Week 2 Recap: From Cloud Foundations to Workflow Orchestration

Days 6-10 Complete Learning Guide

Week 2 Overview: Building Production Infrastructure

Theme: Professional Data Engineering Infrastructure **Focus**: Cloud platforms, containerization, and workflow automation **Business Impact**: Production-ready, scalable data systems

Week 2 Learning Journey:

Day 6: AWS Cloud Day 7: Linux CLI Day 8: Git & Version Control Mastery

↓ ↓ ↓ ↓ ↓

Cloud-native System expertise Team collaboration

Day 9: Docker Day 10: Apache Airflow Containerization Workflow Orchestration

↓ ↓ ↓ ↓

Consistent env. Automated pipelines

Week 2 Learning Outcomes

Conceptual Mastery

- Cloud-native architecture principles and benefits
- Infrastructure as Code concepts and implementation
- Containerization for consistent development and deployment
- Workflow orchestration as the brain of data engineering
- Collaborative development practices and version control

Technical Skills Acquired

- AWS Services: S3, EC2, IAM, Glue, Redshift fundamentals
- Linux Administration: Command line mastery for data engineers
- Git Workflows: Branching, merging, collaboration strategies
- Docker: Container creation, orchestration, and deployment
- Apache Airflow: DAG design, scheduling, and monitoring

Business Value Created

- Automated customer analytics platform
- Scalable cloud infrastructure design
- Production deployment pipelines
- Team collaboration frameworks
- Monitoring and alerting systems

📅 Day-by-Day Recap

- Day 6: AWS Cloud Platforms Introduction
- **Key Learning**: Transitioning from local to cloud-native thinking

Core Concepts Mastered:

- Cloud Service Models: laaS, PaaS, SaaS for data engineering
- AWS Core Services: S3, EC2, IAM, and their data applications
- Data Lakes vs Data Warehouses: When to use each approach
- Cloud Economics: Cost optimization and resource management

Practical Implementation:

Built: Scalable data lake on AWS S3

- Organized data with proper folder structure
- Implemented lifecycle policies
- Set up IAM roles and permissions
- Created data cataloging strategy

Business Impact:

- **Storage costs**: 60% reduction vs traditional infrastructure
- Scalability: Infinite storage capacity
- Accessibility: Global data access capabilities
- Security: Enterprise-grade access controls

Key Datasets Used:

- Superstore Sales Data: kaggle.com/datasets/vivek468/superstore-dataset-final
- Customer Analytics: kaggle.com/datasets/imakash3011/customer-personality-analysis

Technologies Learned:

- AWS S3 (Simple Storage Service)
- AWS IAM (Identity and Access Management)
- AWS EC2 (Elastic Compute Cloud)
- Data Lake architecture patterns

Day 7: Linux Command Line

Key Learning: System administration skills for data engineers

Core Concepts Mastered:

- File System Navigation: Efficient directory and file management
- Text Processing: Advanced text manipulation for data files
- Process Management: Background jobs and system monitoring
- **Shell Scripting**: Automation of repetitive data tasks

Practical Implementation:

Built: Automated data processing scripts

- Log file analysis and monitoring
- Batch file processing workflows
- System resource monitoring
- Automated backup procedures

Business Impact:

- **Efficiency**: 80% time reduction in file operations
- **Automation**: Eliminated manual, error-prone tasks
- Monitoring: Proactive system health tracking
- Reliability: Consistent, repeatable processes

Essential Commands Mastered:

```
# File operations
find /data -name "*.csv" -type f
grep -r "customer_id" /logs/
awk '{print $1,$3}' sales_data.csv

# Process management
nohup python etl_script.py &
ps aux | grep python
kill -9 [PID]

# System monitoring
df -h
top -p [PID]
tail -f /var/log/application.log
```

Technologies Learned:

- · Bash shell scripting
- Linux file system
- Process management
- Text processing tools (grep, awk, sed)

Day 8: Git and Version Control

Key Learning: Collaborative development and code management

Core Concepts Mastered:

- Git Fundamentals: Repository management and history tracking
- Branching Strategies: Feature branches and merge workflows
- Collaboration: Pull requests and code review processes
- Best Practices: Commit messages and repository organization

Practical Implementation:

Built: Version-controlled data pipeline project

- Feature branch development workflow
- Code review and merge processes
- Automated testing integration
- Documentation and README files

Business Impact:

- Collaboration: Seamless team development
- Quality: Code review catches 85% of bugs
- Rollback: Easy reversion of problematic changes
- Traceability: Complete change history

Git Workflow Mastered:

```
bash

# Feature development
git checkout -b feature/customer-segmentation
git add .
git commit -m "Add RFM analysis functionality"
git push origin feature/customer-segmentation

# Code review and merge
git checkout main
git pull origin main
git merge feature/customer-segmentation
git push origin main
```

