# Introduction to
# Data Science and Analytics

## Stephan Sorger
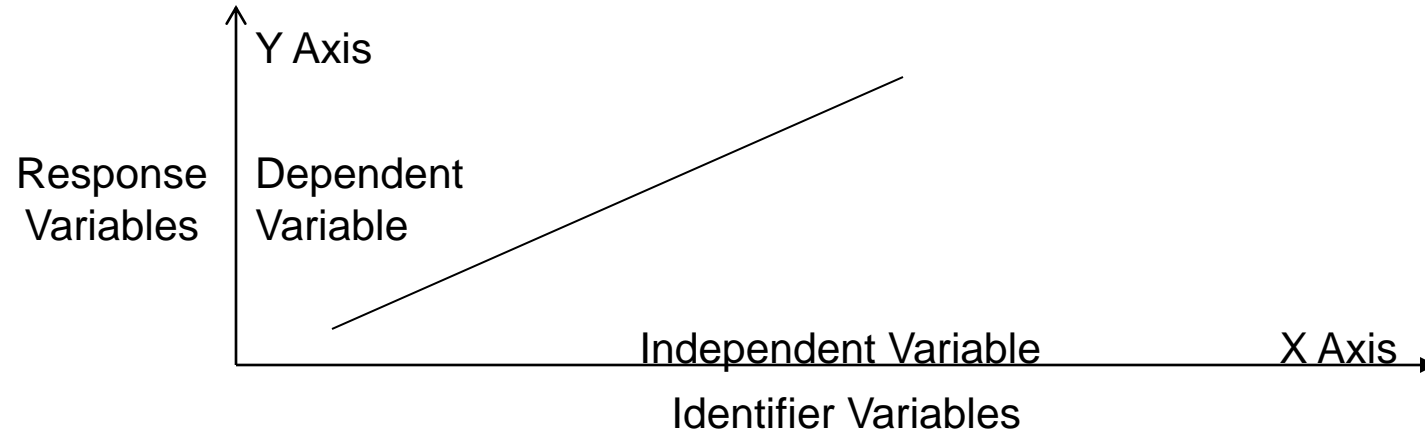### www.StephanSorger.com

# Unit 4. Excel Regression

Disclaimer:
• All images such as logos, photos, etc. used in this presentation are the property of their respective copyright owners and are used here for educational purposes only
• Some material adapted from: Sorger, "Marketing Analytics: Strategic Models and Metrics"
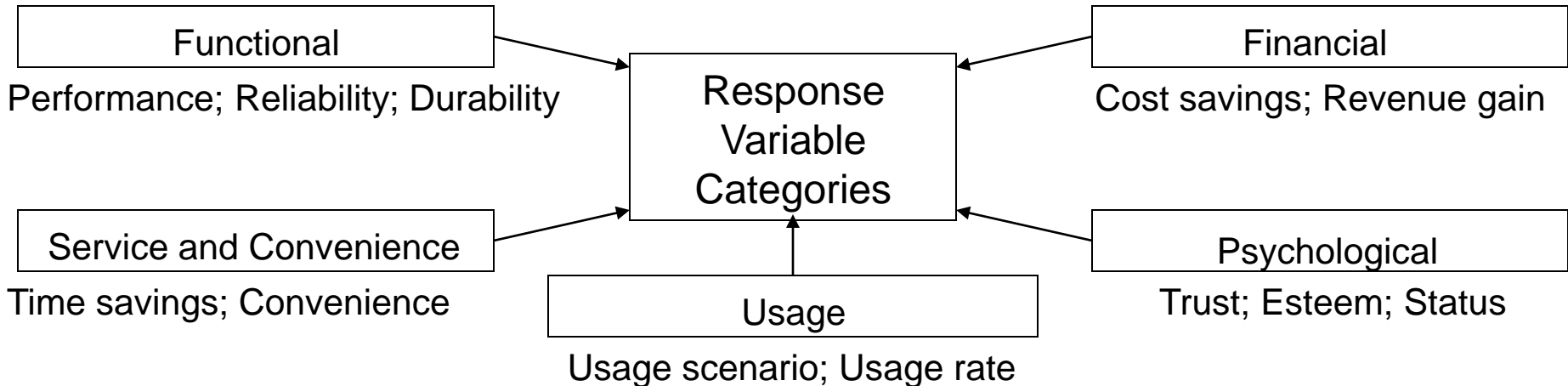
# Outline/ Learning Objectives

