

CertNexus Certified Data Science Practitioner (CDSP) Exam DSP-110

Exam Information

Exam Objective Statement:

The exam will certify that the successful candidate has the knowledge, skills, and abilities required to answer questions by collecting, wrangling, and exploring data sets, applying statistical models and artificial-intelligence algorithms, to extract and communicate knowledge and insights.

Candidate Eligibility

The Certified Data Science Practitioner exam requires no application fee, supporting documentation, or other eligibility verification measures for you to be eligible to take the exam. An exam voucher will come bundled with your training program or can be purchased separately [here](#). Once purchased, you will receive more information about how to register for and schedule your exam through Pearson Vue. You can also purchase a voucher directly through Pearson Vue. Once you have obtained your voucher information, you can register for an exam time [here](#). By registering, you agree to our Candidate Agreement included [here](#).

Exam Prerequisites

There are no formal prerequisites to register for and schedule an exam. Successful candidates will possess the knowledge, skills, and abilities as identified in the domain objectives in this blueprint. It is also strongly recommended that candidates possess the following knowledge, skills and abilities:

- A working level knowledge of programming languages such as Python and R
- Proficiency with a querying language
- Strong communication skills
- Proficiency with statistics and linear algebra
- Demonstrate responsibility based upon ethical implications when sharing data sources
- Familiarity with data visualization

You can obtain this level of skill and knowledge by taking the following courseware, which is available through training providers located around the world, or by attending an equivalent third-party training program:

- Introduction to Programming with Python
- Python Programming Advanced
- Using Data Science Tools in Python
- Data Wrangling with Python
- Applied Data Science with Python and Jupyter
- Big Data Analysis with Python
- Certified Ethical Emerging Technologist

Exam Specifications

Number of Items: 100, of which 75 count toward your score

Passing Score: To be determined

Duration: 120 minutes (**Note:** Published exam times include the 10 minutes you are allotted for reading and signing the Candidate Agreement and reviewing exam instructions.)

Exam Options: Online through Pearson OnVUE or in person at Pearson VUE test centers.

Item Formats: Multiple Choice/Single Response

Exam Description

Target Candidate:

The Certified Data Science Practitioner exam is designed for professionals across different industries seeking to demonstrate the ability to gain insights and build predictive models from data.

The Certified Data Science Practitioner exam will test them on the following domains with the following weightings:

| Domain | # of Items |
|--|-------------------|
| 1.0 Defining the question to be addressed through the application of data science | 6 |
| 2.0 Extracting, Transforming, and Loading Data | 16 |
| 3.0 Performing exploratory data analysis | 23 |
| 4.0 Building models | 18 |
| 5.0 Testing models | 6 |
| 6.0 Communicating findings | 6 |
| Total | 75 |

The information that follows is meant to help you prepare for your certification exam. This information does not represent an exhaustive list of all the concepts and skills that you may be tested on during your exam. The exam domains, identified previously and included in the objectives listing, represent the large content areas covered in the exam. The objectives within those domains represent the specific tasks associated with the job role(s) being tested. The information beyond the domains and objectives is meant to provide examples of the types of concepts, tools, skills, and abilities that relate to the corresponding domains and objectives. All of this information represents the industry-expert analysis of the job role(s) related to the certification and does not necessarily correlate one-to-one with the content covered in your training program or on your exam. We strongly recommend that you independently study to familiarize yourself with any concept identified here that was not explicitly covered in your training program or products.

Objectives

Domain 1.0 Defining the question to be addressed through the application of data science

Objective 1.1 Identify the project scope.

- Identify project specifications, including objectives (metrics/KPIs) and stakeholder requirements
- Identify mandatory deliverables, optional deliverables
 - Determine project timeline
- Identify project limitations (time, technical, resource, data, risks)

Objective 1.2 Understand stakeholder challenges.

- Understand stakeholder terminology
 - Milestone
 - POC (Proof of concept)
 - MVP (Minimal Viable Product)
- Become aware of data privacy, security, and governance policies
 - GDPR
 - HIPAA
 - California Privacy Act
- Obtain permission/access to data

Objective 1.3 Classify a question into a known data science problem.

- Access references
 - Optimization problem
 - Forecasting problem
 - Prediction problem
 - Classification problem
 - Segmentation/Clustering problem
- Identify data sources and type

- Image
- Text
- Numerical
- Categorical
- Select modeling type
 - Regression
 - Classification
 - Forecasting
 - Clustering
 - Optimization
 - Recommender systems

Domain 2.0 Extracting, Transforming, and Loading Data

Objective 2.1 Gather relevant data sets.

- Read data
 - Write a query for a SQL database
 - Write a query for a NoSQL database
 - Load CSV files to dataframes
 - Read data from cloud storage solutions
 - AWS S3
 - Google Storage Buckets
 - Azure Data Lake
- Research third-party data availability
 - Demographic data
 - Bloomberg
- Collect open-source data
 - Use APIs to collect data
 - Scrape the Web

Objective 2.2 Clean data sets.