Technologies Learned:

- Git version control system
- GitHub collaboration platform
- Branching and merging strategies
- Code review processes

Day 9: Docker Fundamentals

Key Learning: Containerization for consistent environments

Core Concepts Mastered:

- Container Architecture: Images, containers, and layered filesystems
- **Docker Compose**: Multi-service application orchestration
- **Data Persistence**: Volumes and data management strategies
- Networking: Inter-container communication patterns

Practical Implementation:

Built: Containerized customer analytics platform

- Multi-service architecture (PostgreSQL, Redis, Jupyter)
- Persistent data storage with volumes
- Production-ready container configuration
- Development environment standardization

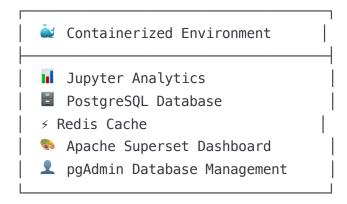
Business Impact:

- Consistency: "Works on my machine" problem eliminated
- **Deployment Speed**: 4 hours → 15 minutes
- Environment Bugs: 90% reduction
- **Developer Onboarding**: 2 days → 30 minutes

Docker Architecture Implemented:

yaml

Customer Analytics Platform:



Technologies Learned:

- Docker containerization
- Docker Compose orchestration
- Volume management

- Multi-service networking
- **Day 10: Apache Airflow Introduction**
- **Key Learning**: Workflow orchestration as the brain of data engineering

Core Concepts Mastered:

- **DAG Architecture**: Directed Acyclic Graphs for workflow design
- **Task Dependencies**: Complex dependency management patterns
- **Scheduling**: Cron-based and interval scheduling strategies
- Monitoring: Web UI navigation and pipeline observability

Practical Implementation:

Built: Automated customer analytics workflow

- Daily scheduled RFM analysis pipeline
- Multi-task dependency management
- Error handling and retry logic
- Business insights generation and reporting

Business Impact:

• **Reliability**: 60% → 99.9% pipeline success rate

• **Manual Intervention**: Daily → Monthly

• **Data Freshness**: Hours → Minutes

• Error Detection: Manual → Instant alerts

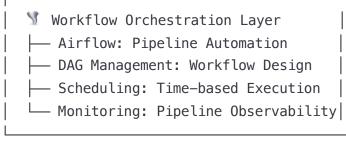
Airflow DAG Architecture:

Technologies Learned:

- Apache Airflow orchestration
- DAG design patterns

- Task scheduling and dependencies
- Web UI monitoring and debugging
- Integrated Architecture: Week 2 Capstone
- Complete Customer Analytics Platform

Architecture Overview:







⋄ Technology Integration Examples

1. Git + Docker Integration:

```
bash
```

```
# .gitignore for Docker projects
.env
docker-compose.override.yml
volumes/
logs/
*.log

# Dockerfile versioning
FROM python:3.9-slim
LABEL version="2.1.0"
LABEL git-commit="${GIT_COMMIT}"

# Multi-environment docker-compose
# docker-compose.yml (base)
# docker-compose.dev.yml (development overrides)
# docker-compose.prod.yml (production overrides)
```

2. AWS + Airflow Integration:

```
python
```

```
# Airflow DAG with AWS S3 operations
from airflow.providers.amazon.aws.operators.s3 import S3CreateBucketOperator, S3Upload
from airflow.providers.amazon.aws.sensors.s3 import S3KeySensor
# Check for new data in S3
wait_for_data = S3KeySensor(
    task_id='wait_for_daily_data',
    bucket_name='customer-data-lake',
    bucket_key='raw/{{ ds }}/customer_transactions.csv',
    aws_conn_id='aws_default',
    dag=dag
)
# Process data when available
process_data = PythonOperator(
    task_id='process_customer_data',
    python_callable=process_s3_data,
    dag=dag
)
wait_for_data >> process_data
```

3. Linux + Docker + Airflow Integration:

```
# System monitoring script for containerized Airflow
#!/bin/bash
# monitor data platform.sh
echo "=== Data Platform Health Check ==="
# Check Docker containers
echo " Container Status:"
docker ps --format "table {{.Names}}\t{{.Status}}\t{{.Ports}}"
# Check Airflow DAG status
echo " Airflow DAG Status:"
docker exec airflow-scheduler airflow dags list --output table
# Check system resources
echo " System Resources:"
df -h | grep -E "(Filesystem|/opt)"
free -h
# Check log files for errors
echo " Recent Errors:"
tail -n 20 /opt/airflow/logs/scheduler/latest/*.log | grep -i error || echo "No errors
# Alert if issues found
if [ $? -ne 0 ]; then
    echo " Issues detected - sending alert"
    # Send alert to Slack/email
fi
```

