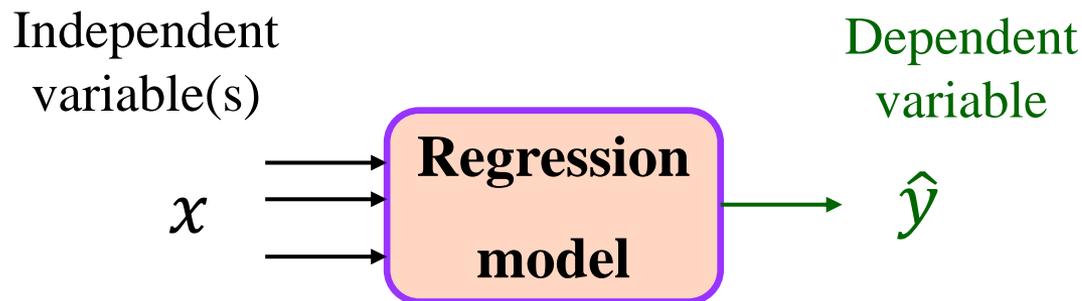


Regression

- Simple Linear Regression
- Multiple Linear Regression
- Polynomial “Linear” Regression

Regression analysis

- In statistical modeling, regression analysis is a set of statistical processes for estimating the relationships among variables.
- It includes many techniques for modeling and analyzing several variables, when the focus is on the **relationship between a dependent variable and one or more independent variables** (or 'predictors').



More specifically, regression analysis helps one understand how the typical value of the dependent variable (or 'criterion variable') changes when any one of the independent variable is varied, while the other independent variables are held fixed.

[Regression analysis, Wikipedia, https://en.wikipedia.org/wiki/Regression_analysis]

- ❑ The regression is an approach to **model** the relationship between a scalar response (**dependent variable / regressor**) and one or more input variables (**independent variables**).
- ❑ Regression models (both linear and non-linear) are **machine-learning models**; used for **predicting/forecasting**.
- ❑ Regression models are used for **predicting a real value**, for example, salary or height. If the independent variable is time, then you are forecasting future values. Otherwise, the model is predicting present but unknown values.
- ❑ A regression model have to **learn the correlation between data**.

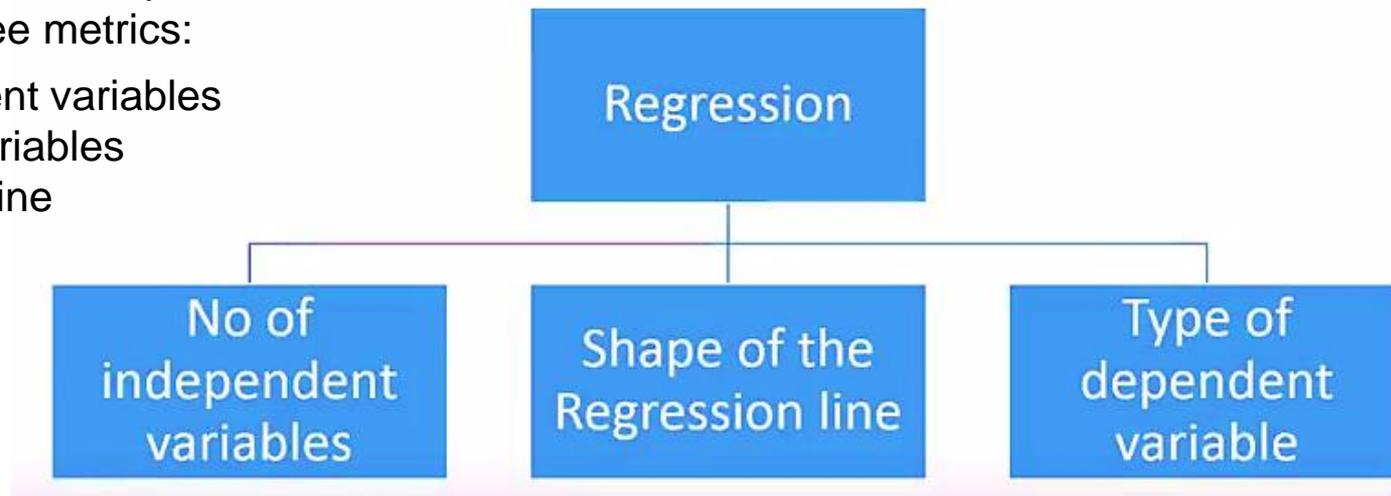
The case of **one input** variable (explanatory variable; independent variable) is called **simple linear regression**.

For **more input** variables (explanatory variables; independent variable), the process is called **multiple linear regression**.

Regression analysis is an important tool for **modelling** and **analyzing** data.

There are various kinds of regression techniques available to make predictions. These techniques are mostly driven by three metrics:

- number of independent variables
- type of dependent variables
- shape of regression line



[7 Types of Regression Techniques you should know!,

<https://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/>]

Case study – build a regression model to predict the salary in a company for a new employee according with his/her years of experience in the workforce.

