

# Introduction to Linear Regression

Arnaud Legrand

MOSIG Lecture

- 1 Simple Linear Regression
  - General Introduction
  - Fitting a Line to a Set of Points
- 2 Linear Model
  - Linear Regression
  - Underlying Hypothesis
  - Checking hypothesis
  - Decomposing the Variance
  - Making Predictions
  - Confidence interval
- 3 Conclusion

- 1 Simple Linear Regression
  - General Introduction
  - Fitting a Line to a Set of Points
- 2 Linear Model
  - Linear Regression
  - Underlying Hypothesis
  - Checking hypothesis
  - Decomposing the Variance
  - Making Predictions
  - Confidence interval
- 3 Conclusion

## 1 Simple Linear Regression

- General Introduction
- Fitting a Line to a Set of Points

## 2 Linear Model

- Linear Regression
- Underlying Hypothesis
- Checking hypothesis
- Decomposing the Variance
- Making Predictions
- Confidence interval

## 3 Conclusion

# What is a regression?

Regression analysis is the most widely used statistical tool for **understanding relationships among variables**. Several possible objectives including:

- 1 **Prediction** of future observations. This includes extrapolation since we all like connecting points by lines when we *expect* things to be continuous
- 2 Assessment of the **effect** of, or **relationship** between, explanatory variables on the response
- 3 A **general description** of data structure (generally expressed in the form of **an equation or a model** connecting the response or dependent variable and one or more explanatory or predictor variable)
- 4 Defining what you should "expect" as it allows you to define and detect what **does not behave as expected**

The linear relationship is the most commonly found one

- we will illustrate how it works
- it is very general and is the basis of many more advanced tools (polynomial regression, ANOVA, ...)

## 1 Simple Linear Regression

- General Introduction
- Fitting a Line to a Set of Points

## 2 Linear Model

- Linear Regression
- Underlying Hypothesis
- Checking hypothesis
- Decomposing the Variance
- Making Predictions
- Confidence interval

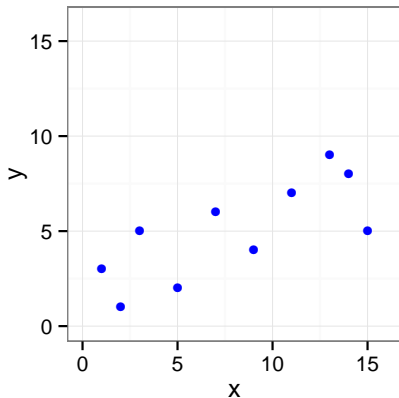
## 3 Conclusion

# Starting With a Simple Data Set

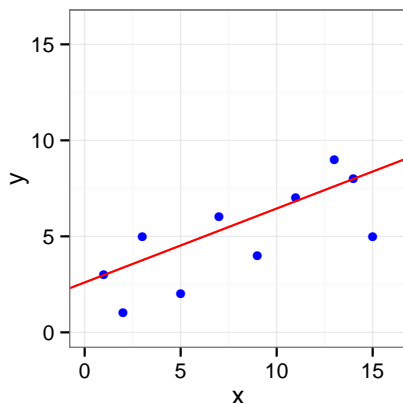
**Descriptive statistics** provides simple summaries about the sample and about the observations that have been made.

How could we summarize the following data set ?

	x	y
1	1.00	3.00
2	2.00	1.00
3	3.00	5.00
4	5.00	2.00
5	7.00	6.00
6	9.00	4.00
7	11.00	7.00
8	13.00	9.00
9	14.00	8.00
10	15.00	5.00



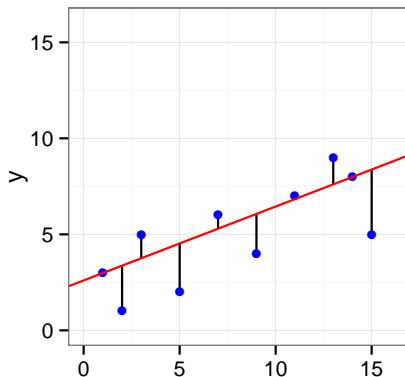
# Eyeball Method



- A straight line drawn through the maximum number of points on a scatter plot balancing about an equal number of points above and below the line
- Some points are rather far from the line. Maybe we should instead try to minimize some kind of *distance to the line*



## Least Squares Line (1): What to minimize?



Intuitively, a large error is *much* more important than a small one. We could try to minimize  $F(\alpha, \beta) = \sum_i \underbrace{(y_i - \alpha - \beta x_i)^2}_{e_i}$ , 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

## Least Squares Line (2): Simple Formula

$$F(\alpha, \beta) = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

$F$  is quadratic in  $\alpha$  and in  $\beta$  so if we simply differentiate  $F$  by  $\alpha$  and by  $\beta$ , we can obtain a closed form for the minimum:

