# Linear Regression II

Presented by David John Baker
November 2019

**// FLATIRON SCHOOL**

# Why Linear Regression II ?

- Today we take second pass at Linear Regression
- Goals: Have better understanding of core components of model
- Procede through entire checklist of running linear regression model
- Clearly define more specific terms and explain how they relate to NHST

# Linear Regression

- Can you predict X given Y if we assume a linear relationship between the variables in question?

- Dependant variable is continuous

- Independent variables can be either continuous or categorical

**Format**

- Assumptions

- Fitting

- Multiple Models

- Interpretation

# Common statistical tests are linear models

Last updated: 29 June, 2019. Also check out the R version!

See worked examples and more details at the accompanying notebook: https://github.com/eigenfoo/tests-as-linear

| | Common name | Function in scipy.stats | Equivalent linear model in smf.ols | Exact? | The linear model in words | Icon |
|---|---|---|---|---|---|---|
| **Simple Regression: (y ~ 1 + x)** | **y is independent of x**<br>P: One-sample t-test<br>N: Wilcoxon signed-rank | scipy.stats.ttest_1samp(y)<br>scipy.stats.wilcoxon(y) | smf.ols("y ~ 1", data)<br>smf.ols("y ~ 1", signed_rank(data)) | ✓<br>for N >14 | One number (intercept, i.e., the mean) predicts **y**.<br>- (Same, but it predicts the *signed rank* of **y**.) | |
| | **P: Paired-sample t-test**<br>N: Wilcoxon matched pairs | scipy.stats.ttest_rel(y1, y2)<br>scipy.stats.wilcoxon(y1, y2) | smf.ols("y2_sub_y1 ~ 1", data)<br>smf.ols("y2_sub_y1 ~ 1",<br>signed_rank(data)) | ✓<br>for N >14 | One intercept predicts the pairwise $y_2-y_1$ differences.<br>- (Same, but it predicts the *signed rank* of $y_2-y_1$.) | |
| | **y ~ continuous x**<br>P: Pearson correlation<br>N: Spearman correlation | scipy.stats.pearsonr(x, y)<br>scipy.stats.spearmanr(x, y) | smf.ols("y ~ 1 + x", data)<br>smf.ols("y ~ 1 + x", rank(data)) | ✓<br>for N >10 | One intercept plus **x** multiplied by a number (slope) predicts **y**.<br>- (Same, but with *ranked* **x** and **y**) | |
| | **y ~ discrete x**<br>P: Two-sample t-test<br>P: Welch's t-test<br>N: Mann-Whitney U | scipy.stats.ttest_ind(y1, y2)<br>N/A in Python, but see R version<br>scipy.stats.mannwhitneyu(y1, y1) | smf.ols("y ~ 1 + group", data)[A]<br>N/A in Python, but see R version<br>smf.ols("y ~ 1 + group",<br>signed_rank(data))[A] | ✓<br>for N >11 | An intercept for **group 1** (plus a difference if **group 2**) predicts **y**.<br>- (Same, but with one variance *per group* instead of one common.)<br>- (Same, but it predicts the *signed rank* of **y**.) | |
| **Multiple regression: (y ~ 1 + x₁ + x₂ + …)** | P: One-way ANOVA<br>N: Kruskal-Wallis | scipy.stats.f_oneway(a, b, c)<br>scipy.stats.kruskal(a, b, c) | smf.ols(y ~ 1 + G₂ + G₃ +…+ Gₙ)[A]<br>smf.ols(rank(y) ~ 1 + G₂ + G₃ +…+ Gₙ)[A] | ✓<br>for N >11 | An intercept for **group 1** (plus a difference if group ≠ 1) predicts **y**.<br>- (Same, but it predicts the *rank* of **y**.) | |
| | P: One-way ANCOVA | N/A in Python, but see R version | smf.ols("y ~ 1 + G₂ + G₃ +…+ Gₙ + x",<br>data)[A] | ✓ | - (Same, but plus a slope on **x**.)<br>*Note: this is discrete AND continuous. ANCOVAs are ANOVAs with a continuous x.* | |
| | P: Two-way ANOVA | N/A in Python, but see R version | smf.ols("y ~ 1 + G₂ + G₃ + … + Gₙ +<br>S₂ + S₃+ … + Sₖ +<br>G₂*S₂+G₃*S₃+…+Gₙ*Sₖ", data) | ✓ | Interaction term: changing **sex** changes the **y ~ group** parameters.<br>*Note: G₂ to N is an indicator (0 or 1) for each non-intercept levels of the group variable. Similarly for S₂ to K for sex. The first line (with Gᵢ) is main effect of group, the second (with Sⱼ) for sex and the third is the group × sex interaction. For two levels (e.g. male/female), line 2 would just be "S₂" and line 3 would be S₂ multiplied with each Gᵢ.* | [Coming] |
| | **Counts ~ discrete x**<br>N: Chi-square test | scipy.stats.chisquare(data) | **Equivalent log-linear model**<br>sm.GLM(y ~ 1 + G₂ + G₃ + … + Gₙ +<br>S₂ + S₃+ … + Sₖ +<br>G₂*S₂+G₃*S₃+…+Gₙ*Sₖ, family=…)[A] | ✓ | Interaction term: (Same as Two-way ANOVA.)<br>*Note: Run glm using the following arguments: glm(model, family=poisson())*<br>*As linear-model, the Chi-square test is log(y) = log(N) + log(αᵢ) + log(βⱼ) + log(αᵢβⱼ) where αᵢ and βⱼ are proportions. See more info in the accompanying notebook.* | Same as Two-way ANOVA |
| | N: Goodness of fit | scipy.stats.chi2_contingency(data) | sm.GLM(y ~ 1 + G₂ + G₃ +…+ Gₙ,<br>family=…)[A] | ✓ | (Same as One-way ANOVA and see Chi-Square note.) | 1W-ANOVA |

