

	MON	TUE	WED	THU	FRI	ENVIRONMENT
Week 1 Advanced Python, Data Engineering & Big Data Fundamentals	Project 1 Python Advanced - Collections & Error Handling Collections <ul style="list-style-type: none"> Counters Named Tuples Ordered Dicts Default Dict ChainMap Advanced Exception Handling <ul style="list-style-type: none"> Exception hierarchy and custom exceptions Context managers (with statement) Exception chaining and traceback Logging & Debugging <ul style="list-style-type: none"> Logging module: handlers, formatters, levels 	Project 1 Python Advanced - Data Visualization & Testing NumPy & Pandas Advanced <ul style="list-style-type: none"> NumPy: array operations, broadcasting, vectorization Pandas: advanced groupby, pivot tables, merging Performance optimization techniques Memory management with large datasets Matplotlib & Data Visualization <ul style="list-style-type: none"> Creating professional visualizations Subplots, multiple axes, and layouts Customization and styling Interactive visualizations basics Unit Testing with Pytest <ul style="list-style-type: none"> Test-driven development (TDD) principles Pytest fixtures and parameterization Mocking and patching Code coverage analysis 	Project 1 Advanced SQL & Database Optimization Advanced Queries - Window Functions & CTEs <ul style="list-style-type: none"> Window functions: ROW_NUMBER, RANK, DENSE_RANK, LEAD, LAG Running totals and moving averages Common Table Expressions (CTEs) and recursive CTEs Query structuring best practices 	Project 1 Big Data & Cloud Fundamentals Big Data Architecture <ul style="list-style-type: none"> Big Data characteristics: 4Vs (Volume, Velocity, Variety, Veracity) Big Data components and ecosystem Distributed computing concepts CAP theorem and eventual consistency Data Lifecycle Stages 	Project 1 Data Warehousing & Modeling File Formats & Data Types <ul style="list-style-type: none"> Structured vs Semi-structured vs Unstructured data File formats: CSV, JSON, XML, Avro, Parquet, ORC Columnar vs Row-based storage Compression algorithms and strategies Schema evolution considerations 	Data Warehousing Concepts Dimensional Modeling <ul style="list-style-type: none"> Conceptual vs Logical vs Physical modeling Star Schema design Snowflake Schema design Fact tables: additive, semi-additive, non-additive measures Dimension tables: conformed dimensions

	MON	TUE	WED	THU	FRI	ENVIRONMENT
			<ul style="list-style-type: none"> Views, Materialized Views Stored Procedures and User Defined Functions <p>ACID Properties & Transaction Management</p> <ul style="list-style-type: none"> Transaction isolation levels Deadlock detection and prevention Transaction commit, rollback patterns CRUD operations optimization 	Synapse, Databricks <ul style="list-style-type: none"> Pricing models and cost optimization 	cleansing strategies <ul style="list-style-type: none"> BigQuery overview: datasets, tables, queries 	
Week 2 PySpark Deep Dive	Project 2 Spark Ecosystem & Architecture Big Data Processing & Spark Evolution <ul style="list-style-type: none"> Hadoop MapReduce limitations Spark vs Hadoop: performance comparison Spark ecosystem: Core, SQL, Streaming, MLlib, GraphX Use cases for each Spark component Spark Architecture Deep Dive <ul style="list-style-type: none"> Driver, Executor, Cluster Manager roles Spark application lifecycle Job, Stage, Task hierarchy Shuffle operations 	Project 2 RDD Advanced Operations RDD Transformations & Actions <ul style="list-style-type: none"> map, flatMap, filter, distinct reduce, fold, aggregate collect, count, first, take, foreach saveAsTextFile and output operations Key-Value Pair	Project 2 DataFrames & Spark SQL Creating & Managing DataFrames <ul style="list-style-type: none"> Creating from: RDDs, JSON, CSV, Parquet, databases Reading with options: header, delimiter, schema Column data types and casting DataFrame 	Project 2 Advanced DataFrame Operations Data Aggregation <ul style="list-style-type: none"> RDDs vs DataFrames vs Datasets Benefits of DataFrames: Catalyst optimizer, Tungsten execution Spark Session API Schema definition: StructType, StructField Schema inference vs explicit schema Window Functions <ul style="list-style-type: none"> agg() with multiple aggregations pivot and unpivot operations Window specification: partitionBy, orderBy Ranking functions: row_number, rank, dense_rank Analytical functions: lead, lag, first, last Aggregate functions with windows 	Project 2 Data I/O & Performance Tuning Data Loading Patterns <ul style="list-style-type: none"> GroupBy operations: groupBy, rollup, cube Built-in aggregation functions: sum, avg, count, min, max agg() with multiple aggregations pivot and unpivot operations Window specification: partitionBy, orderBy Ranking functions: row_number, rank, dense_rank Analytical functions: lead, lag, first, last Aggregate functions with windows Window Functions <ul style="list-style-type: none"> Window specification: partitionBy, orderBy Ranking functions: row_number, rank, dense_rank Analytical functions: lead, lag, first, last Aggregate functions with windows External Data Sources <ul style="list-style-type: none"> GroupBy operations: groupBy, rollup, cube Built-in aggregation functions: sum, avg, count, min, max agg() with multiple aggregations pivot and unpivot operations Window specification: partitionBy, orderBy Ranking functions: row_number, rank, dense_rank Analytical functions: lead, lag, first, last Aggregate functions with windows 	