The model will be built based on a set of data

- 30 observations from that company

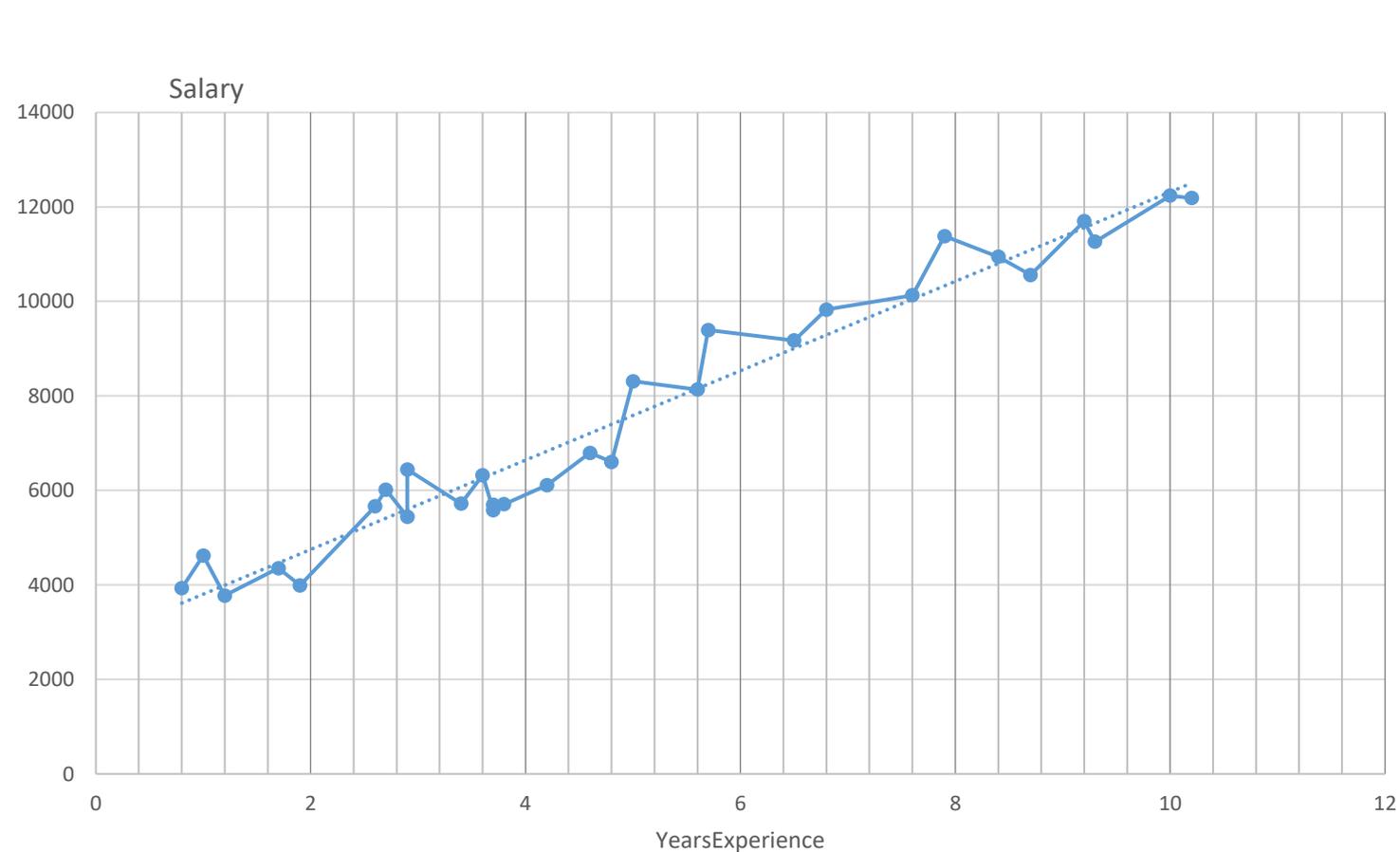
Data set

What is the correlation between years of experience and the salary?

YearsExperience	Salary
0.8	3934.3
1	4620.5
1.2	3773.1
1.7	4352.5
1.9	3989.1
2.6	5664.2
2.7	6015
2.9	5444.5
2.9	6444.5
3.4	5718.9
3.6	6321.8
3.7	5579.4
3.7	5695.7
3.8	5708.1
4.2	6111.1
4.6	6793.8
4.8	6602.9
5	8308.8
5.6	8136.3
5.7	9394
6.5	9173.8
6.8	9827.3
7.6	10130.2
7.9	11381.2
8.4	10943.1
8.7	10558.2
9.2	11696.9
9.3	11263.5
10	12239.1
10.2	12187.2



Data set representation



YearsExperience	Salary
0.8	3934.3
1	4620.5
1.2	3773.1
1.7	4352.5
1.9	3989.1
2.6	5664.2
2.7	6015
2.9	5444.5
2.9	6444.5
3.4	5718.9
3.6	6321.8
3.7	5579.4
3.7	5695.7
3.8	5708.1
4.2	6111.1
4.6	6793.8
4.8	6602.9
5	8308.8
5.6	8136.3
5.7	9394
6.5	9173.8
6.8	9827.3
7.6	10130.2
7.9	11381.2
8.4	10943.1
8.7	10558.2
9.2	11696.9
9.3	11263.5
10	12239.1
10.2	12187.2

Simple Linear Regression

$$\hat{y} = ax + b$$

a – coefficient (**slope**)

b – constant (**intercept**)

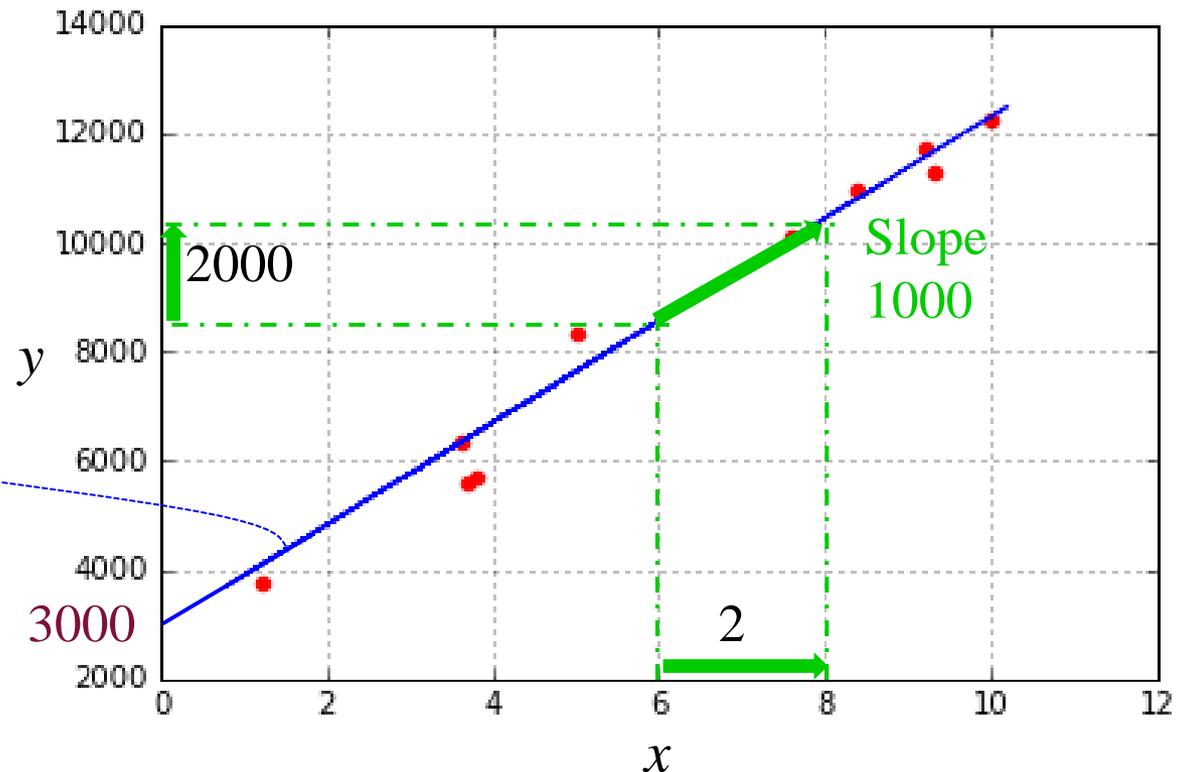
intercept – y value where the line cuts the Y axis

y – the output / dependent variable (DV)

x – the input / independent variable (IV)