List of common parametric (P) non-parametric (N) tests and equivalent linear models. The notation y ~ 1 + x is R shorthand for y = 1·b + a·x which most of us learned in school. Models in similar colors are highly similar, but really, notice how similar they *all* are across colors! For non-parametric models, the linear models are reasonable approximations for non-small sample sizes (see "Exact" column and click links to see simulations). Other less accurate approximations exist, e.g., Wilcoxon for the sign test and Goodness-of-fit for the binomial test. The signed rank function is `signed_rank(df) = np.sign(df) * df.rank()`. The variables Gᵢ and Sᵢ are "dummy coded" indicator variables (either 0 or 1) exploiting the fact that when Δx = 1 between categories the difference equals the slope. Subscripts (e.g., G₂ or y₁) indicate different columns in data. lm requires long-format data for all non-continuous models. All of this is exposed in greater detail and worked examples at https://eigenfoo.xyz/tests-as-linear/.
[A] See the note to the two-way ANOVA for explanation of the notation.

Jonas Kristoffer Lindeløv, George Ho
https://lindeloev.net, https://eigenfoo.xyz

https://eigenfoo.xyz/tests-as-linear/

# Basic Linear Regression Assumptions

1. Independence of Data Points
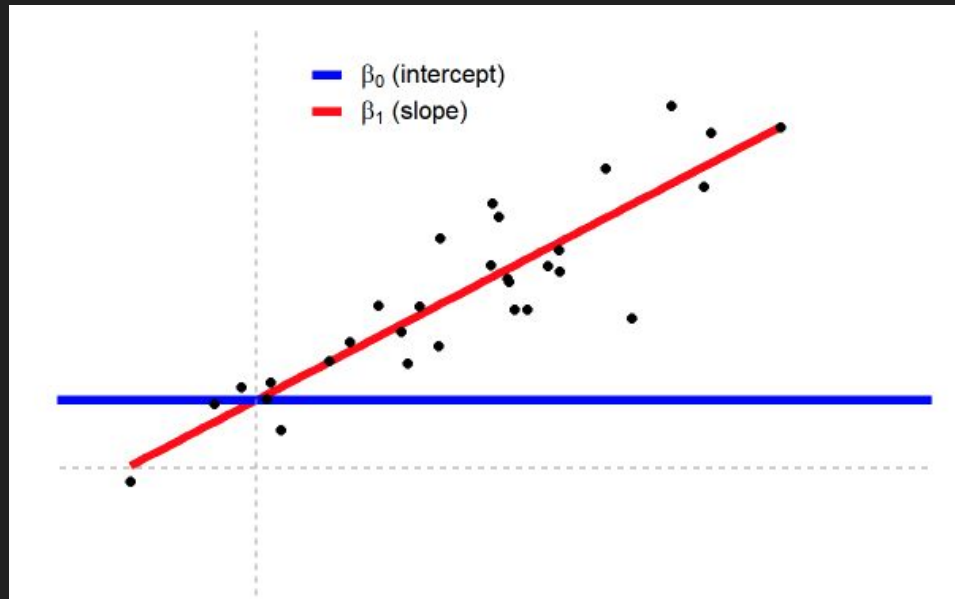2. Normality of Residuals
3. Homoscedasticity

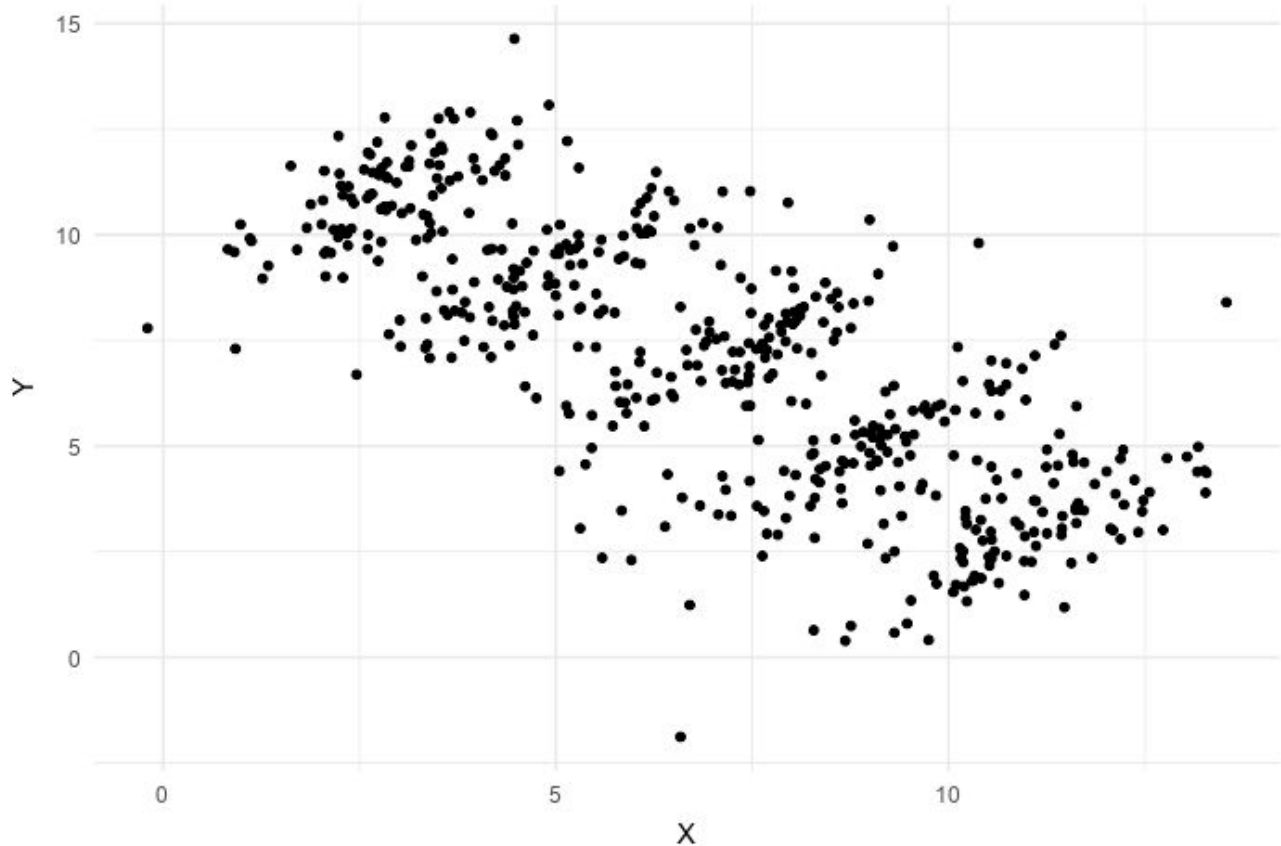# Linear Regression Assumptions

(0. Linearity)

1. Independence of Data Points
2. Normality of Residuals
3. Homoscedasticity
4. Multicollinearity
5. You DO NOT need Normally distributed variables, you DO NEED to Run Diagnostic Plots
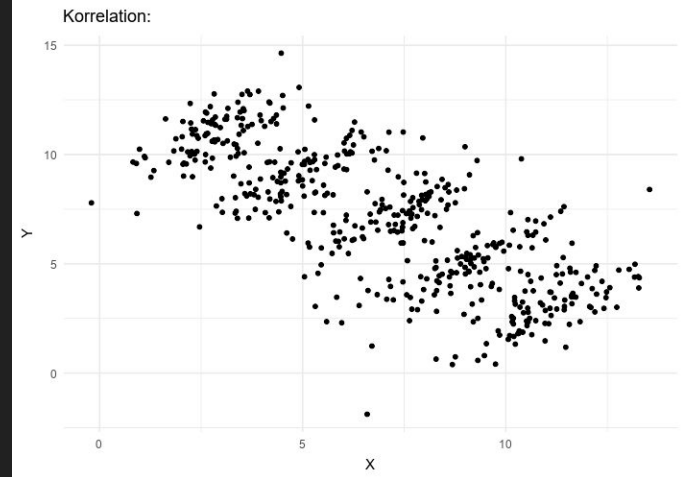
# Independence of Data Points

- Each data point needs to be independently sampled from the parent population

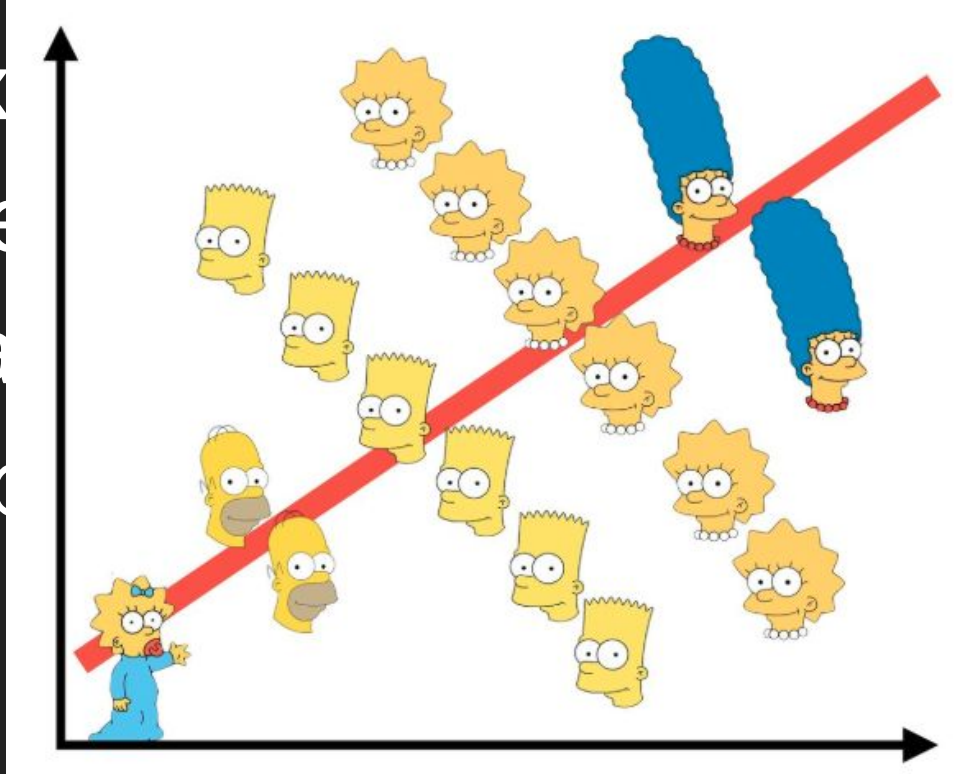- AKA Data points shouldn't be able to influence/talk to one another!!!