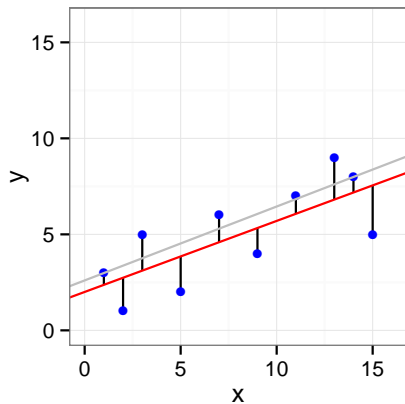
$$\begin{aligned}\hat{\beta} &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \sum_{j=1}^n y_j}{\sum_{i=1}^n (x_i^2) - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \\ &= \frac{\overline{xy} - \bar{x}\bar{y}}{s_x^2 - \bar{x}^2} = \frac{\text{Cov}[x, y]}{\text{Var}[x]} = r_{xy} \frac{s_y}{s_x}\end{aligned}$$

$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$ , where:

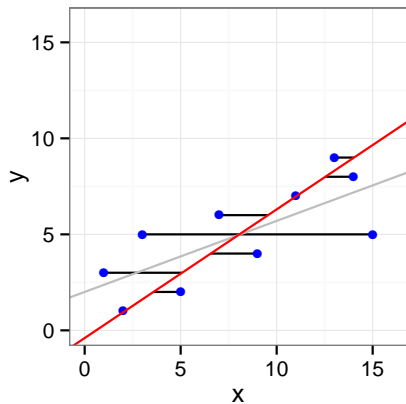
- $\bar{x}$  and  $\bar{y}$  are the sample mean of  $x$  and  $y$
- $r_{xy}$  is the sample correlation coefficient between  $x$  and  $y$
- $s_x$  and  $s_y$  are the sample standard deviation of  $x$  and  $y$

Also has a good geometric interpretation (orthogonal projection)

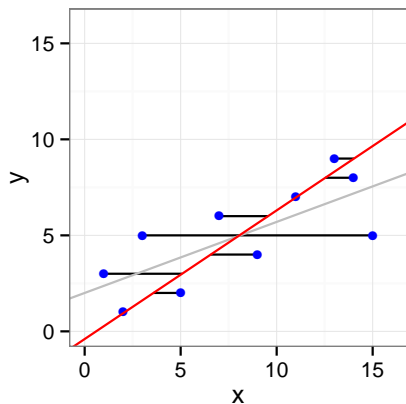
# Least Squares Line (3): $y$ as a function of $x$ or the opposite?



# Least Squares Line (3): $y$ as a function of $x$ or the opposite?

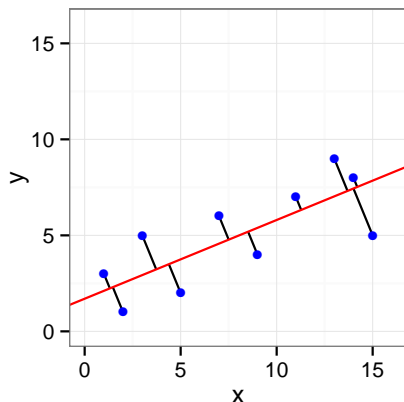


# Least Squares Line (3): $y$ as a function of $x$ or the opposite?



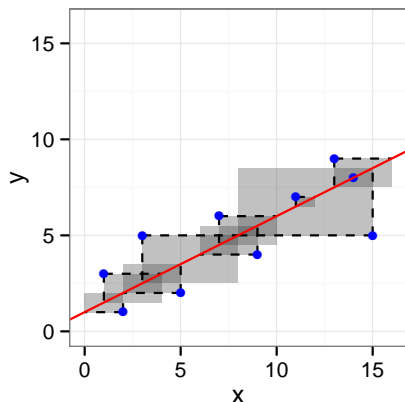
OK, do we have less asymmetrical options?

## Least Distances Line (a.k.a. Deming Regression)



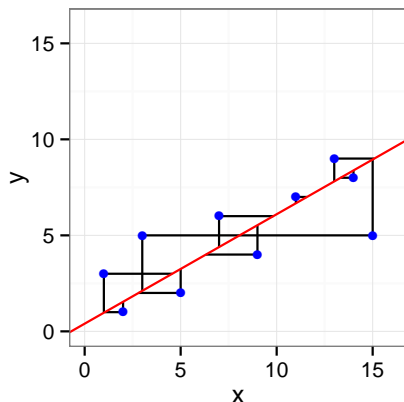
- Note that somehow, this makes sense only if we have a square plot, i.e., if  $x$  and  $y$  have the same units

# Least Rectangles Line



- Minimize  $E(\alpha, \beta) = \sum_{i=1}^n \left| x_i - \frac{y_i - \alpha}{\beta} \right| \cdot |y_i - \alpha - \beta x_i|$
- This leads to the regression line  $y = \frac{s_y}{s_x}(x - \bar{x}) + \bar{y}$ .