	MON	TUE	WED	THU	FRI	ENVIRONMENT
	<p>and impact</p> <ul style="list-style-type: none"> • Memory management: storage vs execution <p>Spark Environment Setup</p> <ul style="list-style-type: none"> • Local vs Cluster mode configuration • Spark Context vs Spark Session • Configuration parameters and tuning • Spark submit options • Deployment modes: client vs cluster <p>Execution Model & DAG</p> <ul style="list-style-type: none"> • Lazy evaluation concepts • Transformation vs Action operations • DAG (Directed Acyclic Graph) visualization • Lineage and fault tolerance 	<p>Operations</p> <ul style="list-style-type: none"> • Creating pair RDDs • groupByKey vs reduceByKey vs combineByKey • aggregateByKey and foldByKey • mapValues and flatMapValues • Join operations: join, leftOuterJoin, rightOuterJoin • cogroup operations <p>RDD Persistence & Shared Variables</p> <ul style="list-style-type: none"> • Storage levels: MEMORY_ONLY, MEMORY_AND_DISK, etc. • Caching strategies and when to use • Broadcast variables: use cases and implementation • Accumulators: built-in and custom • Partitioning: HashPartitioner, RangePartitioner, custom 	<p>schema operations</p> <p>DataFrame Operations</p> <ul style="list-style-type: none"> • select, selectExpr, drop, withColumn, withColumnRenamed • filter, where, distinct • orderBy, sort • limit, sample • Column expressions and functions <p>Spark SQL Integration</p> <ul style="list-style-type: none"> • Creating temporary views: createOrReplaceTempView • Creating global views • SQL queries on DataFrames • Advanced SQL: CTEs, subqueries, window functions • Catalog API 	<ul style="list-style-type: none"> • Frame specification: rows vs range <p>Joins in Detail</p> <ul style="list-style-type: none"> • Inner, outer, left, right, left_semi, left_anti • Cross joins and broadcast joins • Join optimization strategies • Broadcast hints • Skewed join handling <p>Set Operations & UDFs</p> <ul style="list-style-type: none"> • union, unionAll, unionByName • intersect, except (subtract) • User Defined Functions (UDFs): creation and registration • UDF performance considerations • Pandas UDFs (Vectorized UDFs) 	<ul style="list-style-type: none"> • Schema evolution handling <p>Spark Cluster Management</p> <ul style="list-style-type: none"> • Cluster managers: Standalone, YARN, Kubernetes, Mesos • Driver and executor configuration • spark-submit parameters • Dynamic allocation • Resource allocation strategies <p>Performance Tuning Deep Dive</p> <ul style="list-style-type: none"> • Memory management and tuning • Serialization: Java vs Kryo • Partition tuning: repartition vs coalesce • Caching strategies • Avoiding common pitfalls 	

Week 3 Databricks & Delta Lake Mastery	Project 2	Project 2	Project 2	Project 2	Project 2	
	Databricks Platform & Setup	Delta Lake Fundamentals	Advanced Delta Lake and Medallion Architecture	Databricks SQL and Streaming	Integration and Advanced Storage	