Think of situation in your area of expertise where you might come across a violation of the the assumption of independence.

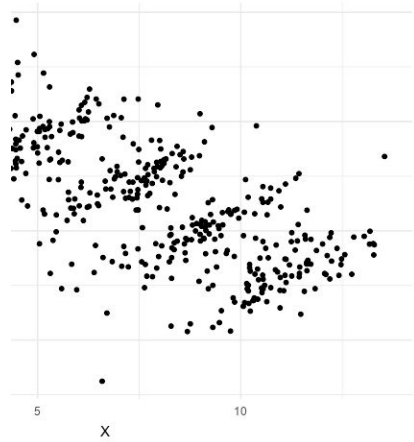Think of situation in your area

of ex                                         t

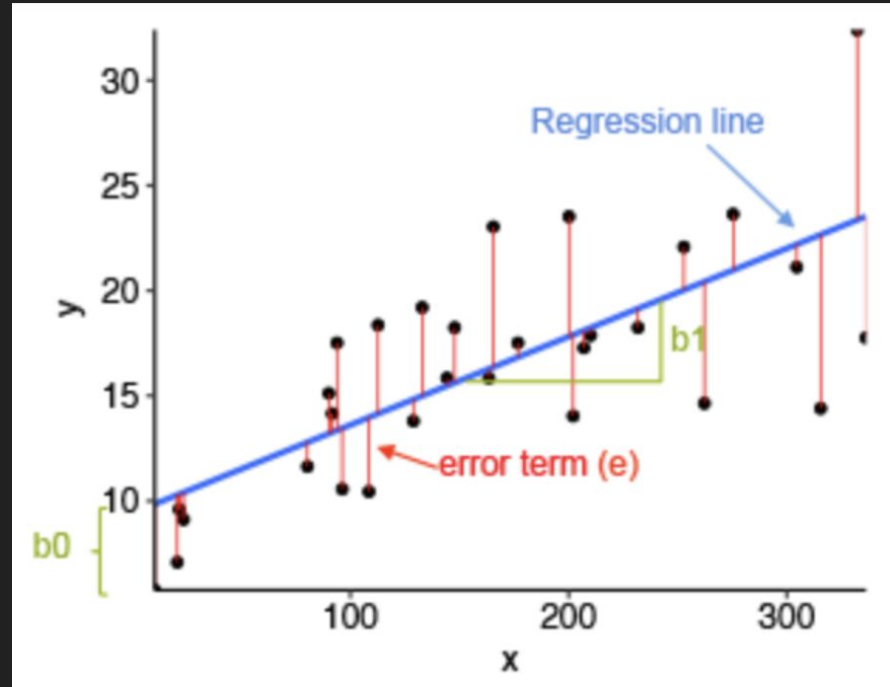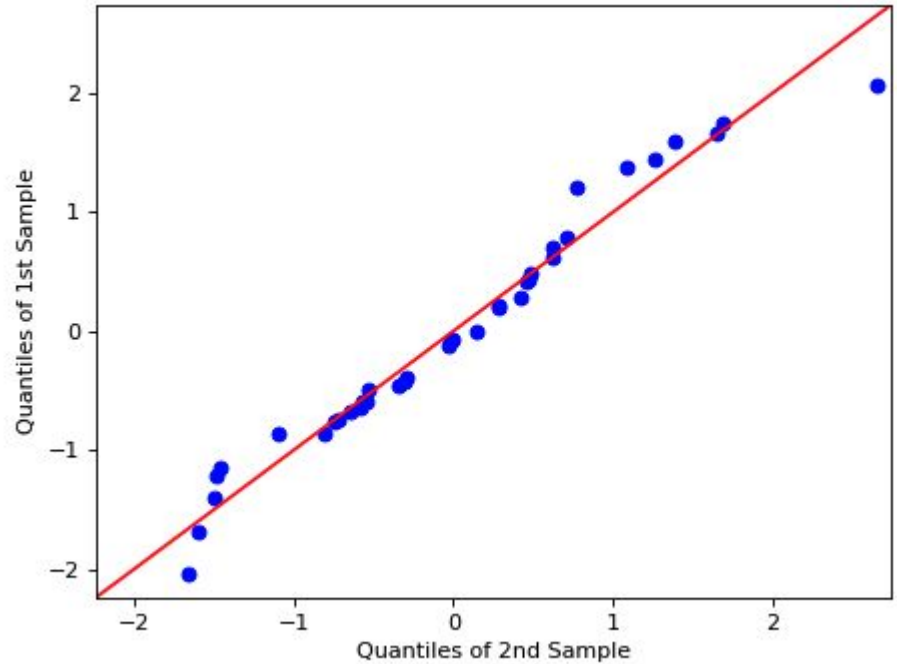come                                          e

the a

inde

# Normally Distributed Residuals

Residuals, or the distances between your regression line and observed data, when plotted should create a normal distribution