Regression line

$$\hat{y} = 1000x + 3000$$



Red dots – facts; blue line – best fits the facts (the data) – linear regression

Linear regression: a trend line that best fits the data

Verify on the dataset

$$\hat{y} = 1000x + 3000$$

$$x = 5, \quad y = 8308.8$$

$$\hat{y} = 1000 \cdot 5 + 3000 = 8000$$

$$y - \hat{y} = 8308.8 - 8000 = -308.8$$

Using the regression model

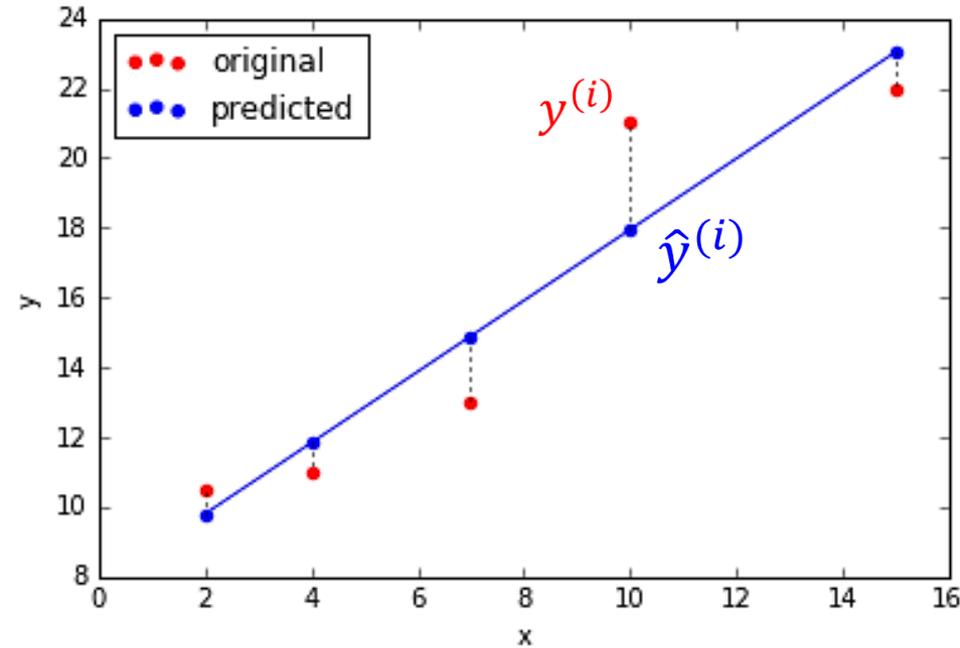
Determine the Salary (dependent variable) for a new value of Years Experience (independent variable)

$$x = 9$$

$$\hat{y} = 1000 \cdot 9 + 3000 = 12000$$

YearsExperience	Salary
0.8	3934.3
1	4620.5
1.2	3773.1
1.7	4352.5
1.9	3989.1
2.6	5664.2
2.7	6015
2.9	5444.5
2.9	6444.5
3.4	5718.9
3.6	6321.8
3.7	5579.4
3.7	5695.7
3.8	5708.1
4.2	6111.1
4.6	6793.8
4.8	6602.9
5	8308.8
5.6	8136.3
5.7	9394
6.5	9173.8
6.8	9827.3
7.6	10130.2
7.9	11381.2
8.4	10943.1
8.7	10558.2
9.2	11696.9
9.3	11263.5
10	12239.1
10.2	12187.2

Ordinary Least Square



$$\sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

to be minimized

A mathematical procedure for finding the best-fitting curve to a given set of points by minimizing the **sum of the squares (Least Square - LS) of the offsets ("the residuals")** of the points from the curve.

The sum of the *squares* of the offsets is used instead of the offset absolute values because this allows the residuals to be treated as a continuous differentiable quantity.

Outlying points can have a disproportionate effect on the fit, a property which may or may not be desirable depending on the problem at hand.

Errors and residuals

y - target (ground truth, target, original, observed)

\hat{y} - predicted (estimated value)

In statistics and optimization, **errors** and **residuals** are two closely related and easily confused measures of the deviation of an observed value of an element of a statistical sample from its "theoretical value".

[https://en.wikipedia.org/wiki/Errors_and_residuals]

The **residual** of an observed value is the difference between the observed value and the *estimated* value of the quantity of interest.

In a linear regression context, residuals applies to the dataset (training, test, validation). (observed value – target value, dependent variable value): $y - \hat{y}$. The residuals are observable.

The **error** (or **disturbance**) of an observed value is the deviation of the observed value from the (unobservable) *true* value of a quantity of interest.

In a linear regression context, error refers to the results in the model utilization phase. (observed value – predicted value) $y_{true} - \hat{y}$. Because we really don't know the true value, the error is unknown.

The above terms come from the statistics realm.

In the context of **machine learning**, the term "**error**" (singular) means the **difference between predicted and observed values**, and the term "residual(s)" is practically almost never used.

$$error = \hat{y} - y$$

Quality of regression

For a dataset with m examples:

$y^{(i)}$ denotes the i^{th} example (the target)
 $\hat{y}^{(i)}$ (the prediction)

Error $\hat{y}^{(i)} - y^{(i)}$

Squared error $(\hat{y}^{(i)} - y^{(i)})^2$

Mean squared-error $\frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$

Root mean squared-error $\sqrt{\frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2}$

R^2 - R squared

In statistics, the **coefficient of determination**, denoted **R^2 (R squared)**, is the proportion of the variance in the dependent variable that is predictable from the independent variable(s)

$$R^2 \equiv 1 - \frac{SS_{res}}{SS_{tot}} = \frac{SS_{tot} - SS_{res}}{SS_{tot}}$$

Best possible
value is **1.0**

$$SS_{res} = \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

Residual sum of squares

$$SS_{tot} = \sum_{i=1}^m (y_i - \bar{y})^2$$

Total sum of squares

$$\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i$$

Mean of the data

[Coefficient of determination, https://en.wikipedia.org/wiki/Coefficient_of_determination]