	MON	TUE	WED	THU	FRI	ENVIRONMENT
	<ul style="list-style-type: none"> Databricks Lakehouse platform overview <p>Databricks Architecture</p> <ul style="list-style-type: none"> Control Plane vs Data Plane Azure Databricks architecture components Account structure and workspace hierarchy Networking and security architecture <p>Azure Databricks Setup</p> <ul style="list-style-type: none"> Account creation and subscription setup Workspace creation and configuration Pricing models and SKUs Cost optimization strategies Resource tagging and management <p>Workspace Components & DBFS</p> <ul style="list-style-type: none"> UI navigation and customization Cluster management and configuration Notebook environment and collaboration <p>Databricks FileStore (DBFS)</p> <ul style="list-style-type: none"> DBFS structure and file management Mount points and external storage File upload and organization 	<p>Traditional Formats</p> <ul style="list-style-type: none"> Performance comparisons and benefits Migration strategies from Parquet File format optimization Creating and managing Delta tables Reading and writing Delta format Basic table operations and metadata <p>Advanced Delta Features</p> <ul style="list-style-type: none"> MERGE operations and UPSERT patterns Slowly Changing Dimensions (SCD) implementation Soft delete with incremental data 	<ul style="list-style-type: none"> operations and retention policies <p>Delta Table Performance Optimization</p> <ul style="list-style-type: none"> Partitioning strategies for Delta Clustering key selection File size optimization Bronze, Silver, Gold layers concept Data quality progression patterns Architecture design principles <p>Medallion Implementation</p> <ul style="list-style-type: none"> Bronze layer: raw data ingestion Silver layer: data cleaning and standardization Gold layer: business-ready datasets 	<p>Writing</p> <ul style="list-style-type: none"> External table creation and management Data source connections and formats Batch processing optimization Auto Loader and Streaming Auto Loader concepts and implementation File detection modes and schema evolution Incremental processing patterns <p>Data Streams</p> <ul style="list-style-type: none"> Structured streaming fundamentals Stream processing and windowing Real-time analytics patterns 	<p>External Integration</p> <ul style="list-style-type: none"> Multi-format support and optimization Compression algorithms and strategies External system integration patterns File organization and naming conventions Access pattern optimization Cross-region data management <p>Data Pipeline Architecture</p> <ul style="list-style-type: none"> End-to-end pipeline design Error handling and recovery strategies Monitoring and alerting basics 	

	MON	TUE	WED	THU	FRI	ENVIRONMENT
Week 4 Unity Catalog and Advanced Governance	Project 2	Project 2	Project 2	Project 2	Project 2	
	Unity Catalog Setup and Configuration	Advanced Delta Operations and Security	Volumes and Advanced Features	Data Sharing and Federation	Performance Optimization and Monitoring	
	Unity Catalog Introduction <ul style="list-style-type: none"> Three-level namespace architecture Centralized governance and security Unity Catalog vs. Legacy Hive Metastore 	Delta Tables <ul style="list-style-type: none"> MERGE statement syntax and patterns Conditional updates and complex logic Change Data Capture (CDC) implementation 	Volumes in Unity Catalog <ul style="list-style-type: none"> Managed vs. External volumes Volume management and use cases File-based operations in volumes 	Delta Sharing <ul style="list-style-type: none"> Cross-organization data sharing concepts Provider and recipient configuration Databricks Express Edition setup 	Advanced Performance Tuning <ul style="list-style-type: none"> Partitioning and bucketing strategies Clustering and optimization techniques Memory management and resource allocation 	
	Unity Catalog Setup <ul style="list-style-type: none"> Account-level metastore creation Workspace assignment and configuration Regional considerations and planning 	Incremental Data Processing <ul style="list-style-type: none"> Incremental loading strategies Delta time travel for auditing Data versioning and rollback 	Advanced Storage Management <ul style="list-style-type: none"> Volume-based file organization Access pattern optimization Integration with external storage systems 	Data Sharing Implementation <ul style="list-style-type: none"> Share creation and management Recipient access and permissions Monitoring and audit trails 	Query Optimization <ul style="list-style-type: none"> Execution plan analysis Index and statistics usage Cost-based optimization 	
	Catalog and External Location Management <ul style="list-style-type: none"> Creating and managing catalogs External location configuration Storage credential management 	Security and Secret Management <ul style="list-style-type: none"> Secret scopes and management Azure Key Vault integration Databricks-backed vs. Azure Key Vault-backed scopes 	Functions in Unity Catalog <ul style="list-style-type: none"> Scalar User Defined Functions (UDFs) Table-valued functions Function versioning and management 	Lakehouse Federation <ul style="list-style-type: none"> Query federation concepts Foreign catalog configuration Cross-platform query execution 	Monitoring and Observability <ul style="list-style-type: none"> Performance monitoring tools Query profiling and analysis Resource utilization tracking 	
	Managed vs. External Tables <ul style="list-style-type: none"> Table types and use cases Location and lifecycle management Performance implications and best practices 	User Management and Access Control <ul style="list-style-type: none"> User, service principal, and group management Role-based access control (RBAC) Object-level permissions and inheritance 	Advanced SQL Functions <ul style="list-style-type: none"> Complex analytical functions Custom aggregation functions Performance optimization techniques 	External Database Integration <ul style="list-style-type: none"> JDBC connection setup and optimization Query pushdown strategies Hybrid architecture patterns 	Troubleshooting and Debugging <ul style="list-style-type: none"> Common performance issues Debug techniques and tools Error handling and recovery 	
	Project 2	Project 2	Project 2	Project 2	Project 2	
	Workflow Orchestration and Job	Delta Live Tables and Streaming	Apache Airflow Integration	Databricks-Airflow Integration and DBT	Production Deployment	
	Project 2	Project 2	Project 2	Project 2	Project 2	
Week 5 Workflows, Orchestration, and Production	Project 2	Project 2	Project 2	Project 2	Project 2	