| Topic | Description |
|---|---|
| Background | The goal of regression analysis |
| Statistics | Basic statistics governing regression performance |
| Tests | F tests, t tests, p tests |
| Procedure | Executing regression analysis in Microsoft Excel |
| Multivariate | Executing cases with two or more independent variables |

# Regression Analysis
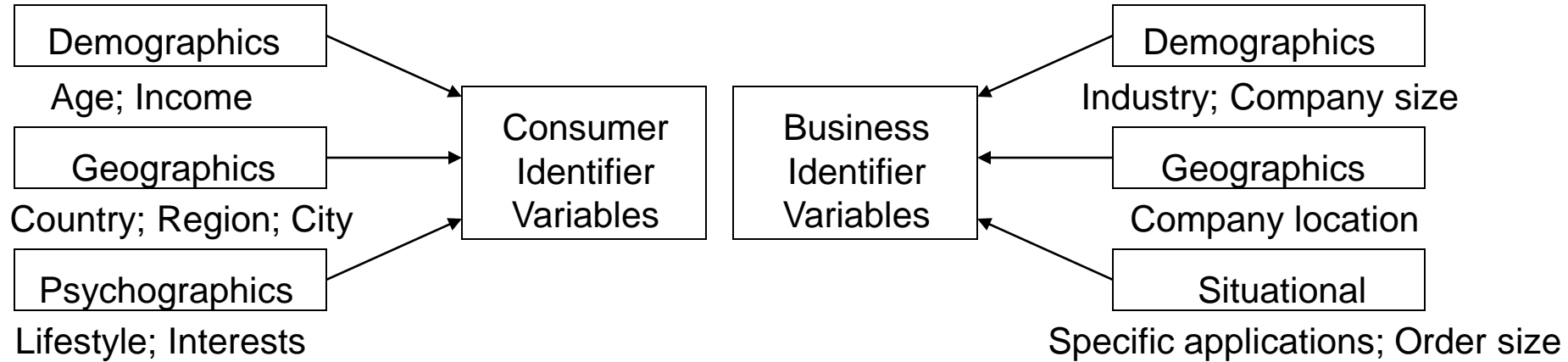
Goal is to establish the relationship between
Independent variables (the "inputs")
and dependent variables (also called response variables)

