# IPM of Date Palm Insect Pests and Diseases
*Training Course*

**Statistical Designs and Analysis of IPM data of Date Palm Pests**
*(Simple and Multiple Regression)*

**Name: Khaled Al-Shamaa**

**Date: 28 February 2017**
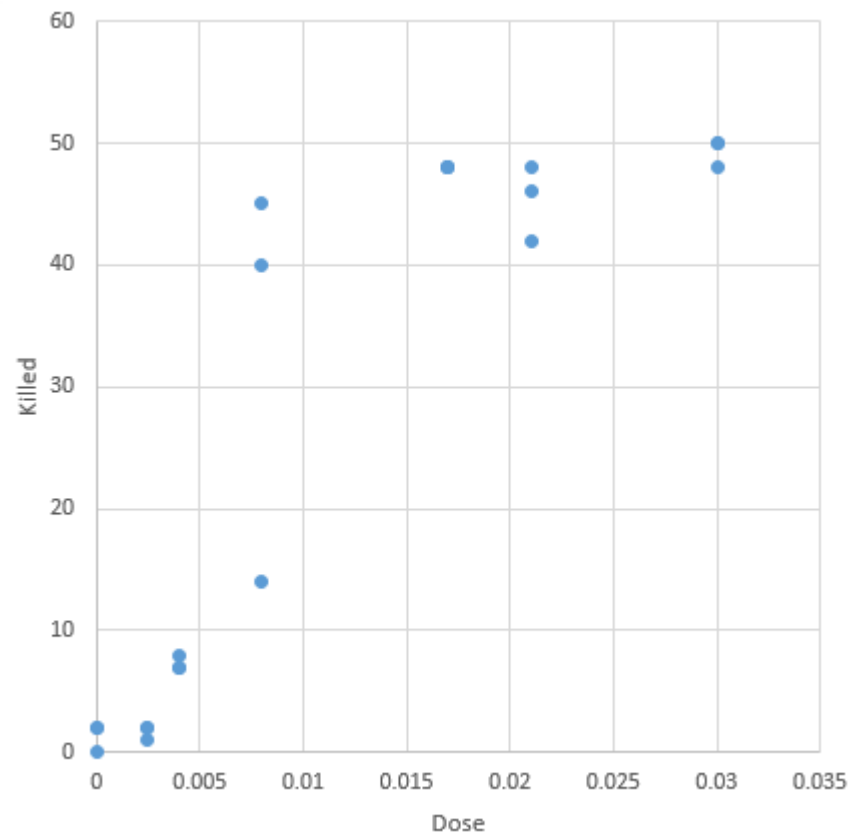
**Venue: Muscat, Oman**

**ICARDA**
Science for Better Livelihoods in Dry Areas

# Correlation

- Quantitative variables, linear relationship.

- Correlation does not imply causation.

- Correlation value vary from -1 to +1
  *-1 indicates perfect negative correlation, and +1 indicates perfect positive correlation. 0 means no correlation.*