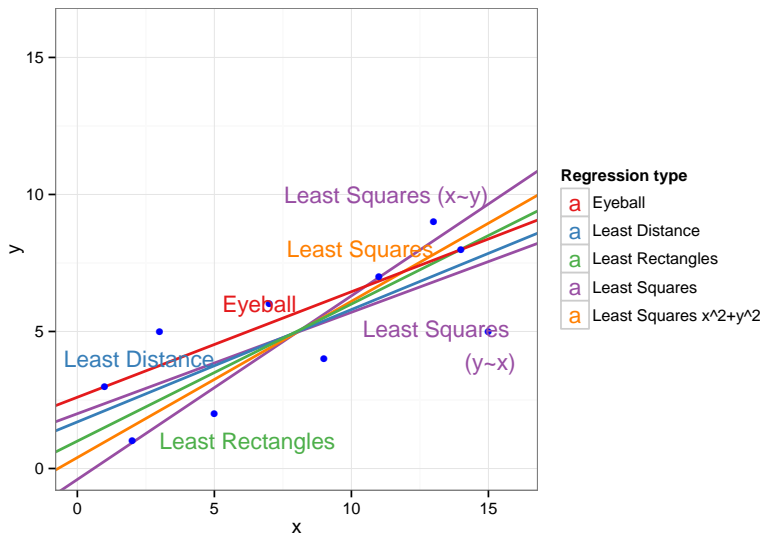
## Least Squares (in Both Directions) Line



- Minimize  $D(\alpha, \beta) = \sum_{i=1}^n \left( x_i - \frac{y_i - \alpha}{\beta} \right)^2 + (y_i - \alpha - \beta x_i)^2$
- Has to be computed analytically



# Which line to choose?



## What does correspond to each line?

- Eyeball: AFAIK nothing
- Least Squares: classical linear regression  $y \sim x$
- Least Squares in both directions: I don't know
- Deming: equivalent to Principal Component Analysis
- Rectangles: may be used when one variable is not "explained" by the other, but are inter-dependent

This is not just a geometric problem. You need a **model** of to decide which one to use

## 1 Simple Linear Regression

- General Introduction
- Fitting a Line to a Set of Points

## 2 Linear Model

- Linear Regression
- Underlying Hypothesis
- Checking hypothesis
- Decomposing the Variance
- Making Predictions
- Confidence interval

## 3 Conclusion

## 1 Simple Linear Regression

- General Introduction
- Fitting a Line to a Set of Points

## 2 Linear Model

- Linear Regression
- Underlying Hypothesis
- Checking hypothesis
- Decomposing the Variance
- Making Predictions
- Confidence interval

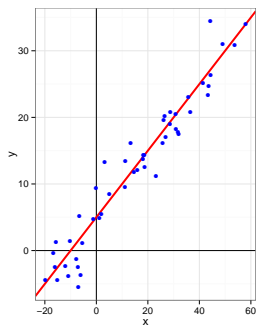
## 3 Conclusion

# The Simple Linear Regression Model

We need to invest in a **probability model**

$$Y = a + bX + \varepsilon$$

- $Y$  is the **response variable**
  - $X$  is a continuous explanatory variable
  - $a$  is the intercept
  - $b$  is the slope
  - $\varepsilon$  is some noise
- 
- $a + bX$  represents the “true line”, the part of  $Y$  that depends on  $X$
  - The error term  $\varepsilon$  is independent “idiosyncratic noise”, i.e., the part of  $Y$  not associated with  $X$

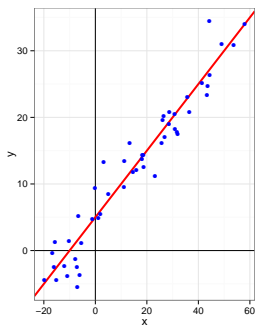


# The Simple Linear Regression Model

We need to invest in a **probability model**

$$Y = a + bX + \varepsilon$$

- $Y$  is the **response variable**
  - $X$  is a continuous explanatory variable
  - $a$  is the intercept
  - $b$  is the slope
  - $\varepsilon$  is some noise
- 
- $a + bX$  represents the “true line”, the part of  $Y$  that depends on  $X$
  - The error term  $\varepsilon$  is independent “idiosyncratic noise”, i.e., the part of  $Y$  not associated with  $X$



## Gauss-Markov Theorem

Under a few assumptions, the least squares regression is the best linear unbiased estimate

- $\mathbb{E}(\hat{\beta}) = b$  and  $\mathbb{E}(\hat{\alpha}) = a$
- $\text{Var}(\hat{\beta})$  and  $\text{Var}(\hat{\alpha})$  are minimal

## Multiple explanatory variables

- The same results hold true when there are **several** explanatory variables:

$$Y = a + b^{(1)}X^{(1)} + b^{(2)}X^{(2)} + b^{(1,2)}X^{(1)}X^{(2)} + \varepsilon$$

The least squares regressions are good estimators of  $a$ ,  $b^{(1)}$ ,  $b^{(2)}$ ,  $b^{(1,2)}$