Python code for linear regression

```
8 # %% Simple linear regression: 30 observations: years of experience - salar
9
10 # Importing the library
11 import numpy as np # math tools
12 import matplotlib.pyplot as plt # for plotting charts
13 import pandas as pd # import and manage datasets
14
15 # %% Importing the dataset (.csv)
16 dataset = pd.read_csv('Salary_Data_ICSDC.csv')
17 # print('\n** The data set is: \n\n', dataset)
18 # x = dataset.iloc[:,0].values # IV (independent variable)
19 x = dataset.iloc[:, :-1].values # both are correct for x
20 y = dataset.iloc[:, 1].values # DV (dependent variable)
21
22 # %% Split the dataset into the Training set and Test set
23 from sklearn.cross_validation import train_test_split
24 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 1/3,
25     random_state=0)
26
27 # %% Fitting the Simple Linear Regression to the Training set
28 from sklearn.linear_model import LinearRegression
29 # from .linear_model library import LinearRegression class
30
31 regressor = LinearRegression() # build our own object named "regressor"
32 regressor.fit(x_train, y_train) # use the method .fit of our object
33
34 # %% Predicting for the x_test
35 y_pred = regressor.predict(x_test)
```



```
37 # %% Plot some graph
38 plt.close('all') # close already existing plots
39 # Trainig set
40 plt.figure()
41 plt.scatter(x_train, y_train, color = 'green')
42 plt.plot(x_train, regressor.predict(x_train),color = 'blue')
43 plt.title('Salary vs. Experience(years) - Training set')
44 plt.xlabel('Years of Experience')
45 plt.ylabel('Salary')
46 plt.show()
47 # Test set
48 plt.figure()
49 plt.scatter(x_test, y_test, color = 'red')
50 plt.plot(x_train, regressor.predict(x_train),color = 'blue')
51 plt.title('Salary vs. Experience(years) - Test set')
52 plt.xlabel('Years of Experience')
53 plt.ylabel('Salary')
54 plt.show()
55 # All data
56 plt.figure()
57 plt.scatter(x_train, y_train, color = 'green')
58 plt.scatter(x_test, y_test, color = 'red')
59 plt.plot(x_train, regressor.predict(x_train),color = 'blue')
60 plt.title('Salary vs. Experience(years) - All data')
61 plt.xlabel('Years of Experience')
62 plt.ylabel('Salary')
63 plt.show()
```



Training
set

Slope = 934.59424431

Intercept = 2961.9974976967906

Salary = 935 * Years + 2962

MSE = 368529.5

RMSE = 607.1

$R^2 = 0.9382$

Results

Test set

MSE = 210260.4

RMSE = 458.5

$R^2 = 0.9749$



test [3773.1, 12239.1, 5708.1, 6321.8, 11696.9, 10943.1, 11263.5, 5579.4, 8308.8, 10130.2]

pred [4083.5, 12307.9, 6513.5, 6326.5, 11560.3, 10812.6, 11653.7, 6420.0, 7635.0, 10064.9]

pred-test [310.4, 68.8, 805.4, 4.7, -136.6, -130.5, 390.2, 840.6, -673.8, -65.2]

Results

All data



Regression in Excel

	A	B
1	YearsExperienc	Salary
2	0.8	3934.3
3	1	4620.5
4	1.2	3773.1
5	1.7	4352.5
6	1.9	3989.1
7	2.6	5664.2
8	2.7	6015
9	2.9	5444.5
10	2.9	6444.5
11	3.4	5718.9
12	3.6	6321.8
13	3.7	5579.4
14	3.7	5695.7
15	3.8	5708.1
16	4.2	6111.1
17	4.6	6793.8
18	4.8	6602.9
19	5	8308.8
20	5.6	8136.3
21	5.7	9394
22	6.5	9173.8
23	6.8	9827.3
24	7.6	10130.2
25	7.9	11381.2
26	8.4	10943.1
27	8.7	10558.2
28	9.2	11696.9
29	9.3	11263.5
30	10	12239.1
31	10.2	12187.2

Data
Data Analysis

The image shows two overlapping dialog boxes from Microsoft Excel. The top dialog is the 'Data Analysis' tool, with 'Regression' selected in the list of analysis tools. The bottom dialog is the 'Regression' configuration window, which is set up with the following parameters:

- Input Y Range:** \$B\$2:\$B\$31
- Input X Range:** \$A\$2:\$A\$31
- Labels:** (unchecked)
- Constant is Zero:** (unchecked)
- Confidence Level:** 95 %
- Output options:**
 - Output Range:** \$D\$2:\$D\$31 (selected)
 - New Worksheet Ply:** (unchecked)
 - New Workbook:** (unchecked)
- Residuals:**
 - Residuals
 - Standardized Residuals
 - Residual Plots
 - Line Fit Plots
- Normal Probability:**
 - Normal Probability Plots

Regression in Excel - Results

RESIDUAL OUTPUT			
<i>Observation</i>	<i>Predicted Y</i>	<i>Residuals</i>	<i>Standard Residuals</i>
1	3618.715875	315.5841248	0.554859429
2	3807.715122	812.7848783	1.42903688
3	3996.714368	-223.6143681	-0.393158371
4	4469.212484	-116.7124842	-0.205203675
5	4658.211731	-669.1117306	-1.176431016
6	5319.709093	344.4909069	0.60568328
7	5414.208716	600.7912837	1.056310132
8	5603.207963	-158.7079627	-0.279040049
9	5603.207963	841.2920373	1.479158115
10	6075.706079	-356.8060788	-0.627335792
11	6264.705325	57.09467477	0.100383752
12	6359.204948	-779.8049484	-1.371051628
13	6359.204948	-663.5049484	-1.166573182
14	6453.704572	-745.6045717	-1.310920589
15	6831.703065	-720.6030645	-1.266962985
16	7209.701557	-415.9015574	-0.731237354
17	7398.700804	-795.8008038	-1.399175512
18	7587.70005	721.0999498	1.267836607
19	8154.69779	-18.39778953	-0.03234696
20	8249.197413	1144.802587	2.012789806
21	9005.194398	168.6056015	0.296442059
22	9288.693268	538.6067319	0.946977367
23	10044.69025	85.50974618	0.150343079
24	10328.18912	1053.010877	1.851401789
25	10800.68724	142.4127605	0.250389854
26	11084.18611	-525.9861092	-0.924787811
27	11556.68423	140.2157748	0.246527118
28	11651.18385	-387.6838485	-0.68162503
29	12312.68121	-73.58121097	-0.12937035
30	12501.68046	-314.4804574	-0.552918963

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.978242
R Square	0.956957
Adjusted R Square	0.955419
Standard Error	578.8315
Observations	30