- Correlation Significance
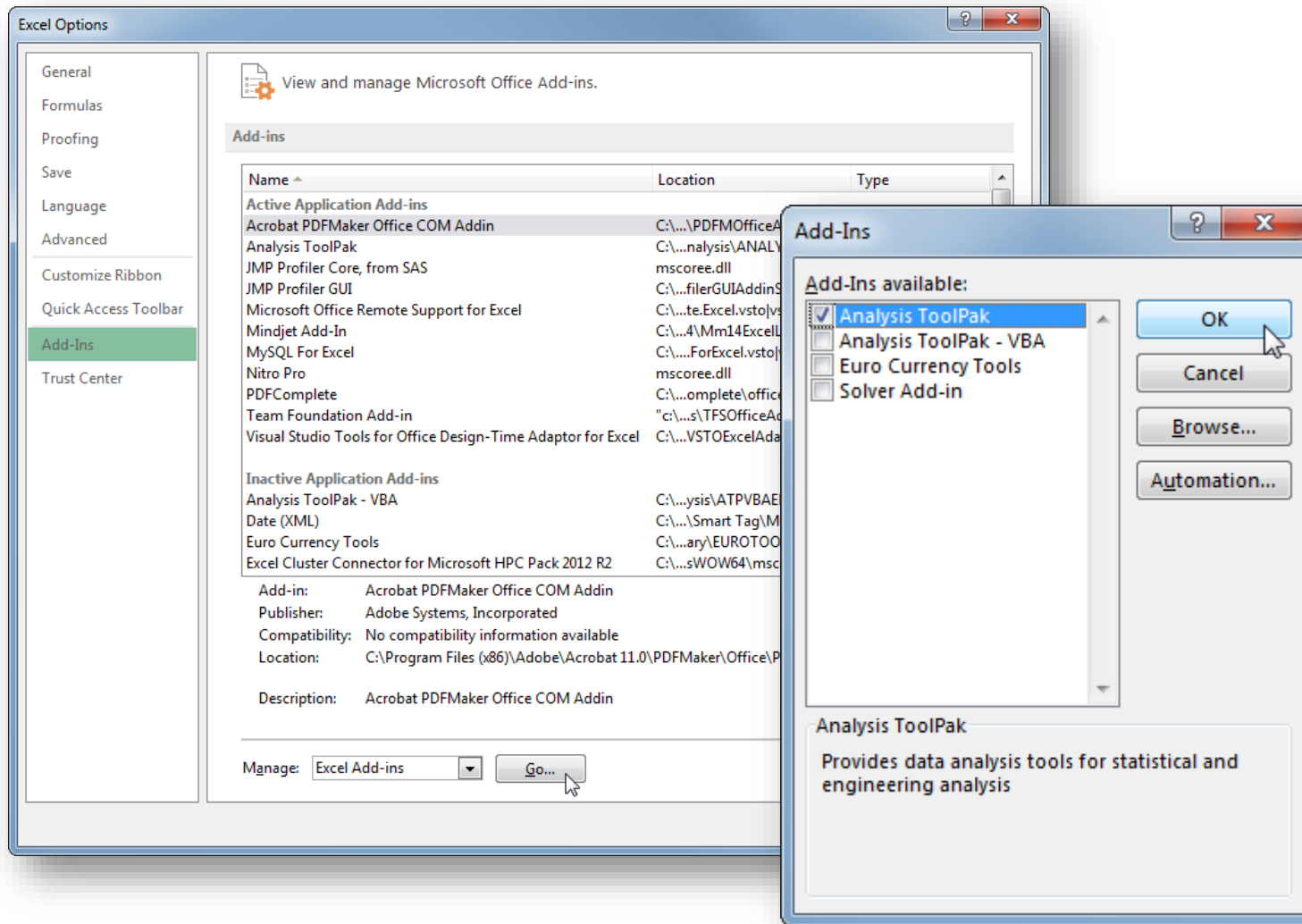  *depends on the correlation value and number of observations (test using t-test against 0).*