- We can use an **arbitrary** linear combination of variables, hence

$$Y = a + b^{(1)}X + b^{(2)}\frac{1}{X} + b^{(3)}X^3 + \varepsilon$$

is also a linear model

- Obviously the closed-form formula are much more complicated but softwares like **R** handle this very well

## 1 Simple Linear Regression

- General Introduction
- Fitting a Line to a Set of Points

## 2 Linear Model

- Linear Regression
- **Underlying Hypothesis**
- Checking hypothesis
- Decomposing the Variance
- Making Predictions
- Confidence interval

## 3 Conclusion



# Important Hypothesis (1)

**Weak exogeneity** The predictor variables  $X$  can be treated as fixed values, rather than random variables: the  $X$  are assumed to be **error-free**, i.e., they are not contaminated with measurement errors

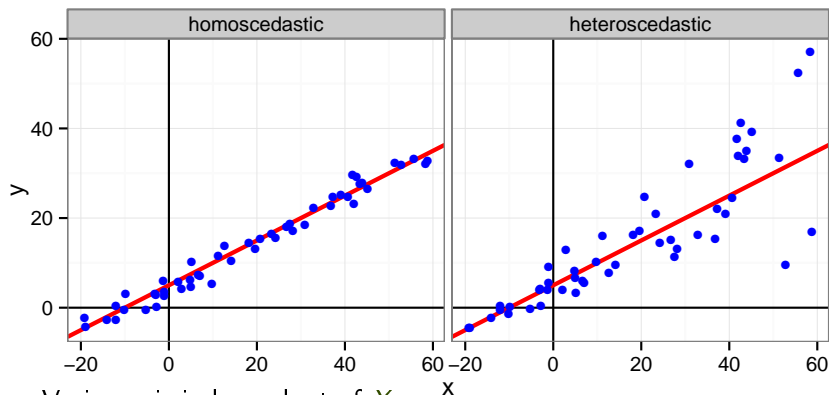
Although not realistic in many settings, dropping this assumption leads to significantly more difficult errors-in-variables models

**Linearity** the mean of the response variable is a linear combination of the parameters (regression coefficients) and the predictor variables  
Since predictor variables themselves can be arbitrarily transformed, this is not that restrictive. This trick is used, for example, in **polynomial regression**, but beware of **overfitting**

**Independence of Errors** if several responses  $Y_1$  and  $Y_2$  are fit,  $\varepsilon_1$  and  $\varepsilon_2$  should be independent

# Other Very Important Hypothesis

## Constant variance (a.k.a. homoscedasticity)

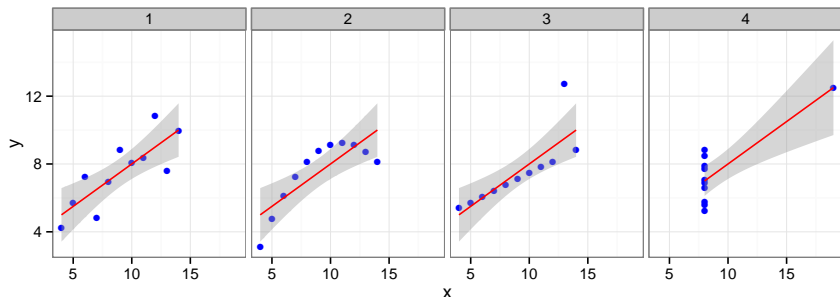


- Variance is independent of  $X$
- If several responses  $Y_1$  and  $Y_2$  are fit,  $\varepsilon_1$  and  $\varepsilon_2$  should have the same variance
- Either normalize  $Y$  or use an other estimator

## Other Classical Hypothesis (3)

**Normal and iid errors** This is **not** an assumption of the Gauss Markov Theorem. Yet, it is quite convenient to build confidence intervals of the regression

**Arrangement of the predictor variables  $X$**  it has a major influence on the precision of estimates of  $\beta$  (remember Anscombe's quartet).



This is part of your design of experiments:

- If you want to test linearity,  $X$  should be uniformly distributed
- If you want the best estimation, you should use extreme values of  $X$

## 1 Simple Linear Regression

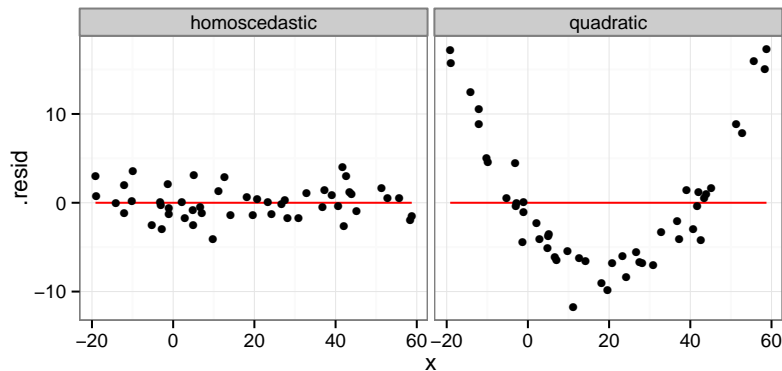
- General Introduction
- Fitting a Line to a Set of Points