- Identify and eliminate irregularities in data
 - Nulls
 - Duplicates
 - Corrupt values
- Parse the data
- Check for corrupted data
- Correct the data format for storing/querying purposes
- Deduplicate data

Objective 2.3 Merge data sets.

- Join data from different sources
 - Make sure a common key exists in all datasets
 - Unique identifiers

Objective 2.4 Apply problem-specific transformations to data sets.

- Apply word embeddings
 - Word2vec
 - TF-IDF
 - Glove
- Generate latent representations for image data

Objective 2.5 Load data

- Load into DB
- Load into dataframe
- Export to CSV files
- Load into visualization tool
- Make an endpoint

Domain 3.0 Performing exploratory data analysis

Objective 3.1 Examine data

- Generate summary statistics
- Examine feature types
- Visualize distributions
- Identify outliers
- Find correlations
- Identify target feature(s)

Objective 3.2 Preprocess data

- Identify missing values
- Make decisions about missing values (e.g., imputing method, record removal)
- Normalize, standardize, or scale data

Objective 3.3 Carry out feature engineering

- Apply encoding to categorical data
 - One hot encoding
 - Target encoding
 - Label Encoding or Ordinal Encoding
 - Dummy Encoding
 - Effect Encoding
 - Binary Encoding
 - BaseN Encoding
 - Hash Encoding
- Assign feature values to bins or groups
 - Combining categories of data
 - Binning into a category
- Split features
 - Text manipulation

- Split
 - Trim
 - Reverse
- Take year from a date
- Split names
- Extract year from title
- Convert dates to useful features
- Apply feature reduction methods
 - PCA
 - t-SNE
 - Random Forest
 - Backward Feature Elimination,
 - Forward Feature Selection
 - Factor Analysis
 - Missing Value Ratio
 - Low Variance Filter
 - High Correlation Filter

Domain 4.0 Building models

Objective 4.1 Prepare data sets for modeling.

- Decide proportion of data set to use for training, testing, and (if applicable) validation
- Split data to train, test, and (if applicable) validation sets

Objective 4.2 Build training models

- Define algorithms to try
 - Regression,
 - Linear regression,
 - Random forest
 - XGBoost
 - Classification:
 - Logistic regression
 - Random forest classification
 - XGBoost classifier
 - Naive Bayes
 - Forecasting
 - ARIMA
 - Clustering:
 - K-means
 - Latent class
 - Hierarchical clustering
- Train model
- Tune hyperparameters, if applicable

- Grid Search
- Hyper opt

Objective 4.3 Evaluate models

- Define evaluation metric
- Compare model outputs
 - Confusion matrix
 - Learning curve
- Select best-performing model
- Store model for operational use
 - MLFlow
 - Kubeflow

Domain 5.0 Testing models

Objective 5.1 Test hypotheses

- Design A/B tests
 - Experimental Design
 - Design Use Cases
 - Test Creation
 - Statistics
- Define success criteria for test
- Evaluate test results

Objective 5.2 Test pipelines

- Put model into production
 - AWS SageMaker,
 - Azure ML
 - Docker
 - Kubernetes
- Ensure model works operationally
- Monitor pipeline for performance of model over time
 - MLFlow
 - Kubeflow
 - Data Dog

Domain 6.0 Communicating findings

Objective 6.1 Report findings

- Implement model in a basic web application for demonstration (POC implementation)
 - Web frameworks (Flask, Django)
 - Basic HTML
 - CSS
- Derive insights from findings

- Identify features that drive outcomes (e.g., explainability; variable importance plot)
- Show model results
- Generate lift or gain chart

Recertification Requirements

The Certified Data Science Practitioner certification is valid for 3 years from the date that it is initially granted. In order to maintain a continuously valid certification, candidates can recertify by retaking the most recent version of the exam before their certification expires.

Certified Data Science Professional Acronyms

| Acronym | Expanded Form |
|----------------|---|
| AI | Artificial Intelligence |
| ANN | Artificial Neural Network |
| AUC | Area Under the Curve |
| CNN | Convolutional Neural Network |
| ETL | Extract, Transform, and Load |
| GAN | Generative Adversarial Network |
| GDPR | General Data Protection Regulation |
| HIPAA | Health Insurance Portability and Accountability Act |
| KNN | K-Nearest Neighbors |
| KPI | Key Performance Indicator |
| LSTM | Long Short-Term Memory |
| ML | Machine Learning |
| NLP | Natural Language Processing |
| NLU | Natural Language Understanding |
| PCA | Principal Component Analysis |
| RBF | Radial Basis Function Kernel |
| REST API | Representational State Transfer Application Programming Interface |
| ROC | Receiver Operating Characteristic |
| RNN | Recurrent Neural Network |
| SVM | Support Vector Machines |
| SQL | Structured Query Language |
| TF-IDF | Term Frequency–Inverse Document Frequency |
| CSV | Comma-separated Values |