Inspect this with a QQ plot

# Normally Distributed Residuals
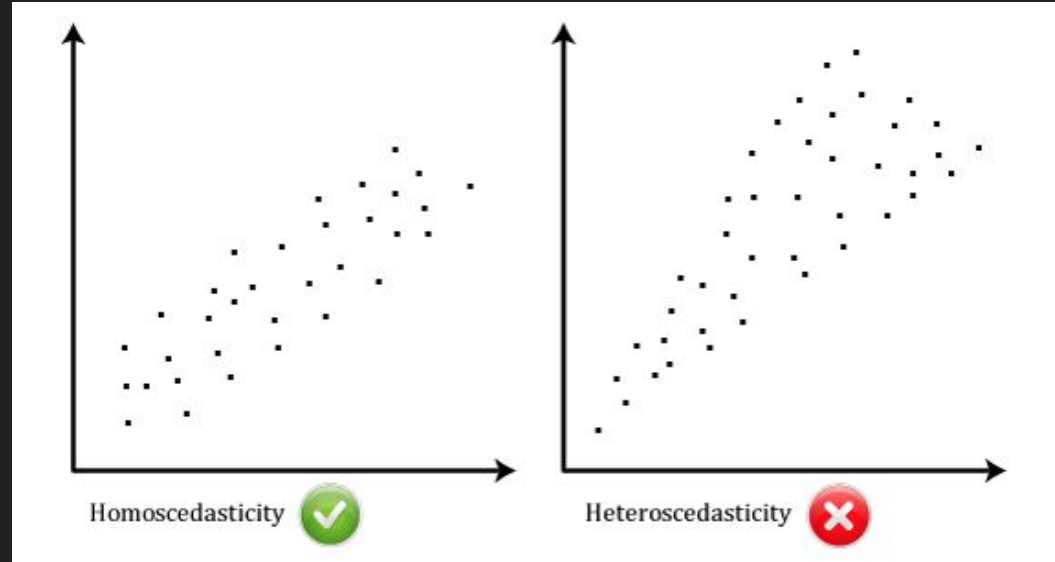
# Homoscedasticity

# Homoscedasticity (ho-mo-ske-das-ti-ci-tee)

# Heteroscedasticity (het-er-o-ske-das-tis-i-tee)

# Homoscedasticity (ho-mo-ske-das-ti-ci-tee) Homogeneity of Variance

- Variance around errors is uniform

# Residual Plots

**Linear Regression Checklist**

- **Plot all variables**
- **Check for Multicollinearity**
- **Check for Outliers (Univariate and Multivariate)**

//

# Regression Output

Out[12]:

## OLS Regression Results

| | | | |
|---|---|---|---|
| **Dep. Variable:** | Sepal.Length | **R-squared:** | 0.867 |
| **Model:** | OLS | **Adj. R-squared:** | 0.862 |
| **Method:** | Least Squares | **F-statistic:** | 155.8 |
| **Date:** | Mon, 18 Nov 2019 | **Prob (F-statistic):** | 3.86e-60 |
| **Time:** | 11:57:26 | **Log-Likelihood:** | -32.558 |
| **No. Observations:** | 150 | **AIC:** | 79.12 |
| **Df Residuals:** | 143 | **BIC:** | 100.2 |
| **Df Model:** | 6 | | |
| **Covariance Type:** | nonrobust | | |

# Regression Output

**R-Squared:**

**Coefficient of Determination**

**0 -- 1**

**Variance explained**

**Literally, r, squared //**



```
Out[12]:    OLS Regression Results
```

| Dep. Variable: | Sepal.Length | R-squared: | 0.867 |
|---|---|---|---|
| Model: | OLS | Adj. R-squared: | 0.862 |
| Method: | Least Squares | F-statistic: | 155.8 |
| Date: | Mon, 18 Nov 2019 | Prob (F-statistic): | 3.86e-60 |
| Time: | 11:57:26 | Log-Likelihood: | -32.558 |
| No. Observations: | 150 | AIC: | 79.12 |
| Df Residuals: | 143 | BIC: | 100.2 |
| Df Model: | 6 | | |
| Covariance Type: | nonrobust | | |

https://rpsychologist.com/d3/correlation/

# Regression Output

F statistic:

Omnibus test,

Just like ANOVA

Variance / Error
//



Out[12]:

OLS Regression Results

| Dep. Variable: | Sepal.Length | R-squared: | 0.867 |
| Model: | OLS | Adj. R-squared: | 0.862 |
| Method: | Least Squares | **F-statistic:** | **155.8** |
| Date: | Mon, 18 Nov 2019 | Prob (F-statistic): | 3.86e-60 |
| Time: | 11:57:26 | Log-Likelihood: | -32.558 |
| No. Observations: | 150 | AIC: | 79.12 |
| Df Residuals: | 143 | BIC: | 100.2 |
| Df Model: | 6 | | |
| Covariance Type: | nonrobust | | |

# Regression Output

Model Fit Comparison Metrics

In statistics, the **likelihood function** (often simply called the **likelihood**) expresses how likely particular values of statistical model parameters are for a given sample of data.[a]