## 2 Linear Model

- Linear Regression
- Underlying Hypothesis
- **Checking hypothesis**
- Decomposing the Variance
- Making Predictions
- Confidence interval

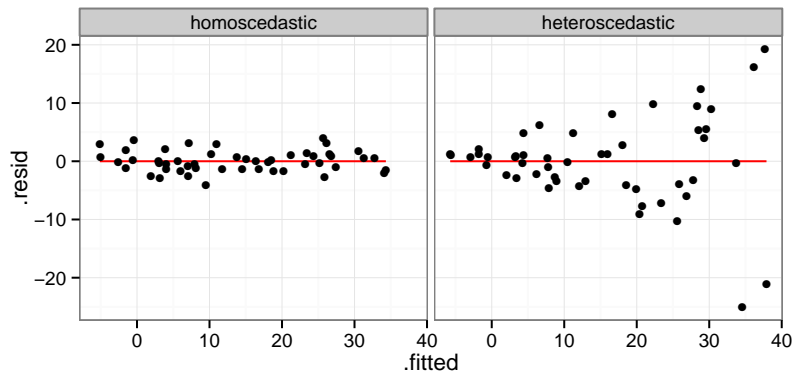
## 3 Conclusion

## Linearity: Residuals vs. Explanatory Variable

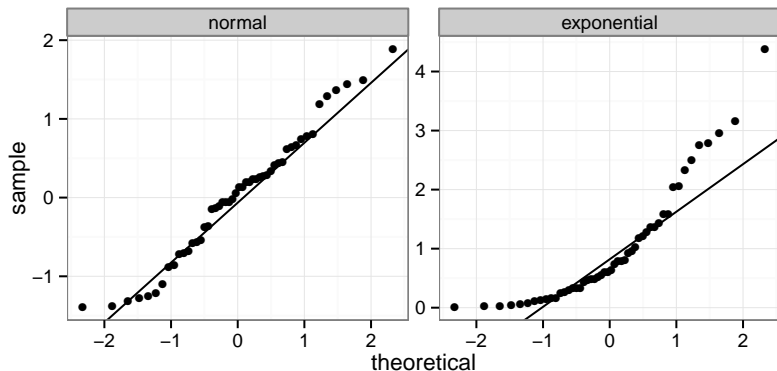


When there are several factors, you have to check for every dimension...

# Homoscedasticity: Residuals vs. Fitted values



## Normality: qqplots



A quantile-quantile plot is a graphical method for comparing two probability distributions by plotting their quantiles against each other

# Model Formulae in R

The structure of a model is specified in the formula like this:

response variable ~ explanatory variable(s)

~ reads "is modeled as a function of" and  $\text{lm}(y \sim x)$  means  $y = \alpha + \beta x + \varepsilon$

On the right-hand side, one should specify how the explanatory variables are combined. The symbols used here have a **different meaning** than in arithmetic expressions

- + indicates a variable inclusion (not an addition)
- - indicates a variable deletion (not a subtraction)
- \* indicates inclusion of variables and their interactions
- : means an interaction

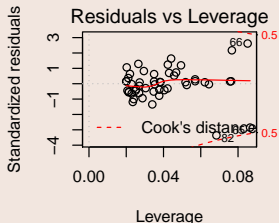
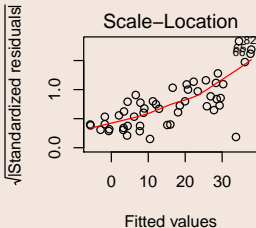
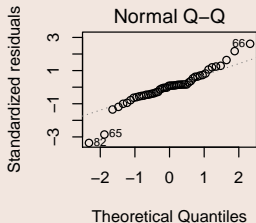
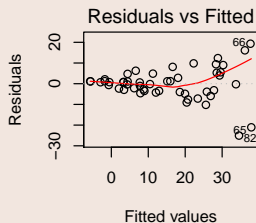
Therefore

- $z \sim x+y$  means  $z = \alpha + \beta_1 x + \beta_2 y + \varepsilon$
- $z \sim x*y$  means  $z = \alpha + \beta_1 x + \beta_2 y + \beta_3 xy + \varepsilon$
- $z \sim (x+y)^2$  means the same
- $\log(y) \sim I(1/x)+x+I(x^2)$  means  $z = \alpha + \beta_1 \times \frac{1}{x} + \beta_2 x + \beta_3 x^2 + \varepsilon$