ICARDA
Science for Better Livelihoods in Dry Areas

# Excel - Correlation
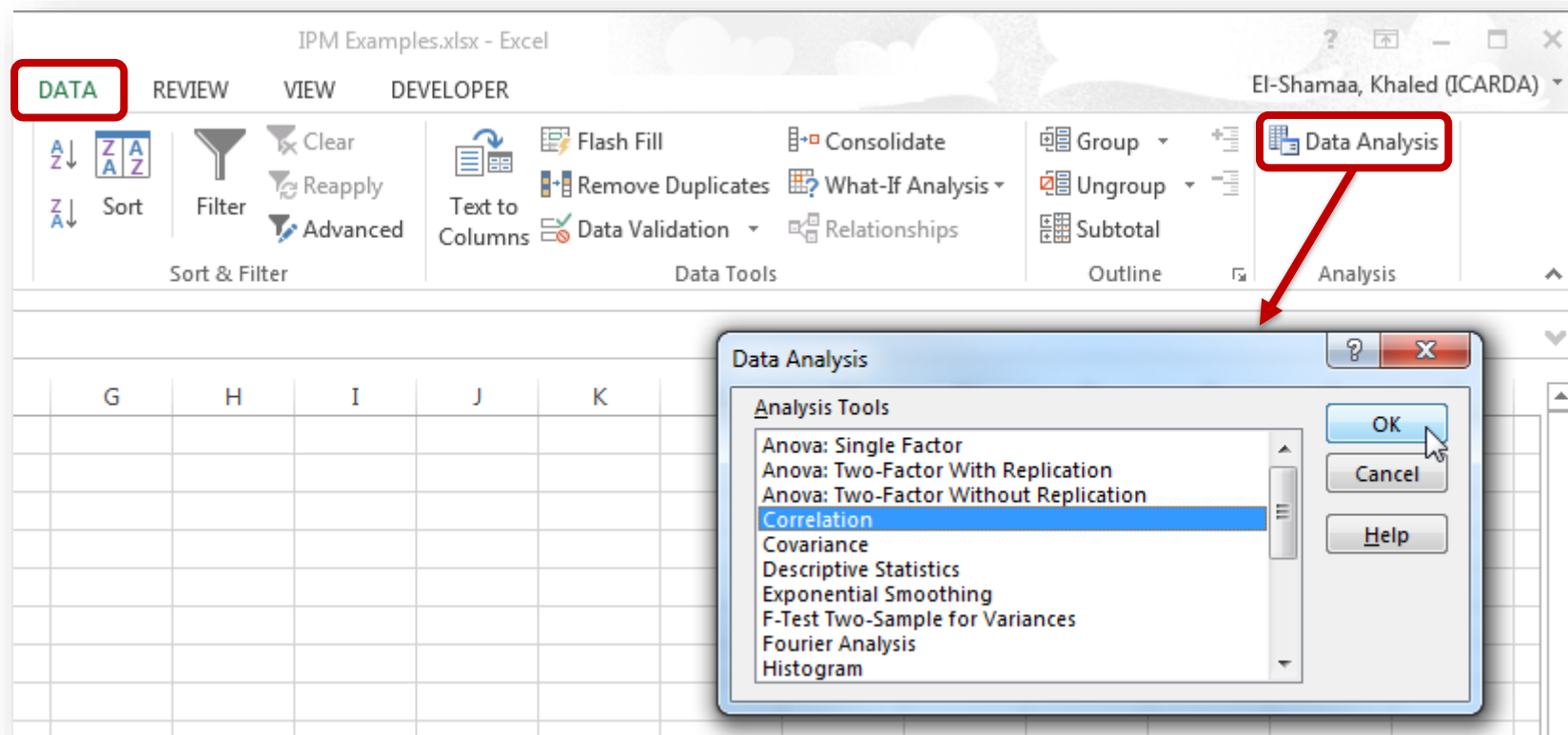


Treatment: Phosphine, PH3 (mg/l)
Subject: Storage Pests (*R. dominica)*

# Excel – Analysis ToolPak

# Excel – Analysis ToolPak *(continue)*

# GenStat - Correlation

# Correlation Abuse!



Correlation does not provide any information about the slope of the linear dependency

# Regression Analysis

- The goal of regression analysis is to use the data on some objects to predict values for another object.

- If X predicts Y it does not mean that X causes Y.

- Accurate prediction depends heavily on measuring the right variables.

# R – Squared (R²)

- R-squared is a statistical measure of how close the data are to the fitted regression line.

- It is the percentage of the response variable variation that is explained by a linear model.

- Yes! It is squared of Correlation R value.

$$R^2 = 1 - \frac{SS_{residuals}}{SS_{total}}$$

# Linear Regression Analysis

# Importance of Graphics!



All four sets are identical when examined using simple summary statistics (i.e. mean, variance, correlation, and regression), but vary considerably when graphed

# Linear Regression

- When the regression model contains one dependent variable and one independent variable, we call the approach simple linear regression.