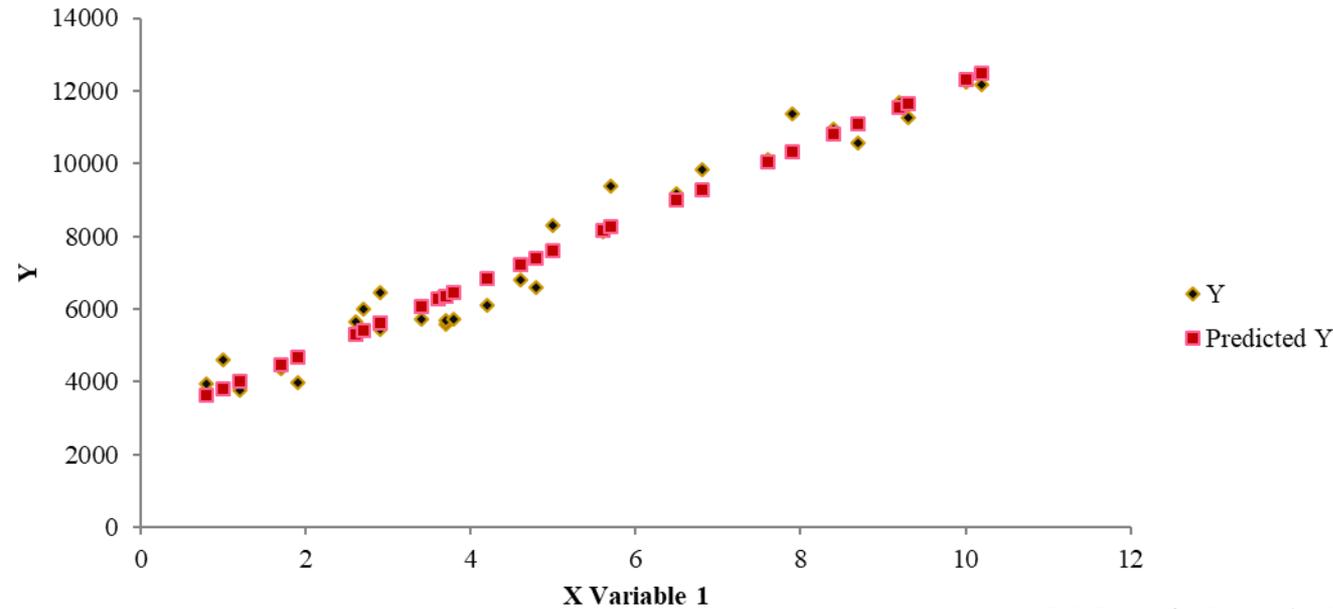
ANOVA					
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	1	208568493	2.09E+08	622.5072	1.1431E-20
Residual	28	9381285.517	335045.9		
Total	29	217949778.5			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	2862.71889	217.3096462	13.17346	1.6E-13	2417.58026	3307.857521
X Variable 1	944.9962321	37.87545742	24.95009	1.14E-20	867.411875	1022.58059