# Checking the model with R

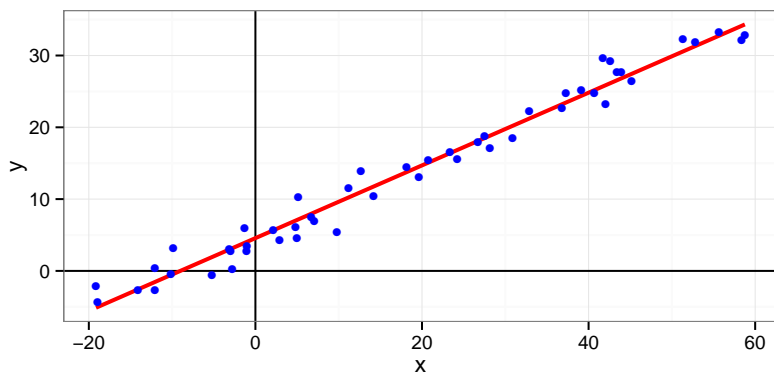
```
reg <- lm(data=df[df$type=="heteroscedastic"],y~x)
par(mfrow=c(2,2)); plot(reg); par(mfrow=c(1,1))
```



- 1 Simple Linear Regression
  - General Introduction
  - Fitting a Line to a Set of Points
- 2 Linear Model
  - Linear Regression
  - Underlying Hypothesis
  - Checking hypothesis
  - **Decomposing the Variance**
  - Making Predictions
  - Confidence interval
- 3 Conclusion

# Decomposing the Variance

How well does the least squares line explain variation in  $Y$ ?



# Decomposing the Variance

How well does the least squares line explain variation in  $Y$ ?

Remember that  $Y = \hat{Y} + \varepsilon$  ( $\hat{Y}$  is the "true mean").

Since  $\hat{Y}$  and  $\varepsilon$  are uncorrelated, we have

$$\text{Var}(Y) = \text{Var}(\hat{Y} + \varepsilon) = \text{Var}(\hat{Y}) + \text{Var}(\varepsilon)$$

$$\frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \frac{1}{n-1} \sum_{i=1}^n (\hat{Y}_i - \bar{\hat{Y}})^2 + \frac{1}{n-1} \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2$$

Since  $\bar{\varepsilon} = 0$  and  $\bar{Y} = \bar{\hat{Y}}$ , we have

$$\underbrace{\sum_{i=1}^n (Y_i - \bar{Y})^2}_{\text{Total Sum of Squares}} = \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{\text{Regression SS}} + \underbrace{\sum_{i=1}^n \varepsilon_i^2}_{\text{Error SS}}$$

- $SSR$  = Variation in  $Y$  explained by the regression line
- $SSE$  = Variation in  $Y$  that is left unexplained

$$SSR = SST \Rightarrow \text{perfect fit}$$

## A Goodness of Fit Measure: $R^2$

The **coefficient of determination**, denoted by  $R^2$ , measures goodness of fit:

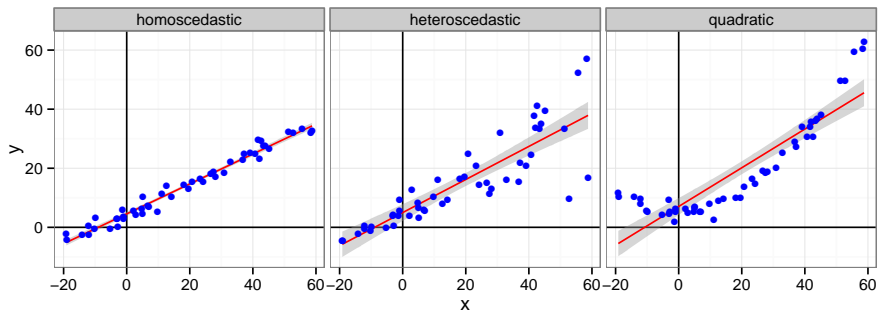
$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

- $0 \leq R^2 \leq 1$
- The closer  $R^2$  is to 1, the better the fit

Warning:

- A not so low  $R^2$  may mean important noise or bad model
- As you add parameters to a model, you inevitably improve the fit. There is a trade-off between model simplicity and fit. Strive for simplicity!

# Illustration with R (homoscedastic data)



## Illustration with R (homoscedastic data)

```
reg <- lm(data=df[df$type=="homoscedastic",],y~x)
summary(reg)
```

## Illustration with R (homoscedastic data)

```
reg <- lm(data=df[df$type=="homoscedastic"],y~x)
summary(reg)
```

```
Call:
lm(formula = y ~ x, data = df[df$type == "homoscedastic", ])
```

Residuals:

	Min	1Q	Median	3Q	Max
	-4.1248	-1.3059	-0.0366	1.0588	3.9965

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	4.56481	0.33165	13.76	<2e-16 ***
x	0.50645	0.01154	43.89	<2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

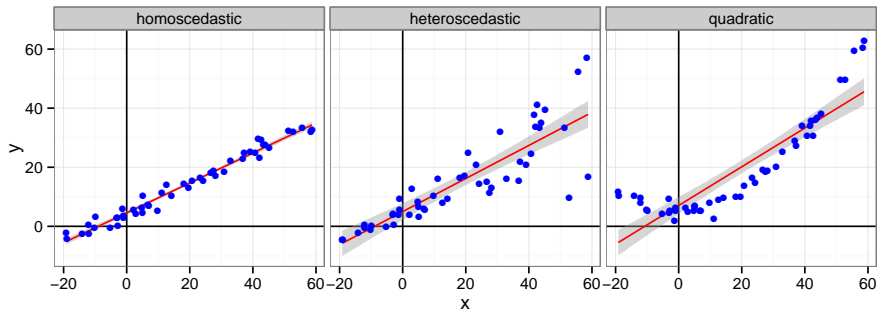
Residual standard error: 1.816 on 48 degrees of freedom