Nonlinear Data



Transformed Nonlinear Data

ICARDA
Science for Better Livelihoods in Dry Areas

# Nonlinear Regression

- When there's one predictor variable but powers of the variable are included (e.g. $X^2$, $X^3$, etc.), we call it polynomial regression (e.g. quadratic or cubic regression).

- When there's more than one predictor variable (e.g. $X_1$, $X_2$, $X_3$, etc.), we call it multiple linear regression.

Essentially, all models are wrong, but some are useful



**HuffPost Model**

Our model of the polls suggests Clinton was **very likely leading**. (In >99% of simulations, Clinton led Trump.)



George E. P. Box

# Excel – Scatter Plot & Add Trendline

# Excel – Linear Regression, No Intercept

# Excel - Nonlinear Regression, Polynomial

# Excel (Analysis ToolPak) - Regression

# Excel (Analysis ToolPak) – Regression *(continue)*

SUMMARY OUTPUT

| *Regression Statistics* | |
|---|---|
| Multiple R | 0.87889039 |
| R Square | 0.77244831 |
| Adjusted R Square | 0.76047191 |
| Standard Error | 10.7075596 |
| Observations | 21 |

ANOVA

| | df | SS | MS | F | Significance F |
|---|---|---|---|---|---|
| Regression | 1 | 7394.75805 | 7394.75805 | 64.4975133 | 1.58222E-07 |
| Residual | 19 | 2178.384807 | 114.651832 | | |
| Total | 20 | 9573.142857 | | | |

| | Coefficients | Standard Error | t Stat | P-value | Lower 95% | Upper 95% |
|---|---|---|---|---|---|---|
| Intercept | 5.17107232 | 3.544048002 | 1.45908642 | 0.16087667 | -2.24670539 | 12.58885 |
| Dose | 1817.99143 | 226.370769 | 8.03103438 | 1.5822E-07 | 1344.191964 | 2291.79089 |

ICARDA
Science for Better Livelihoods in Dry Areas

# GenStat – Simple Linear Regression

# GenStat – Simple Linear Regression *(continue)*

```
25   "Simple Linear Regression"
26   MODEL Killed
27   TERMS Dose
28   FIT [PRINT=model,summary,estimates; CONSTANT=estimate;
```

## Regression analysis

Response variate:  Killed
Fitted terms:  Constant, Dose

## Summary of analysis

| Source | d.f. | s.s. | m.s. | v.r. | F pr. |
|---|---|---|---|---|---|
| Regression | 1 | 7395. | 7394.8 | 64.50 | <.001 |
| Residual | 19 | 2178. | 114.7 | | |
| Total | 20 | 9573. | 478.7 | | |

Percentage variance accounted for 76.0
Standard error of observations is estimated to be 10.7.

*Message: the following units have large standardized residuals.*

| Unit | Response | Residual |
|---|---|---|
| 11 | 45.0 | 2.43 |

*Message: the residuals do not appear to be random; for example, fitted values in the range 5.2 to 19.7 are consistently larger than observed values and fitted values in the range 19.7 to 43.3 are consistently smaller than observed values.*

*Message: the following units have high leverage.*

| Unit | Response | Leverage |
|---|---|---|
| 19 | 50.0 | 0.196 |
| 20 | 48.0 | 0.196 |
| 21 | 50.0 | 0.196 |

## Estimates of parameters

| Parameter | estimate | s.e. | t(19) | t pr. |
|---|---|---|---|---|
| Constant | 5.17 | 3.54 | 1.46 | 0.161 |
| Dose | 1818. | 226. | 8.03 | <.001 |

24

# GenStat – Regression, Standard Curves

```
40  "Logistic (s-shaped or inverse s-shaped curve)"
41  MODEL Killed
42  TERMS Dose
43  FITCURVE [PRINT=model,summary,estimates; CURVE=logistic; CONSTANT=estimate; FPROB=yes]\
44   Dose
```
------------------------------------------------------------------------------------------

## Nonlinear regression analysis

Response variate:  Killed
Explanatory:  Dose
Fitted Curve:  A + C/(1 + EXP(-B*(X - M)))

## Summary of analysis

| Source | d.f. | s.s. | m.s. | v.r. | F pr. |
|---|---|---|---|---|---|
| Regression | 3 | 8963.5 | 2987.85 | 83.32 | <.001 |
| Residual | 17 | 609.6 | 35.86 | | |
| Total | 20 | 9573.1 | 478.66 | | |

Percentage variance accounted for 92.5
Standard error of observations is estimated to be 5.99.