Out[12]: sults

| | Sepal.Length | R-squared: | 0.867 |
| --- | --- | --- | --- |
| | OLS | Adj. R-squared: | 0.862 |
| | Least Squares | F-statistic: | 155.8 |
| Date: | Mon, 18 Nov 2019 | Prob (F-statistic): | 3.86e-60 |
| Time: | 11:57:26 | Log-Likelihood: | -32.558 |
| No. Observations: | 150 | AIC: | 79.12 |
| Df Residuals: | 143 | BIC: | 100.2 |
| Df Model: | 6 | | |
| Covariance Type: | nonrobust | | |

# Regression Output

Model Fit Comparison Metrics

In statistics, the **likelihood function** (often simply called the **likelihood**) expresses how likely particular values of statistical model parameters are for a given sample of data.[a]

The **Akaike information criterion** (**AIC**) is an estimator of out-of-sample prediction error and thereby relative quality of statistical models for a given set of data.[1][2] Given a

In statistics, the **Bayesian information criterion** (**BIC**) or **Schwarz information criterion** (also **SIC**, **SBC**, **SBIC**) is a criterion for model selection among a finite set of models; the model with the lowest BIC is preferred. It is based, in part, on the likelihood function and it is closely related to the Akaike information criterion (AIC).

Out[12]:                      ...sults

| | | | |
|---|---|---|---|
| Sepal.Length | **R-squared:** | | 0.867 |
| OLS | **Adj. R-squared:** | | 0.862 |
| Least Squares | **F-statistic:** | | 155.8 |
| Mon, 18 Nov 2019 | **Prob (F-statistic):** | | 3.86e-60 |
| 11:57:26 | **Log-Likelihood:** | | -32.558 |
| 150 | **AIC:** | | 79.12 |
| 143 | **BIC:** | | 100.2 |
| 6 | | | |
| nonrobust | | | |

# Regression Output -- Unstandardized Beta

For every change in the indicator variable, you have y changes in the dependant variable, all others equal

|  | coef | td err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Unnamed: 0 | -3.606e-05 | 0.002 | -0.020 | 0.984 | -0.004 | 0.003 |
| Sepal.Width | 0.4960 | 0.086 | 5.737 | 0.000 | 0.325 | 0.667 |
| Petal.Length | 0.8290 | 0.070 | 11.888 | 0.000 | 0.691 | 0.967 |
| Petal.Width | -0.3150 | 0.152 | -2.075 | 0.040 | -0.615 | -0.015 |
| Species_setosa | 2.1722 | 0.285 | 7.628 | 0.000 | 1.609 | 2.735 |
| Species_versicolor | 1.4511 | 0.327 | 4.443 | 0.000 | 0.805 | 2.097 |
| Species_virginica | 1.1531 | 0.441 | 2.614 | 0.010 | 0.281 | 2.025 |

//

# Regression Output -- Standardized Beta

**Everything on z score**

**Need to standardize Variables ahead Of time!!!**

|  | coef | td err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Unnamed: 0 | 606e-5 | 0.002 | -0.020 | 0.984 | -0.004 | 0.003 |
| Sepal.Width | 4960 | 0.086 | 5.737 | 0.000 | 0.325 | 0.667 |
| Petal.Length | 0.8 90 | 0.070 | 11.888 | 0.000 | 0.691 | 0.967 |
| Petal.Width | -0.315 | 0.152 | -2.075 | 0.040 | -0.615 | -0.015 |
| Species_seto | 2.1722 | 0.285 | 7.628 | 0.000 | 1.609 | 2.735 |
| Species_versi olor | 1.4511 | 0. 27 | 4.443 | 0.000 | 0.805 | 2.097 |
| Species_ rginica | 1.1531 | 0.44 | 2.614 | 0.010 | 0.281 | 2.025 |

# Regression Output

|  | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Unnamed: 0 | -3.606e-05 | 0.002 | -0.020 | 0.984 | -0.004 | 0.003 |
| Sepal.Width | 0.4960 | 0.086 | 5.737 | 0.000 | 0.325 | 0.667 |
| Petal.Length | 0.8290 | 0.070 | 11.888 | 0.000 | 0.691 | 0.967 |
| Petal.Width | -0.3150 | 0.152 | -2.075 | 0.040 | -0.615 | -0.015 |
| Species_setosa | 2.1722 | 0.285 | 7.628 | 0.000 | 1.609 | 2.735 |
| Species_versicolor | 1.4511 | 0.327 | 4.443 | 0.000 | 0.805 | 2.097 |
| Species_virginica | 1.1531 | 0.441 | 2.614 | 0.010 | 0.281 | 2.025 |

# Regression Output

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Unnamed: 0 | -3.606e-05 | 0.002 | -0.020 | 0.984 | -0.004 | 0.003 |
| Sepal.Width | 0.4960 | 0.086 | 5.737 | 0.000 | 0.325 | 0.667 |
| Petal.Length | 0.8290 | 0.070 | 11.888 | 0.000 | 0.691 | 0.967 |
| Petal.Width | -0.3150 | 0.152 | -2.075 | 0.040 | -0.615 | -0.015 |
| Species_setosa | 2.1722 | 0.285 | 7.628 | 0.000 | 1.609 | 2.735 |
| Species_versicolor | 1.4511 | 0.327 | 4.443 | 0.000 | 0.805 | 2.097 |
| Species_virginica | 1.1531 | 0.441 | 2.614 | 0.010 | 0.281 | 2.025 |