Multiple R-squared: 0.9757, Adjusted R-squared: 0.9752

F-statistic: 1926 on 1 and 48 DF, p-value: < 2.2e-16



# Illustration with R (heteroscedastic data)



## Illustration with R (heteroscedastic data)

```
reg <- lm(data=df[df$type=="heteroscedastic",],y~x)
summary(reg)
```

## Illustration with R (heteroscedastic data)

```
reg <- lm(data=df[df$type=="heteroscedastic",],y~x)
summary(reg)
```

```
Call:
lm(formula = y ~ x, data = df[df$type == "heteroscedastic", ])
```

```
Residuals:
```

	Min	1Q	Median	3Q	Max
	-25.063	-3.472	0.663	3.707	19.327

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	4.98800	1.41061	3.536	0.000911 ***
x	0.56002	0.04908	11.411	2.83e-15 ***

```
---
```

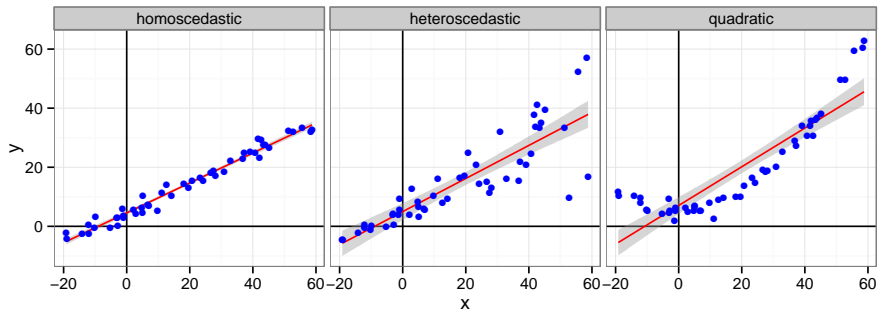
```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 7.722 on 48 degrees of freedom
```

```
Multiple R-squared:  0.7306, Adjusted R-squared:  0.725
```

```
F-statistic: 130.2 on 1 and 48 DF,  p-value: 2.83e-15
```

# Illustration with R (quadratic data)



## Illustration with R (quadratic data)

```
reg <- lm(data=df[df$type=="quadratic"],y~x)
summary(reg)
```

## Illustration with R (quadratic data)

```
reg <- lm(data=df[df$type=="quadratic"],y~x)
summary(reg)
```

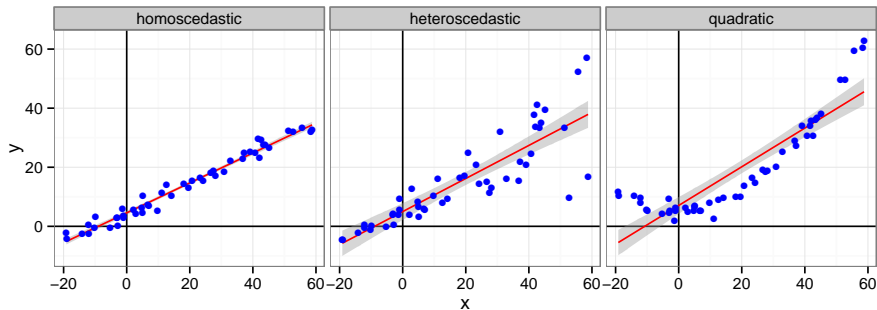
```
Call:
lm(formula = y ~ x, data = df[df$type == "quadratic", ])

Residuals:
    Min       1Q   Median       3Q      Max
-11.759  -5.847  -2.227   3.746  17.346

Coefficients:
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  7.05330    1.41238   4.994 8.23e-06 ***
x            0.65517    0.04914  13.333 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 7.732 on 48 degrees of freedom
Multiple R-squared:  0.7874, Adjusted R-squared:  0.783
F-statistic: 177.8 on 1 and 48 DF,  p-value: < 2.2e-16
```

# Illustration with R (quadratic data, polynomial regression)



## Illustration with R (quadratic data, polynomial regression)

```
df$x2=df$x^2  
reg_quad <- lm(data=df[df$type=="quadratic",],y~x+x2)  
summary(reg_quad)
```



# Illustration with R (quadratic data, polynomial regression)

```
df$x2=df$x^2
reg_quad <- lm(data=df[df$type=="quadratic",],y~x+x2)
summary(reg_quad)
```

```
Call:
lm(formula = y ~ x + x2, data = df[df$type == "quadratic", ])
```

Residuals:

