

ARTIFICIAL INTELLIGENCE

Lab8. Correlation and regression

Caution! All the data are available in respective *csv* files!

- The following sample observations have been obtained by a chemical engineer investigating the relationship between weight of final product *Y* (in pounds) and volume of raw materials *X* (in gallons):

<i>X</i> (volume of raw materials)	14	23	9	17	10	22	5	12	6	16
<i>Y</i> (weight of final product)	68	105	40	79	81	95	31	72	45	93

- compute the sample covariance between the weight of final product and the volume of raw materials
138.49
- compute and interpret the sample correlation coefficient
0.895
- draw the data on a scatter plot
- determine regression line of the final yield of a product with respect to the volume of raw materials used for production; represent regression line on the graph with scatter plot
 $y = 22.4 + 3.62x$
- describe, how the final yield of a product will change if the volume of raw materials increases by 1 gallon
increase by 3.62 pounds
- compute, what will be the final yield of a product if the volume of raw materials is equal to 15 gallons;
77
- compute, what will be the final yield of a product if the volume of raw materials is equal to 20 gallons
95
- compute, how good is approximation by linear function;
80%
- test the hypothesis about significance of regression coefficient. Assume 5% of significance level. What is the interpretation of the decision?
 $F = 32.22$, reject H_0

- Arsenic is found in many ground-waters and some surface waters. Recent health effects research has prompted the Environmental Protection Agency to reduce allowable arsenic levels in drinking water so that many water systems are no longer compliant with standards. This has spurred interest in the development of methods to remove arsenic. The accompanying data on pH (*X*) and percentage arsenic removed (*Y*) by a particular process was collected:

<i>X</i>	7.01	7.11	7.12	7.24	7.94	7.94	8.04	8.05	8.07	8.90	8.94	8.95	8.97	8.98	9.85	9.86	9.86	9.87
<i>Y</i>	60	67	66	52	50	45	52	48	40	23	20	40	31	26	9	22	13	7

- compute the covariance between the pH of water and percentage arsenic removed
-18.32
- compute and interpret the correlation coefficient
-0.95
- draw the data on a scatter plot
- determine regression line of the percentage arsenic removed with respect to the pH of water; present the line on the graph with scatter plot
 $y = 190.27 - 18.03x$
- describe, how the percentage of arsenic removed will change if the pH of water increases by 1
decreases by 18.03 (%)
- compute, what will be the percentage of arsenic removed if the pH of water is equal to 7.5
55
- compute, what will be the percentage of arsenic removed if the pH of water is equal to 9
28
- compute, how good is approximation by linear function;
90%
- test the hypothesis about significance of regression coefficient. Assume 1% of significance level. What is the interpretation of the decision?
 $F = 149.7$, reject H_0

3. Variations in clay brick masonry weight have implications not only for structural and acoustical design but also for design of heating, ventilating, and air conditioning systems. The following table represents X – mortar air content (%) and Y – mortar dry density (lb/ft³):

X	5.7	6.8	9.6	10.0	10.7	12.6	14.4	15.0	15.3	16.2	17.8	18.7	19.7	20.6	25.0
Y	119.0	121.3	118.2	124.0	112.3	114.1	112.2	115.1	111.3	107.2	108.9	107.8	111.0	106.2	105.0

- compute the covariance between the mortar air content and mortar dry density -26.6
 - compute and interpret the correlation coefficient -0.87
 - draw the data on a scatter plot
 - determine regression line of the mortar dry density with respect to the mortar air content; present the line on the same graph as scatter plot $y = 126.25 - 0.92x$
 - describe, how the mortar dry density will change if the mortar air content increases by 1%
decreases by 0.92 (lb/ft³)
 - compute, what will be the mortar dry density if the mortar air content is equal to 11% 116
 - compute, what will be mortar dry density if the mortar air content is equal to 23%; 105
 - compute, how good is approximation by linear function; 75%
 - test the hypothesis about significance of regression coefficient. Assume 1% of significance level. What is the interpretation of the decision? $F = 39.51$, reject H_0
4. It was decided to test the processor temperature [°C] depending on the computer operation time [in h]. The results are presented in the table:

X (operation time)	1	4	7	8	10	13	18	22	25	28
Y (temperature)	35	45	59	77	62	50	39	47	54	71

- Determine and interpret covariance and correlation coefficient between temperature and the operation time.
- On the basis of the scatter plot, verify possible relationship between temperature and operation time.
- Determine polynomial regression for the temperature depending on operation time; try various degrees of the polynomial.
- Add a polynomial regression functions to an existing plot.
- Verify significance of regression; assume significance level 5%.
- Predict the temperature after 15 and 48 hours.
- Determine and interpret the coefficient of determination.

Caution! To determine square polynomial regression use function $\text{lm}(Y \sim X + I(X^2))$