# Regression Output

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Unnamed: 0 | -3.606e-05 | 0.002 | -0.020 | 0.984 | -0.004 | 0.003 |
| Sepal.Width | 0.4960 | 0.086 | 5.737 | 0.000 | 0.325 | 0.667 |
| Petal.Length | 0.8290 | 0.070 | 11.888 | 0.000 | 0.691 | 0.967 |
| Petal.Width | -0.3150 | 0.152 | -2.075 | 0.040 | -0.615 | -0.015 |
| Species_setosa | 2.1722 | 0.285 | 7.628 | 0.000 | 1.609 | 2.735 |
| Species_versicolor | 1.4511 | 0.327 | 4.443 | 0.000 | 0.805 | 2.097 |
| Species_virginica | 1.1531 | 0.441 | 2.614 | 0.010 | 0.281 | 2.025 |

# Regression Output

| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Unnamed: 0 | -3.606e-05 | 0.002 | -0.020 | 0.984 | -0.004 | 0.003 |
| Sepal.Width | 0.4960 | 0.086 | 5.737 | 0.000 | 0.325 | 0.667 |
| Petal.Length | 0.8290 | 0.070 | 11.888 | 0.000 | 0.691 | 0.967 |
| Petal.Width | -0.3150 | 0.152 | -2.075 | 0.040 | -0.615 | -0.015 |
| Species_setosa | 2.1722 | 0.285 | 7.628 | 0.000 | 1.609 | 2.735 |
| Species_versicolor | 1.4511 | 0.327 | 4.443 | 0.000 | 0.805 | 2.097 |
| Species_virginica | 1.1531 | 0.441 | 2.614 | 0.010 | 0.281 | 2.025 |

# Regression Output

| | | | |
|---|---|---|---|
| **Omnibus:** | 0.414 | **Durbin-Watson:** | 1.966 |
| **Prob(Omnibus):** | 0.813 | **Jarque-Bera (JB):** | 0.567 |
| **Skew:** | -0.060 | **Prob(JB):** | 0.753 |
| **Kurtosis:** | 2.723 | **Cond. No.** | 1.98e+03 |

In statistics, the **Durbin–Watson statistic** is a test statistic used to detect the presence of autocorrelation at lag 1 in the residuals (prediction errors) from a regression analysis. It is named after James Durbin and Geoffrey Watson. The small

# Regression Output

| | | | |
|---|---|---|---|
| **Omnibus:** | 0.414 | **Durbin-Watson:** | 1.966 |
| **Prob(Omnibus):** | 0.813 | **Jarque-Bera (JB):** | **0.567** |
| **Skew:** | -0.060 | **Prob(JB):** | 0.753 |
| **Kurtosis:** | 2.723 | **Cond. No.** | 1.98e+03 |

In statistics, the **Jarque–Bera test** is a goodness-of-fit test of whether sample data have the skewness and kurtosis matching a normal distribution. The test is named after Carlos Jarque and Anil K. Bera. The test statistic is always nonnegative. If it is far from zero, it signals the data do not have a normal distribution.

# Regression Output

| | | | |
|---|---|---|---|
| **Omnibus:** | 0.414 | **Durbin-Watson:** | 1.966 |
| **Prob(Omnibus):** | 0.813 | **Jarque-Bera (JB):** | 0.567 |
| **Skew:** | -0.060 | Prob(JB): | 0.753 |
| **Kurtosis:** | 2.723 | **Cond. No.** | 1.98e+03 |

In the field of numerical analysis, the **condition number** of a function measures how much the output value of the function can change for a small change in the input argument. This is used to measure how sensitive a function is to changes or errors in the input, and how much error in the output

# Regression Output

# Regression Output



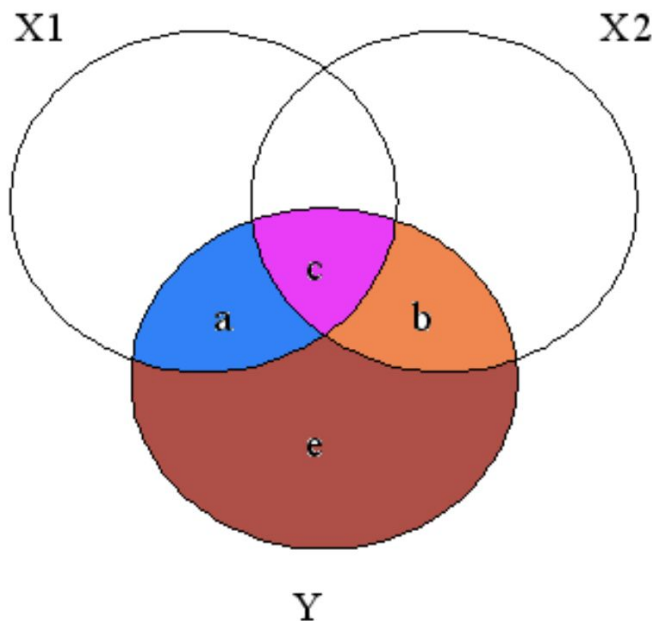Concept of multiple $R^2$ for the full regression model

X1          X2

$$r^2_{Y \cdot 12} = \frac{a+b+c}{a+b+c+e=1} = a+b+c$$

What of Y can be accounted for by X1, X2, and any redundancy among X1 and X2?