	Min	1Q	Median	3Q	Max
	-4.7834	-0.8638	-0.0480	1.1312	3.9913

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.3065389	0.3348067	15.850	<2e-16 ***
x	0.0036098	0.0252807	0.143	0.887
x2	0.0164635	0.0005694	28.913	<2e-16 ***

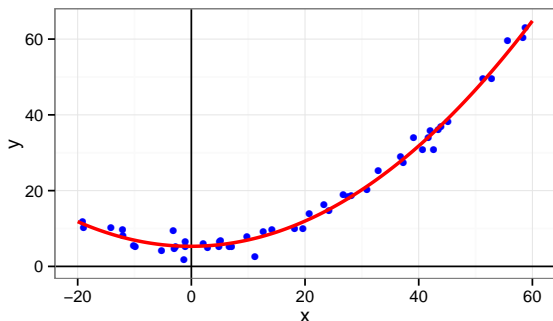
---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.803 on 47 degrees of freedom  
Multiple R-squared: 0.9887, Adjusted R-squared: 0.9882  
F-statistic: 2053 on 2 and 47 DF, p-value: < 2.2e-16

- 1 Simple Linear Regression
  - General Introduction
  - Fitting a Line to a Set of Points
- 2 Linear Model
  - Linear Regression
  - Underlying Hypothesis
  - Checking hypothesis
  - Decomposing the Variance
  - **Making Predictions**
  - Confidence interval
- 3 Conclusion

# Making Predictions

```
xv <- seq(-20,60,.5)
yv <- predict(reg_quad,list(x=xv,x2=xv^2))
ggplot(data=df[df$type=="quadratic",]) + theme_bw() +
  geom_hline(yintercept=0) + geom_vline(xintercept=0) +
  geom_point(aes(x=x,y=y),color="blue") +
  geom_line(data=data.frame(x=xv,y=yv),aes(x=x,y=y),color="red",size=2)
```



## 1 Simple Linear Regression

- General Introduction
- Fitting a Line to a Set of Points

## 2 Linear Model

- Linear Regression
- Underlying Hypothesis
- Checking hypothesis
- Decomposing the Variance
- Making Predictions
- Confidence interval

## 3 Conclusion

Remember that

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \sum_{j=1}^n y_j}{\sum_{i=1}^n (x_i^2) - \frac{1}{n} (\sum_{i=1}^n x_i)^2}$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$\hat{\beta}$  and  $\hat{\alpha}$  are sums of the  $\varepsilon_i$ 's and it is thus possible to compute confidence intervals assuming:

- the linear model holds true
- either the errors in the regression are normally distributed
- or the number of observations is sufficiently large so that the actual distribution of the estimators can be approximated using the central limit theorem

# Illustration with R

## The Anscombe quartet

```
head(a,10)
```

	idx	set	x	y
1	1	1	10	8.04
2	1	2	10	9.14
3	1	3	10	7.46
4	1	4	8	6.58
5	2	1	8	6.95
6	2	2	8	8.14
7	2	3	8	6.77
8	2	4	8	5.76
9	3	1	13	7.58
10	3	2	13	8.74

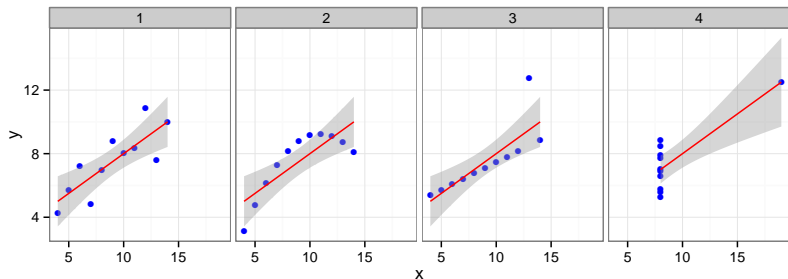
# Illustration with R

## The Anscombe quartet

```
head(a, 10)
```

## Confidence intervals with ggplot

```
ggplot(data=a, aes(x=x, y=y)) + theme_bw() +  
  facet_wrap(~set, nrow=1) + geom_point(color="blue") +  
  geom_smooth(method='lm', color="red")
```



- 1 Simple Linear Regression
  - General Introduction
  - Fitting a Line to a Set of Points
- 2 Linear Model
  - Linear Regression
  - Underlying Hypothesis
  - Checking hypothesis
  - Decomposing the Variance
  - Making Predictions
  - Confidence interval
- 3 Conclusion



# Conclusion

- 1 You need a model to perform your regression
- 2 You need to **check** whether the underlying **hypothesis** of this model are reasonable or not

This model will allow you to:

- 1 **Assess** and **quantify the effect** of parameters on the response
- 2 **Extrapolate within the range** of parameters you tried
- 3 Detect **outstanding** points (those with a high residual and/or with a high lever)

This model will guide on how to design your experiments:

- e.g., the linear model assumes some **uniformity** of interest over the parameter space range
- if your system is heteroscedastic, you will have to perform more measurements for parameters that lead to higher variance