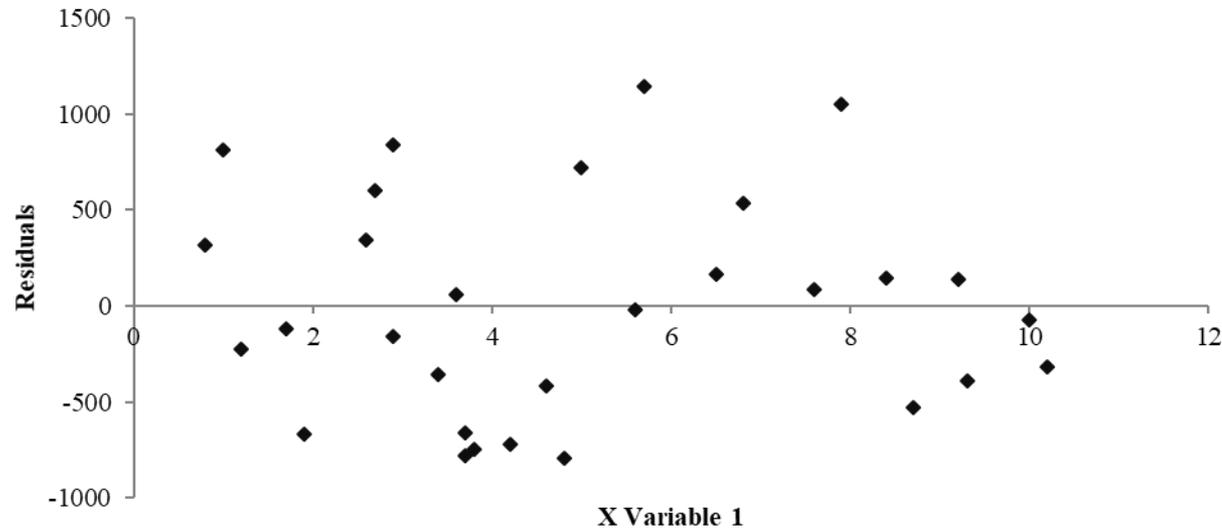


Regression in Excel Results

X Variable 1 Line Fit Plot



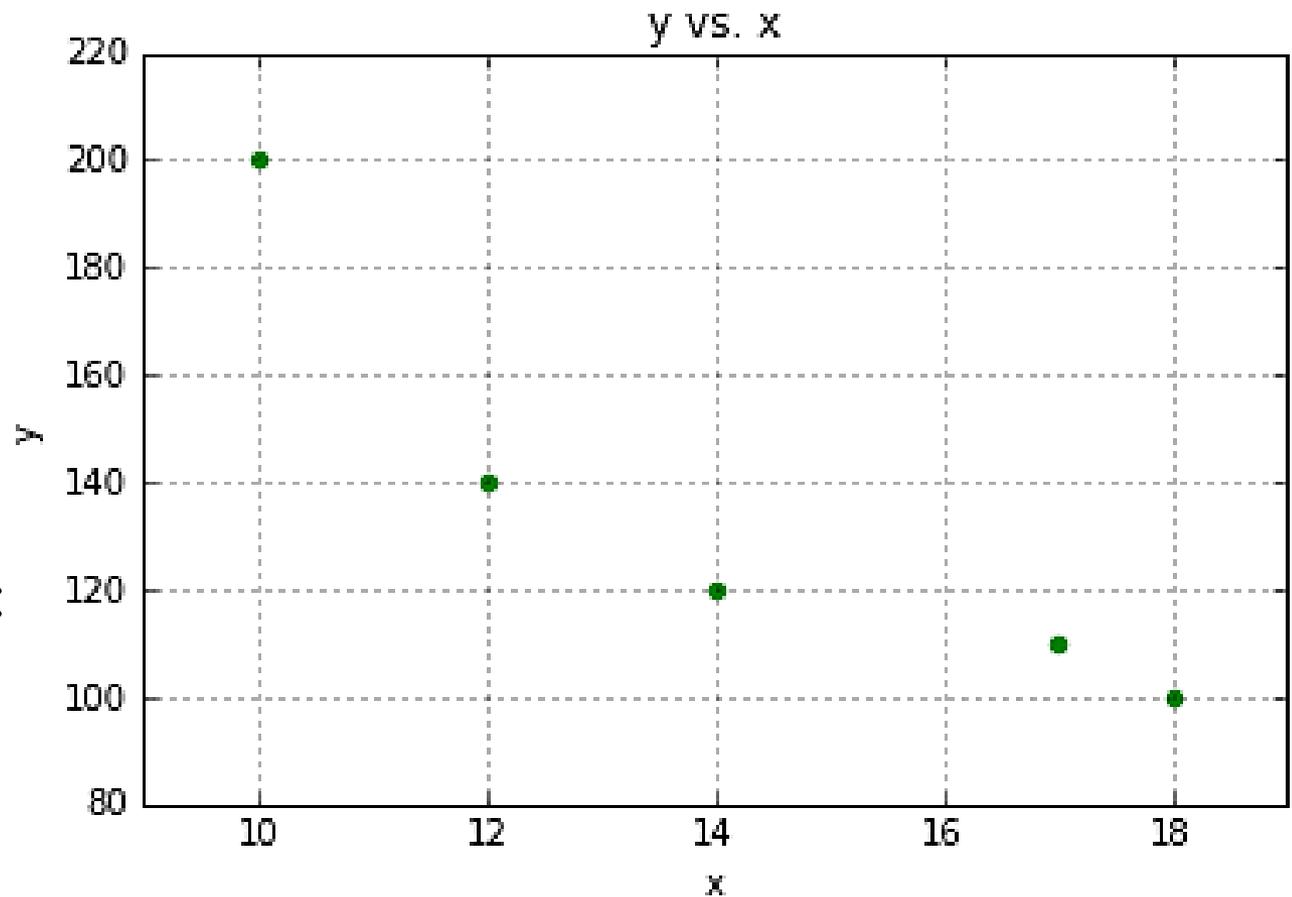
X Variable 1 Residual Plot



Problem

Using the OLS linear regression to model the $y(x)$ relation, the result is:

slope = -10.8;
intercept = 287.



- What is the equation of the linear regression model?
- Plot the linear model on the same diagram.
- What are the errors?
- What is the interpretation of the following error measures:

$$\text{MSE} = 218.2; \text{RMSE} = 14.77; R^2 = 0.8274$$

Multiple Linear Regression

$$\hat{y} = a_1x_1 + a_2x_2 + \dots a_nx_n + b$$

y – dependent variable (DV) / regressor

$x_1, x_2, x_3, \dots, x_n$ - independent variables (IVs) / predictors

$a_1, a_2, a_3, \dots, a_n$ - coefficients

b - constant

5 methods of building multiple linear regression models:

1. All-in
 2. Backward Elimination
 3. Forward Selection
 4. Bidirectional Elimination
 5. Score Comparison
- } Stepwise regression

[Kirill Ermenko, Building a Model (Step-By-Step), Data Science Training,



Multiple Linear Regression - *Backward elimination*

Usually, we are using all dependent variables; but is this the optimal model?

Some independent variables (IV) can be highly statistically significant with great impact (effect) on the DV (dependent variable)

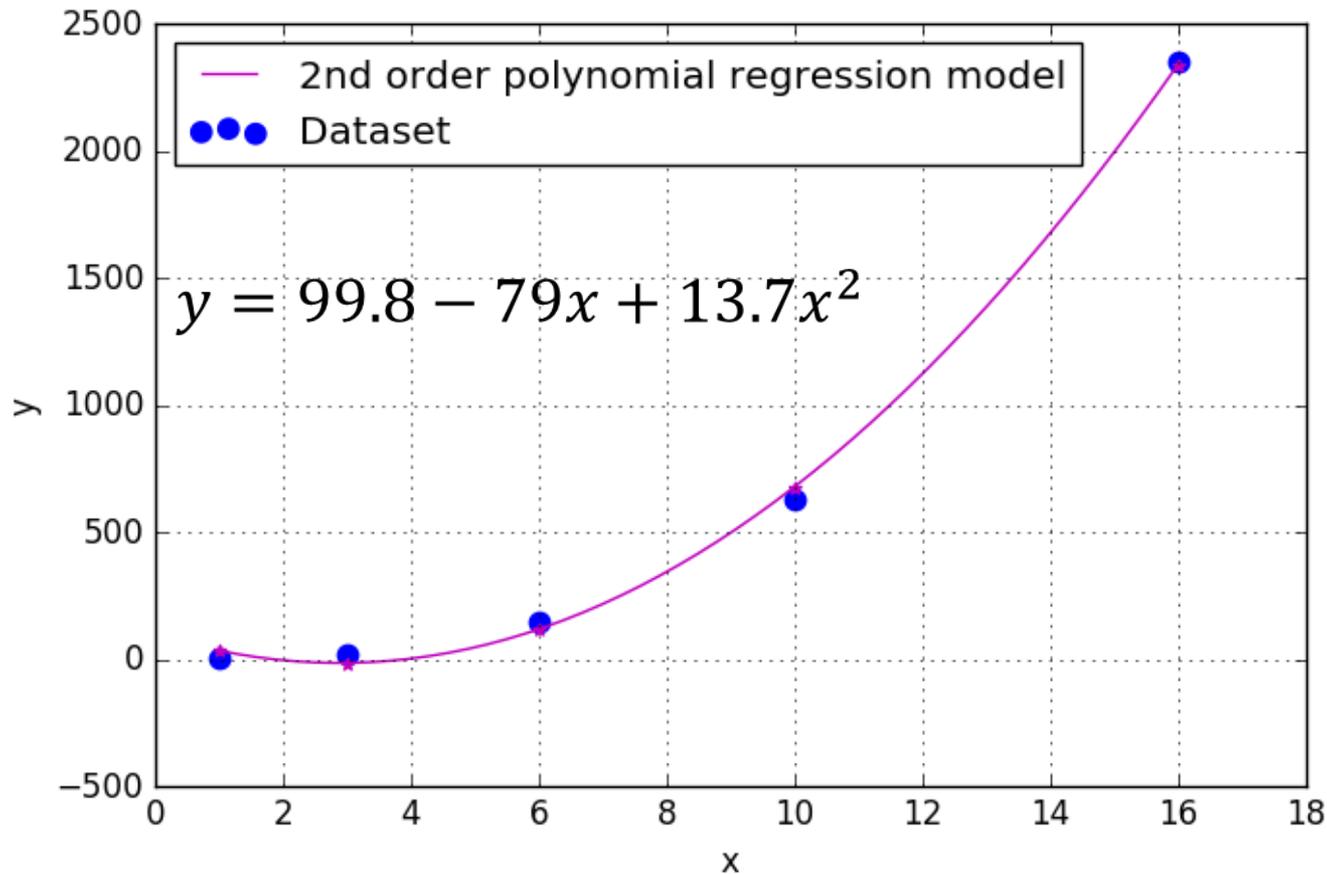
Some IVs are not statistically significant at all – should be removed from the model.