*Message: the following units have large standardized residuals.*
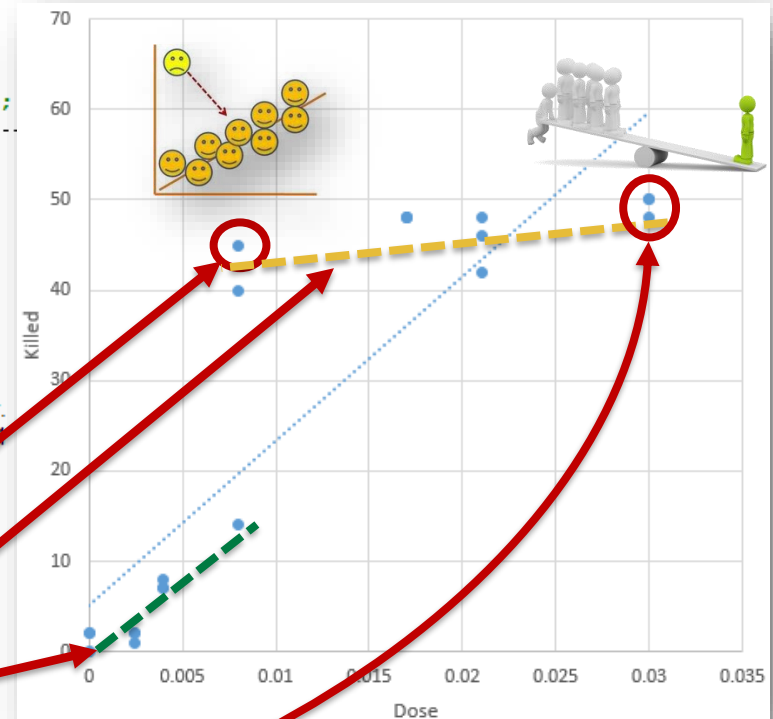
| Unit | Response | Residual |
|---|---|---|
| 10 | 14.00 | -3.90 |
| 11 | 45.00 | 2.44 |

## Estimates of parameters

| Parameter | estimate | s.e. |
|---|---|---|
| B | 660. | 235. |
| M | 0.006758 | 0.000668 |
| C | 47.36 | 4.44 |
| A | 0.20 | 3.90 |

Fitted and observed relationship

# Thank You

# Questions?

**Japanese attitude for work:**

*If one can do it, I can do it. If no one can do it, I must do it.*

**Middle Eastern attitude for work:**

*Wallahi… if one can do it, let him do it.*
*If no one can do it, ya-habibi how can I do it?*

ICARDA
Science for Better Livelihoods in Dry Areas

$$x_1, x_2, \ldots, x_n \qquad \sim N(\mu, \sigma^2)$$

$$\bar{x} = \frac{\sum x_i}{n}$$

$$Var(x) = \frac{\sum(x_i - \bar{x})^2}{n}$$

$$SD(x) = \sigma = \sqrt{Var(x)}$$

$$Z_i = \frac{x_i - \bar{x}}{SD(x)} \qquad \sim N(0, 1)$$

$$t = \frac{\bar{x} - \mu}{SD(x)/\sqrt{n}} \qquad \sim t(n-1)$$

$\tilde{x}$

$s$

ICARDA
Science for Better Livelihoods in Dry Areas

$$Cov(x, y) = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$Cor(x, y) = \frac{Cov(x, y)}{\sigma_x \, \sigma_y}$$

$$b_1 = Cor(x, y) \frac{SD(y)}{SD(x)}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$\hat{y} = b_0 + b_1.x$$

ICARDA
Science for Better Livelihoods in Dry Areas