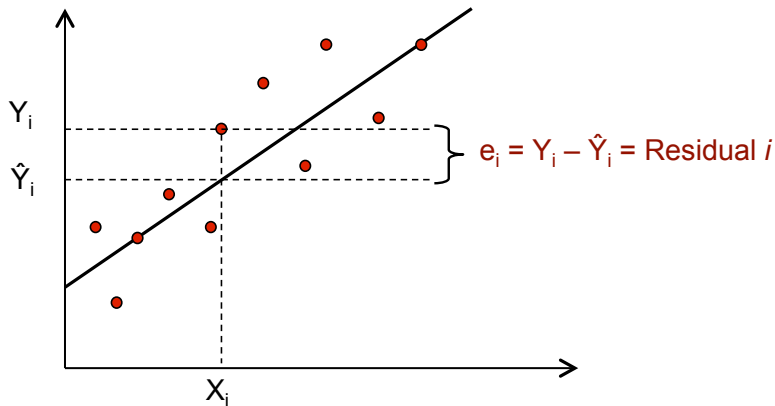


The dots are the observed values and the line represents our fitted values given by  $\hat{Y}_i = b_0 + b_1 X_1$ .

# Linear prediction



What is the “residual” for the  $i$ th observation?



We can write  $Y_i = \hat{Y}_i + (Y_i - \hat{Y}_i) = \hat{Y}_i + e_i$ .



Ideally, we want to minimize the size of all residuals:

- If they were all zero we would have a perfect line.
- Trade-off between moving closer to some points and at the same time moving away from other points.



Ideally, we want to minimize the size of all residuals:

- If they were all zero we would have a perfect line.
- Trade-off between moving closer to some points and at the same time moving away from other points.

The line fitting process:

- Give weights to all of the residuals.
- Minimize the “total” of residuals to get best fit.



Ideally, we want to minimize the size of all residuals:

- If they were all zero we would have a perfect line.
- Trade-off between moving closer to some points and at the same time moving away from other points.

The line fitting process:

- Give weights to all of the residuals.
- Minimize the “total” of residuals to get best fit.

Least Squares chooses  $b_0$  and  $b_1$  to minimize  $\sum_{i=1}^N e_i^2$

$$\sum_{i=1}^N e_i^2 = e_1^2 + e_2^2 + \cdots + e_N^2 = (Y_1 - \hat{Y}_1)^2 + (Y_2 - \hat{Y}_2)^2 + \cdots + (Y_N - \hat{Y}_N)^2$$

# Least squares – R output



```
data = read.csv('housedata.csv')
fit = lm(Price~Size,data)
summary(fit)

##
## Call:
## lm(formula = Price ~ Size, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -30.425  -8.618   0.575  10.766  18.498
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   38.885     9.094   4.276 0.000903 ***
## Size          35.386     4.494   7.874 2.66e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.14 on 13 degrees of freedom
## Multiple R-squared:  0.8267, Adjusted R-squared:  0.8133
## F-statistic:    62 on 1 and 13 DF,  p-value: 2.66e-06
```