Find a team of optimal IVs, where each IV of the team has great impact on the DV (statistically significant)

Polynomial „linear” regression

One independent variable x

$$\hat{y} = a_1x + a_2x^2 + \dots a_nx^n + b$$



Can be considered as a special case of a multiple linear regression – from the point of view of a_i coefficient

In statistics, **polynomial regression** is a form of regression analysis in which the relationship between the independent variable x and the dependent variable y is modelled as an m^{th} degree polynomial in x .

Polynomial regression fits a **nonlinear relationship** between the value of x and the corresponding conditional mean of y , denoted $E(y | x)$.

Although polynomial regression fits a nonlinear model to the data, **as a statistical estimation problem it is linear**, in the sense that the regression function $E(y | x)$ is linear in the unknown parameters (a_0, a_1, \dots, a_m) that are estimated from the data.

For this reason, polynomial regression is considered to be a special case of multiple linear regression [https://en.wikipedia.org/wiki/Polynomial_regression]

A linear combination from the coefficient point a view.

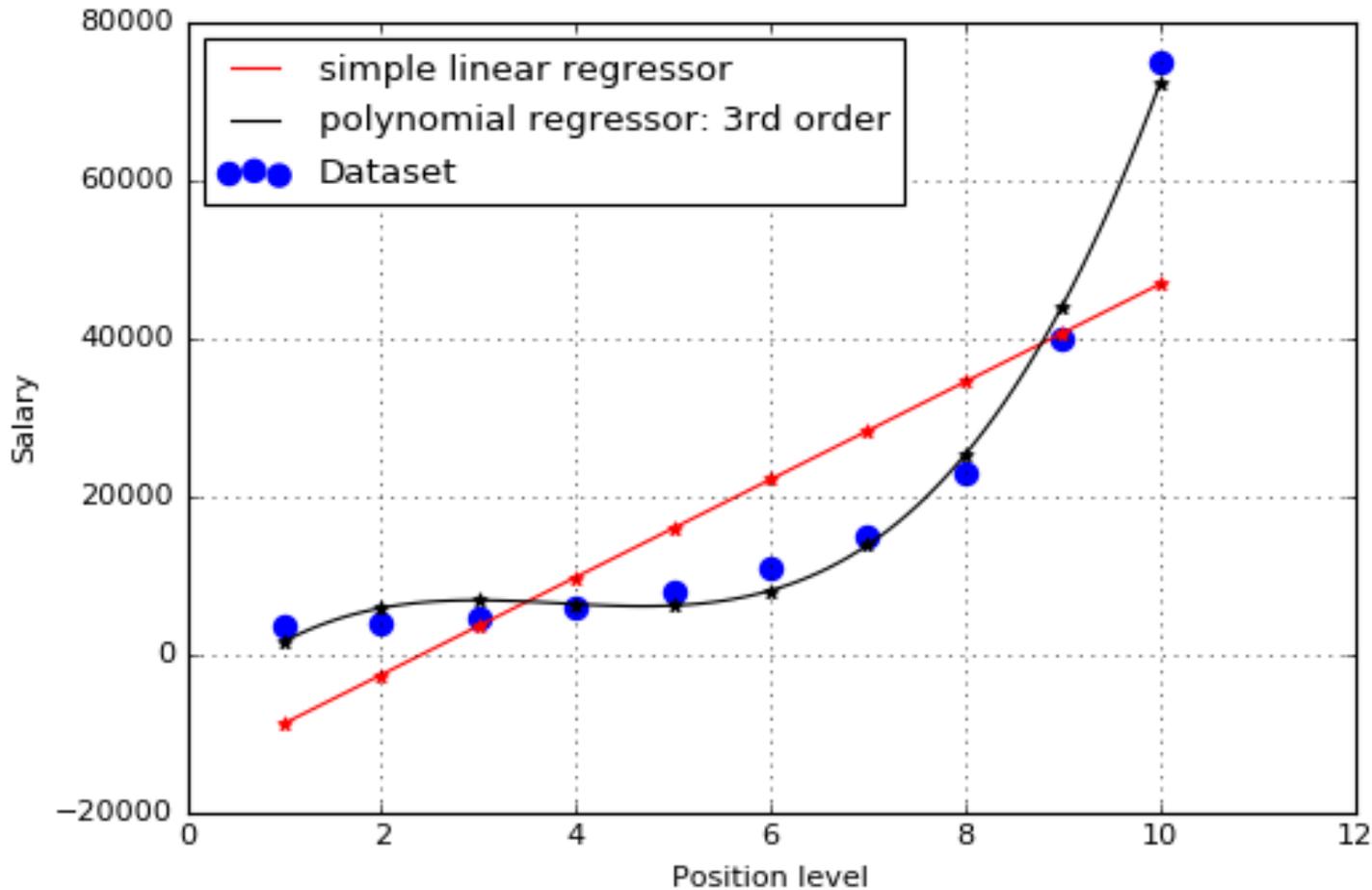
In fact, **the problem is to determine the coefficients.**

Data set: Position- salary

Position	Level	Salary
Business Analyst	1	3500
Junior Consultant	2	3900
Senior Consultant	3	4500
Manager	4	5800
Country Manager	5	8000
Region Manager	6	11000
Partner	7	15000
Senior Partner	8	23000
C-level	9	40000
CEO	10	75000

What is the best polynomial model Position (Level) – Salary?