Area $c$ is not unique to either X1 or X2 alone, but contributes to the full $R^2$

# Regression Output



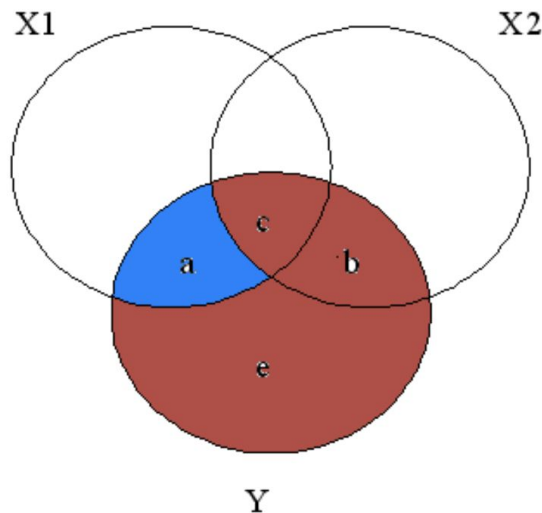Concept of multiple $R^2$ for the full regression model

X1    X2

$$r^2_{Y\cdot12} = \frac{a+b+c}{a+b+c+e=1} = a+b+c$$

What of Y can be accounted for by X1, X2, and any redundancy among X1 and X2?

Area $c$ is not unique to either X1 or X2 alone, but contributes to the full $R^2$

Y

# Regression Output

Influence of predictor X1 when X2 is already in the model: variance shared between X1 & Y **beyond** that accounted for by X2

X1                                    X2



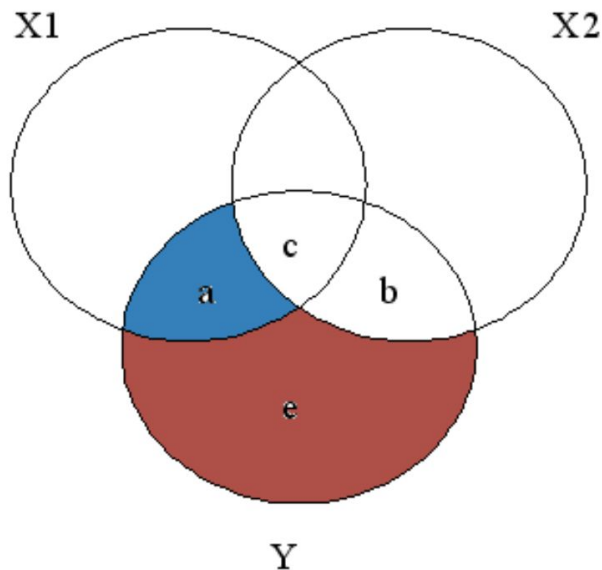Y

"part" or "semipartial" correlation of X1 with Y

$$r^2_{1(Y \cdot 2)} = \frac{a}{a + b + c + e = 1} = a$$

Part correlation: the correlation of Y with that part of X1 that is independent of X2

The squared semipartial correlation is ONE way to determine how much influence (i.e., importance) each predictor has to the full equation

# Regression Output



Partial correlation: X1 with Y when X2's overlap with Y AND X1 is controlled
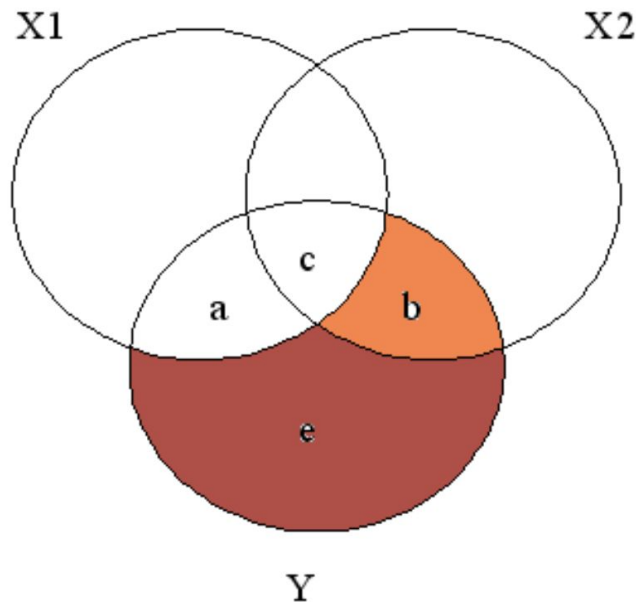
$$r^2_{1Y \cdot 2} = \frac{a}{a+e}$$

Influence of $b$ and $c$ are removed from consideration altogether

<u>Partial</u> correlation: the correlation between the residual of one relationship ($e_{2Y}$) and the residual of another ($e_{12}$) **in the variance of Y**.

# Regression Output



Partial correlation: X2 with Y when X1's overlap with Y AND X2 is controlled

$$r^2_{2Y \cdot 1} = \frac{b}{b + e}$$

Influence of $a$ and $c$ are removed from consideration altogether

Partial correlation: the correlation between the residual of one relationship ($e_{1Y}$) and the residual of another ($e_{12}$) **in the variance of Y**.