Y Axis

Response
Variables

Dependent
Variable

Independent Variable          X Axis

Identifier Variables

# Response (Dependent) Variable Categories

Functional
Performance; Reliability; Durability

Financial
Cost savings; Revenue gain

Response Variable Categories

Service and Convenience
Time savings; Convenience

Psychological
Trust; Esteem; Status

Usage
Usage scenario; Usage rate

# Independent (Input) Variables

Demographics

Age; Income

Geographics

Country; Region; City

Psychographics

Lifestyle; Interests

Consumer Identifier Variables

Business Identifier Variables

Demographics

Industry; Company size

Geographics

Company location

Situational

Specific applications; Order size

Many other independent variables possible:  See next slide

# Regression Example

**Scenario: Moving into a New Apartment**
(regular apt: not mansion; not rent control)

Reponse Variable: S.F. Monthly Rent Paid

Independent Variables:
(want to predict how much people will pay)
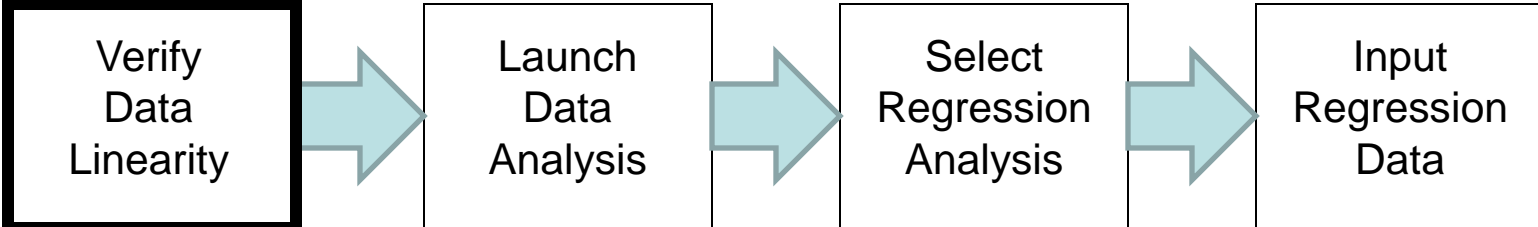
Demographics: Age
Demographics: Income
Geographics: Location of workplace
Psychographics: Status required
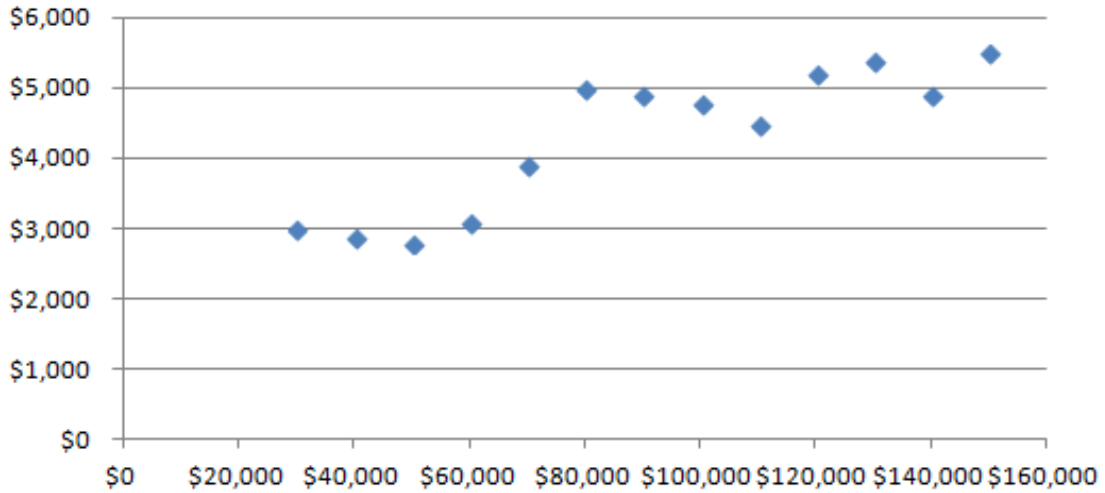Psychographics: Entertaining requirements
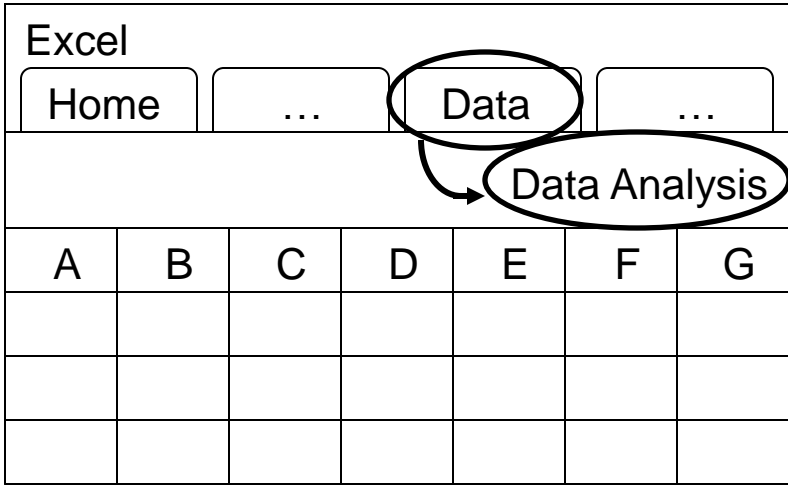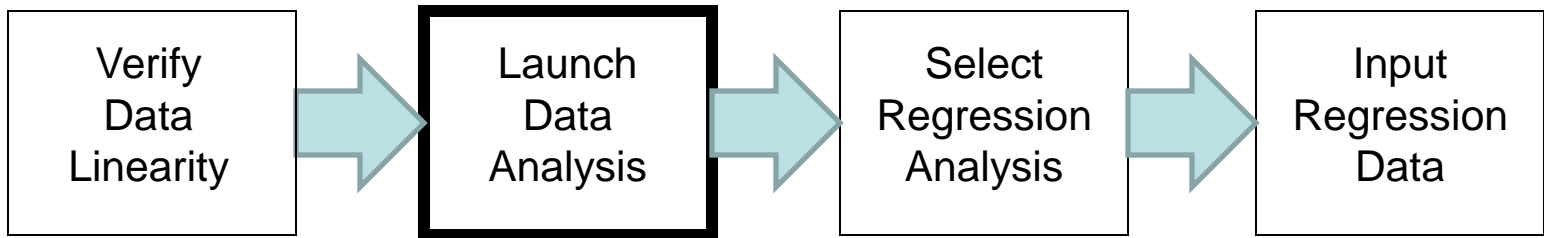
# Regression Analysis: Process