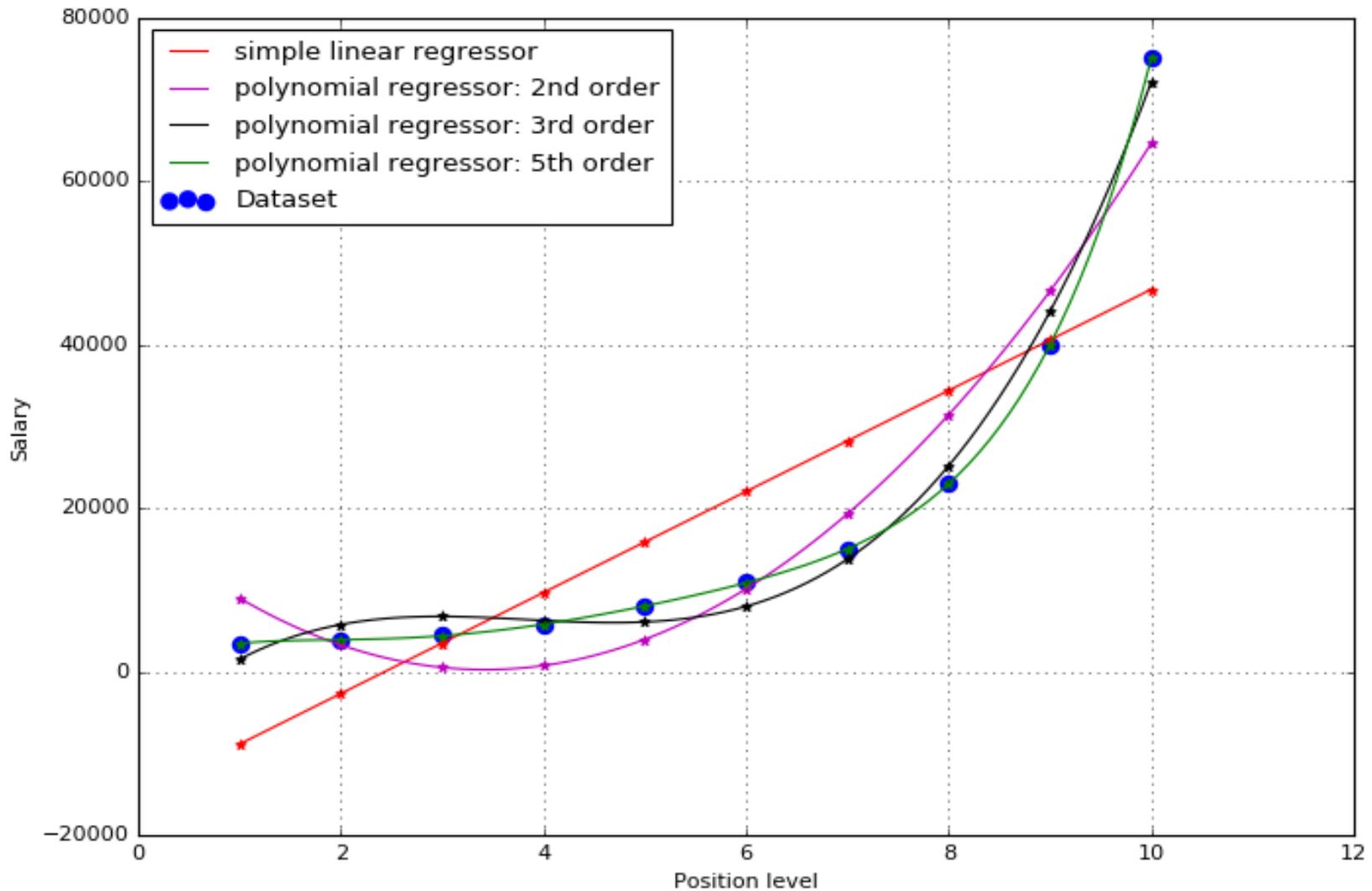
1st order vs 3rd order polynomial model



1st : $\text{Salary} = -15006 + 6178 \text{ Level}$

3rd : $\text{Salary} = -7747 + 12408 \text{ Level} - 3412 \text{ Level}^2 + 297 \text{ Level}^3$

Comparison – different polynomial model

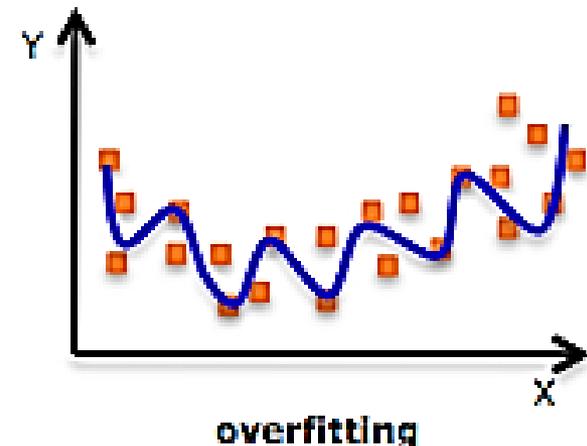
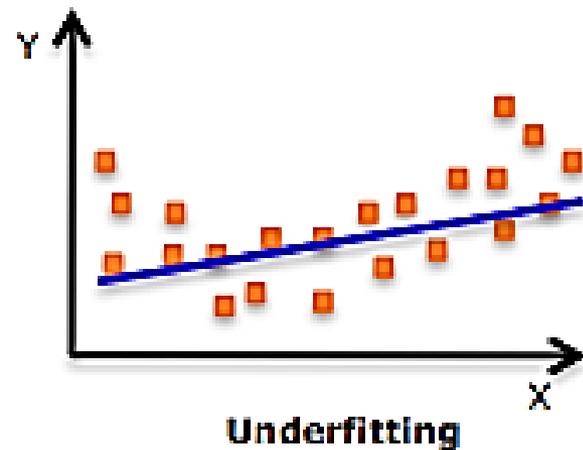


Order	1 st (linear)	2 nd	3 rd	5 th
R^2	0.6775504	0.9292263	0.9878567	0.9999947



Important aspects

- ❑ While there might be a temptation to fit a **higher degree polynomial** to get lower error, this can result in **over-fitting**.
- ❑ Always plot the relationships to see the fit and focus on making sure that the curve **fits the nature of the problem**
- ❑ Look out for curve towards the ends and see whether those shapes and trends **make sense**. Higher polynomials can end up producing weird results on extrapolation.



[7 Types of Regression Techniques you should know!,
<https://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/>]



Problem

For the dataset presented in the next plot, a linear regression model $\hat{y}_l = ax + b$, $b = 175$, $a = -8$, and a quadratic regression model $\hat{y}_q = a_1x + a_2x^2 + b$ $b = 470$; $a_1 = -56$; $a_2 = 1.9$, were developed.

- 0.5p** Which is the equation of the linear model? Plot the regression line.
- 0.5p** Which is the equation of the quadratic model? Plot the regression curve.
- 1p** Which is the predicted value of the dependent variable y for the independent variable $x = 12$ for both models. Plot this data points. Which of the two models is more accurate? Why?

