

Worksheet

21) When implementing linear regression of some dependent variable y on the set of independent variables $\mathbf{x} = (x_1, \dots, x_r)$, where r is the number of predictors, which of the following statements will be true?

- a) $\beta_0, \beta_1, \dots, \beta_r$ are the regression coefficients.**
- b) Linear regression is about determining the best predicted weights by using the method of ordinary least squares.**
- c) E is the random interval**
- d) Both a and b**

Ans. d) Both a and b are true statements when implementing linear regression:

a) $\beta_0, \beta_1, \dots, \beta_r$ are indeed the regression coefficients, which represent the intercept and slopes of the linear relationship between the dependent variable y and each independent variable x_i .

b) Linear regression aims to determine the best predicted weights by using the method of ordinary least squares (OLS), which minimizes the sum of squared residuals between the observed and predicted values of y . In other words, OLS finds the line that best fits the data by minimizing the sum of the vertical distances between the data points and the line.

c) However, statement c) is not necessarily true. E typically refers to the residual error, which is the difference between the observed value of y and the predicted value of y . It is not a random interval, but rather a random variable that follows a normal distribution with mean 0 and constant variance.

22) What indicates that you have a perfect fit in linear regression?

- a) The value $R^2 < 1$, which corresponds to $SSR = 0$**
- b) The value $R^2 = 0$, which corresponds to $SSR = 1$**
- c) The value $R^2 > 0$, which corresponds to $SSR = 1$**
- d) The value $R^2 = 1$, which corresponds to $SSR = 0$**

Ans. d) The value $R^2 = 1$, which corresponds to SSR (sum of squared residuals) = 0, indicates that you have a perfect fit in linear regression.

R^2 , or the coefficient of determination, is a statistical measure that represents the proportion of variance in the dependent variable (y) that is explained by the independent variables (x). R^2 ranges from 0 to 1, where 0 indicates that the model does not explain any of the variability of the dependent variable, and 1 indicates that the model explains all of the variability of the dependent variable.

When $R^2 = 1$, it means that the model perfectly fits the data, as all of the variability in the dependent variable is explained by the independent variables. This corresponds to $SSR = 0$, because there are no residuals or errors left in the model to explain. Therefore, statement d) is correct.

23) In simple linear regression, the value of what shows the point where the estimated regression line crosses the y axis?

- a) Y
- b) B_0
- c) B_1
- d) F

Ans. The point where the estimated regression line crosses the y-axis is called the intercept or the constant term, denoted by b_0 in the equation of simple linear regression. Therefore, the answer is b) B_0 .

The slope of the regression line, denoted by b_1 , represents the change in the dependent variable (y) associated with a unit change in the independent variable (x). The value of F is typically used to test the overall significance of the regression model. Y is the dependent variable in the regression equation and represents the value we are trying to predict or explain.

24)

25) There are five basic steps when you're implementing linear regression: • a. Check the results of model fitting to know whether the model is satisfactory. • b. Provide data to work with, and eventually do appropriate transformations. • c. Apply the model for predictions. • d. Import the packages and classes that you need. • e. Create a regression model and fit it with existing data. However, those steps are currently listed in the wrong order. What's the correct order? a) e, c, a, b, d b) e, d, b, a, c c) d, e, c, b, a d) d, b, e, a, c

Ans. The correct order of the steps for implementing linear regression is: b) Provide data to work with, and eventually do appropriate transformations. d) Import the packages and classes that you need. e) Create a regression model and fit it with existing data. a) Check the results of model fitting to know whether the model is satisfactory. c) Apply the model for predictions.

So, the correct order is b, d, e, a, c.

26) Which of the following are optional parameters to LinearRegression in scikit-learn?

a) Fit

b) fit_intercept

c) normalize

d) copy_X

e) n_jobs

f) reshape

Ans. b) fit_intercept c) normalize d) copy_X e) n_jobs

The optional parameters for LinearRegression in scikit-learn are:

- fit_intercept: a boolean parameter indicating whether to calculate the intercept for this model. If set to False, no intercept will be used in calculations.
- normalize: a boolean parameter indicating whether to normalize the features before fitting the model.
- copy_X: a boolean parameter indicating whether to copy the input X array before fitting the model.
- n_jobs: an integer parameter indicating the number of CPU cores to use for calculations. If set to -1, all available cores will be used.

Fit and reshape are not optional parameters for LinearRegression in scikit-learn. Fit is the method used to train the model, and reshape is a NumPy method used to change the shape of an array.

27) While working with scikit-learn, in which type of regression do you need to transform the array of inputs to include nonlinear terms such as x^2 ?

a) Multiple linear regression

b) Simple linear regression

c) Polynomial regression

Ans. c) Polynomial regression.

In polynomial regression, the input array is transformed to include nonlinear terms such as x^2 , x^3 , etc. This allows the model to capture nonlinear relationships between the input features and the target variable. Scikit-learn provides a PolynomialFeatures transformer that can be used to create polynomial features. The transformed feature matrix can then be used with linear regression models to fit a polynomial regression.

28) You should choose statsmodels over scikit-learn when:

A) You want graphical representations of your data.

b) You're working with nonlinear terms.

c) You need more detailed results.

d) You need to include optional parameters.

Ans. The choice between statsmodels and scikit-learn depends on the specific tasks and requirements of your project. Here are some guidelines to help you decide:

A) If you want graphical representations of your data, statsmodels is a better choice. Statsmodels provides extensive support for visualization and plotting, which can be helpful for exploring and interpreting your data.

B) If you're working with nonlinear terms, statsmodels is a better choice. Statsmodels has a wide range of regression models that can handle nonlinear relationships between variables, including polynomial regression, spline regression, and generalized additive models.

C) If you need more detailed results, statsmodels is a better choice. Statsmodels provides detailed statistical analysis of your models, including hypothesis testing, confidence intervals, and p-values. This can be especially useful if you need to report your results in a publication or presentation.

D) If you need to include optional parameters, scikit-learn is a better choice. Scikit-learn provides a wide range of machine learning models, and many of these models have a large number of optional parameters that can be tuned to improve performance. Scikit-learn also provides tools for feature selection, preprocessing, and model selection, which can be helpful for building more complex machine learning pipelines.

Overall, both statsmodels and scikit-learn are powerful libraries for statistical analysis and machine learning. The choice between them depends on the specific requirements of your project.

29) _____ is a fundamental package for scientific computing with Python. It offers comprehensive mathematical functions, random number generators, linear algebra routines, Fourier transforms, and more. It provides a high-level syntax that makes it accessible and productive.

- a) Pandas
- b) Numpy
- c) Statsmodel
- d) scipy

Ans. Numpy

30) _____ is a Python data visualization library based on Matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics that allow you to explore and understand your data. It integrates closely with pandas data structures.

- a) Bokeh
- b) Seaborn
- c) Matplotlib
- d) Dash

Ans. Seaborn