| Verify Data Linearity | → | Launch Data Analysis | → | Select Regression Analysis | → | Input Regression Data |
|---|---|---|---|---|---|---|

| Income | Rent |
|---|---|
| $30,000 | $3,000 |
| $40,000 | $2,900 |
| $50,000 | $2,800 |
| $60,000 | $3,100 |
| $70,000 | $3,900 |
| $80,000 | $5,000 |
| $90,000 | $4,900 |
| $100,000 | $4,800 |
| $110,000 | $4,500 |
| $120,000 | $5,200 |
| $130,000 | $5,400 |
| $140,000 | $4,900 |
| $150,000 | $5,500 |



**Rent vs. Income**

# Regression Analysis: Process

| Verify Data Linearity | → | **Launch Data Analysis** | → | Select Regression Analysis | → | Input Regression Data |
|---|---|---|---|---|---|---|

Excel

| Home | … | Data | … |
|---|---|---|---|

Data Analysis

| A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|
|  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |

# Regression Analysis: Process

```
┌─────────────┐      ┌─────────────┐    ┏━━━━━━━━━━━━━┓      ┌─────────────┐
│   Verify    │      │   Launch    │    ┃   Select    ┃      │    Input    │
│    Data     │ ──▶  │    Data     │──▶ ┃ Regression  ┃ ──▶  │ Regression  │
│  Linearity  │      │  Analysis   │    ┃  Analysis   ┃      │    Data     │
└─────────────┘      └─────────────┘    ┗━━━━━━━━━━━━━┛      └─────────────┘
```

Data Analysis

Analysis Tools

**Regression**                    OK

# Regression Analysis: Process

| Verify Data Linearity | → | Launch Data Analysis | → | Select Regression Analysis | → | **Input Regression Data** |
|---|---|---|---|---|---|---|

### Regression

Input Y Range [＿＿＿]    OK

Input X Range [＿＿＿]

[X] Labels

[ ] Constant is Zero

[X] Confidence Level: [95] %

X →

| Income | Rent |
|---|---|
| $30,000 | $3,000 |
| $40,000 | $2,900 |
| $50,000 | $2,800 |
| $60,000 | $3,100 |
| $70,000 | $3,900 |
| $80,000 | $5,000 |
| $90,000 | $4,900 |
| $100,000 | $4,800 |
| $110,000 | $4,500 |
| $120,000 | $5,200 |
| $130,000 | $5,400 |
| $140,000 | $4,900 |
| $150,000 | $5,500 |