	MON	TUE	WED	THU	FRI	ENVIRONMENT
	Management Databricks Workflows <ul style="list-style-type: none"> • Job creation and configuration • Multi-task workflows and dependencies • Conditional execution and branching Notebook Orchestration <ul style="list-style-type: none"> • Notebook parameters and communication • Dynamic parameter passing Advanced Job Features <ul style="list-style-type: none"> • For-Each loops and conditional logic • Task failures and recovery strategies • Job scheduling and triggers Job Monitoring and Management <ul style="list-style-type: none"> • Job history and performance tracking • Alert configuration and notifications • Resource optimization for jobs 	Delta Live Tables (DLT) Introduction <ul style="list-style-type: none"> • Declarative ETL concepts • DLT vs. traditional ETL approaches • Pipeline architecture and design DLT Pipeline Development <ul style="list-style-type: none"> • Table and view definitions • Data quality constraints • Error handling and retry logic • Incremental processing patterns Advanced DLT Features <ul style="list-style-type: none"> • Change Data Capture (CDC) with DLT • Pipeline monitoring and debugging • Performance optimization techniques Streaming with DLT <ul style="list-style-type: none"> • Real-time data processing • Stream-batch integration • Windowing and aggregation 	Airflow Introduction <ul style="list-style-type: none"> • Workflow orchestration concepts • DAGs, operators, and tasks • Airflow architecture overview Airflow Setup and Configuration <ul style="list-style-type: none"> • Installation and database setup • Web UI overview and navigation • Connection and variable management DAG Development <ul style="list-style-type: none"> • DAG design principles and patterns • Task dependencies and relationships • Dynamic and parameterized DAGs Operators and Hooks <ul style="list-style-type: none"> • Common operators usage • Custom operator development • Airflow connections and hooks 	Databricks with Airflow <ul style="list-style-type: none"> • Databricks operators and integration • Job submission and monitoring • Parameter passing and error handling Hybrid Orchestration <ul style="list-style-type: none"> • Combining Databricks workflows with Airflow • Architecture patterns and best practices • Performance optimization strategies DBT Introduction <ul style="list-style-type: none"> • Analytics engineering concepts • DBT vs. traditional ETL approaches • Modern data stack integration DBT with Databricks <ul style="list-style-type: none"> • Databricks adapter setup • Model development and testing • Deployment strategies 	Production Architecture <ul style="list-style-type: none"> • Environment management (dev/test/prod) • CI/CD pipeline setup • Infrastructure as Code concepts Monitoring and Observability <ul style="list-style-type: none"> • Logging and metrics collection • Alert configuration and response • Performance monitoring strategies 	
Week 6 Project	Project 2	Project 2	Project 2	Project 2	Project Presentation	

PROJECT	TECHNOLOGIES
Project 1	Python, SQL, ETL
Project 2	Spark, ETL, Databricks



Copyright © 2025 Revature, LLC. All Rights Reserved.

By viewing this document, you agree that under copyright law all content displayed is the sole intellectual property of Revature, LLC, a technology advancement and consulting company based in Reston, VA. All content generated by a representative of Revature which is used for the company's advancement, development, or have otherwise been developed at the company's request, are the sole property of the company. No intellectual property may be reproduced, distributed, altered, or shared without the explicit permission from a representative of Revature.