← Y

# Excel Output



R-Square

F

P value
T stat
Standard Error
Coefficients

**SUMMARY OUTPUT**

| Regression Statistics | |
|---|---|
| Multiple R | 0.896671074 |
| R Square | 0.804019015 |
| Adjusted R Sq | 0.786202561 |
| Standard Error | 471.1613386 |
| Observations | 13 |

**ANOVA**

| | df | SS | MS | F | Significance F |
|---|---|---|---|---|---|
| Regression | 1 | 10018076.92 | 10018076.92 | 45.12789416 | 3.29676E-05 |
| Residual | 11 | 2441923.077 | 221993.007 | | |
| Total | 12 | 12460000 | | | |

| | Coefficients | Standard Error | t Stat | P-value | Lower 95% | Upper 95% | Lower 95.0% | Upper 95.0% |
|---|---|---|---|---|---|---|---|---|
| Intercept | 2188.461538 | 340.4048629 | 6.428996107 | 4.88126E-05 | 1439.235487 | 2937.68759 | 1439.235487 | 2937.68759 |
| Income | 0.023461538 | 0.003492443 | 6.717729837 | 3.29676E-05 | 0.015774641 | 0.031148436 | 0.015774641 | 0.031148436 |

# Regression Analysis: R-Squared

| Scenario | R-Squared |
|---|---|
| No Relationship | 0.0 |
| Social Science Studies | 0.3 |
| Marketing Research | 0.6 |
| Scientific Applications | 0.9 |
| Perfect Relationship | 1.0 |

R-Squared, the Coefficient of Determination
Also known as "Goodness of Fit",
from 0 (no fit) to 1 (perfect fit)

# Hypothesis Testing: t-Stat and P-value

| Statistic | Description |
|---|---|
| Standard Error | Estimate of standard deviation of the coefficient |
| t-Stat | Coefficient divided by the Standard Error |
| P-value | Probability of encountering equal t value in random data<br>P-value should be 5% or lower |

Hypothesis Testing: Test $H_0$ (null hypothesis)
Null hypothesis: No correlation between x and y

| | Coefficients | Standard Error | t Stat | P-value |
|---|---|---|---|---|
| Intercept | 2188.461538 | 340.4048629 | 6.428996107 | 4.88126E-05 |
| Income | 0.023461538 | 0.0349248 | 6.717729837 | 3.29676E-05 |

Less than 5% → OK

# Hypothesis Testing: F value

| Statistic | Description |
|---|---|
| F value | Tests overall significance of the regression model |
| $H_0$ | Tests null hypothesis that all regression coefficients = 0<br>Tests full model against a model with no variables |
| Significance F | Associated P value; Less than 0.05 to invalidate $H_0$ |

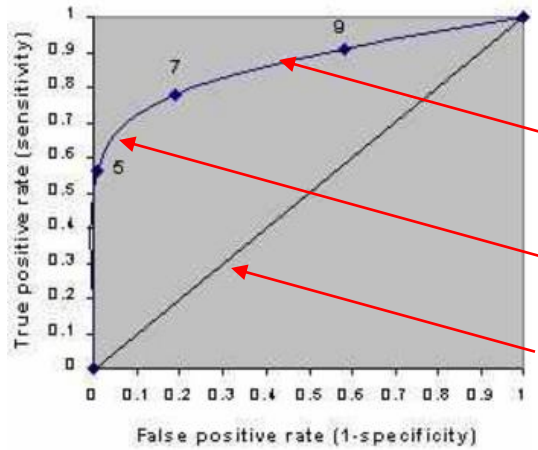Hypothesis Testing: Test $H_0$ (null hypothesis)
Null hypothesis: No correlation between x and y

| ANOVA | | | | | |
|---|---|---|---|---|---|
| | df | SS | MS | F | Significance F |
| Regression | 1 | 10018076.92 | 10018076.92 | 45.1278946 | 3.29676E-05 |
| Residual | 11 | 2441923.077 | 221993.007 | | |
| Total | 12 | 12460000 | | | |

Less than 5% → OK

# Regression Analysis: ROC Curves

| Topic | Description |
|---|---|
| ROC | Receiver Operating Characteristic<br>Plot of true positive rate against false positive rate at different cutpoints |
| History | Developed during World War II by RADAR engineers |
| Tradeoff | Shows tradeoffs, such as sensitivity and specificity for experiments |



Good→ close to top edge

Good→ close to left edge

Bad → close to diagonal

| Cutpoint | Sensitivity | Specificity |
|---|---|---|
| 5 | 0.56 | 0.01 |
| 7 | 0.78 | 0.19 |
| 9 | 0.91 | 0.58 |

| Cutpoint | True Pos. | False Pos. |
|---|---|---|
| 5 | 0.56 | 0.01 |
| 7 | 0.78 | 0.19 |
| 9 | 0.91 | 0.58 |

# Regression Analysis: ROC Curves

| Topic | Description |
|-------|-------------|
| Tests | Test predictive performance of model; how to select cutoffs |
| | At 95% level of confidence, we test $H_0$ at 5% (alpha = 5%) |
| True positive | Correctly identified; High income people rent expensive apts. |
| True negative | Correctly rejected; Low income people rent cheap apartments |
| False positive | Null hypothesis is true; No correlation (type I error) |
| | To address type I error, reduce alpha (in our case, 5%) |
| False negative | Failing to reject null hypothesis which is false (type II error) |
| | We thought model doesn't work, but it does |
| Tradeoff | As we decrease alpha from 5% to 1%... |
| | Type I error decreases, but Type II error increases (typical) |
| | Selecting cutoff a business decision; alpha = 5% usually good |

# Regression Analysis: Coefficients

|            | Coefficients |
|------------|--------------|
| Intercept  | 2188.461538  |
| Income     | 0.023461538  |

**Rent vs. Income**



Slope: 0.023 / 1

Y-intercept: 2,188

# Regression Analysis: Multivariate

| Period | Sales Level | Market Awareness | Number of Locations |
|--------|-------------|------------------|---------------------|
| Q1 2012 | $1.0 million | 80% | 5 |
| Q2 2012 | $1.1 million | 80% | 5 |
| Q3 2012 | $1.3 million | 85% | 6 |
| Q4 2012 | $1.2 million | 85% | 6 |
| Q1 2013 | $1.3 million | 85% | 7 |
| Q2 2013 | $1.5 million | 90% | 8 |
| Q3 2013 | $1.5 million | 90% | 8 |
| Q4 2013 | $1.4 million | 90% | 8 |

*What would happen if we opened 2 new stores, while holding awareness at 90%?*

# Regression Analysis: Multivariate

Y Range: Sales
X Range: Awareness & Locations

# Regression Analysis: Multivariate

R – squared: 0.92

Y-intercept = -1.44286

Coefficient for Awareness: 2.857143
at a P-value of 24.2% (not very good)

Coefficient for Locations: 0.042857
at a P-value of 56.2% (poor)

# Regression Analysis: Multivariate

| Output | Description | Values in Our Sales Example |
|---|---|---|
| R-Square | Goodness of fit of model to data | 0.93 |
| Intercept | Point where line crosses Y axis | -1.44 |
| Coefficient 1 | Coefficient for Market Awareness | 2.857 |
| Coefficient 2 | Coefficient for Number of Locations | 0.043 |

Sales = (Intercept) + (Coefficient 1) * (Market Awareness) + (Coefficient 2) * (Number of Locations)

= (- 1.44) + (2.857) * (Market Awareness) + (0.043) * (Number of Locations)

Example: Maintain brand awareness at 90%; Open two new retail stores (10 total)

= (-1.44) + (2.857) * (0.90) + (0.043) * (10) = $1.56 Million

# Outline/ Learning Objectives

| Topic | Description |
|---|---|
| Background | The goal of regression analysis |
| Statistics | Basic statistics governing regression performance |
| Tests | F tests, t tests, p tests |
| Procedure | Executing regression analysis in Microsoft Excel |
| Multivariate | Executing cases with two or more independent variables |