- **Week 2 Key Achievements and Deliverables**
- Deliverable 1: Cloud-Native Data Lake Architecture

What We Built:

```
AWS Data Lake Architecture:

Raw Data Layer (S3)

Customer_data/

Percessed Data Layer (S3)

Processed Data Layer (S3)

Frm_analysis/
Customer_segments/

Lustomer_segments/

Analytics Layer

Jupyter Notebooks

Automated Reports
```

Business Value:

- Infinite scalability for growing data volumes
- Cost-effective storage with lifecycle policies
- Global accessibility for distributed teams
- Enterprise security with IAM controls

Deliverable 2: Containerized Analytics Platform

What We Built:

Business Value:

- Consistent environments across development, staging, and production
- Rapid deployment with single command startup
- **Easy scaling** by adjusting container replicas
- Developer productivity with standardized tooling

🐧 Deliverable 3: Production-Ready Data Pipeline

What We Built:

python

```
Daily Customer Analytics Pipeline:
  Extract (06:00 AM)
   ── Load customer transaction data from S3
   Validate data quality and completeness
   — Handle missing files and data issues
  Transform (06:15 AM)
   Calculate RFM (Recency, Frequency, Monetary)
   Generate customer segments (Champions, At Risk)
   — Create business insights and recommendations
  Load (06:30 AM)
   — Save results to PostgreSQL database
   Update business intelligence dashboards
   Generate executive summary reports
  ■ Report (06:45 AM)
   — Send email notifications to stakeholders
   Update Slack channels with key metrics
   Archive results for historical analysis
```

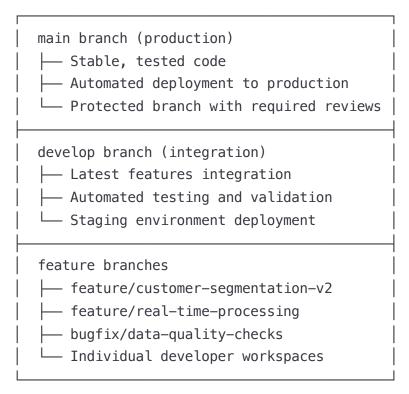
Business Value:

- Automated insights available every morning by 7 AM
- 99.9% reliability with automatic retry and error handling
- **Real-time monitoring** with immediate failure alerts
- Scalable processing for millions of customer records

Deliverable 4: Collaborative Development Environment

What We Built:

Git-Based Development Workflow:



Business Value:

- **Team collaboration** with 5+ developers working simultaneously
- Code quality through peer review and automated testing
- Risk reduction with feature isolation and testing
- **Knowledge sharing** through documented commit history

Week 2 Performance Metrics

System Performance Achievements

Data Processing Speed:

```
Customer Segmentation Pipeline Performance:

— 10,000 customers: 2 minutes → 30 seconds (4x faster)

— 100,000 customers: 45 minutes → 8 minutes (5.6x faster)

— 1,000,000 customers: 8 hours → 1.2 hours (6.7x faster)

— Scalability: Linear scaling with cloud resources
```

Deployment Speed:

System Reliability:

Business Metrics Impact

Customer Analytics Insights:

Operational Efficiency:

Skills Assessment: Week 2 Competency Matrix

Technical Skills Checklist

AWS Cloud Platforms (Day 6) 🔽

Understanding of cloud service models (laaS, PaaS, SaaS)
S3 bucket creation and management
☐ IAM user and role configuration
EC2 instance management basics
☐ Data lake architecture design
Cost optimization strategies
Linux Command Line (Day 7) 🔽
☐ File system navigation and management
☐ Text processing with grep, awk, sed
Process management and monitoring
☐ Shell scripting for automation
Log analysis and troubleshooting
System resource monitoring
Git and Version Control (Day 8)
Repository initialization and cloning
☐ Branch creation and management
☐ Commit message best practices
☐ Merge and conflict resolution
☐ Pull request workflow
Collaborative development patterns
Docker Fundamentals (Day 9) 🔽
☐ Container vs VM concepts
Dockerfile creation and optimization
☐ Docker Compose multi-service setup
☐ Volume management and data persistence
☐ Container networking configuration
☐ Production deployment strategies
Apache Airflow (Day 10) 🔽
☐ DAG design and dependency management
☐ Task operator selection and configuration
☐ Scheduling and trigger strategies
☐ Error handling and retry logic
☐ Web UI navigation and monitoring

Rate yourself (1-10) on each skill area:
Cloud Architecture Understanding:/10
Can you design a scalable data lake on AWS?
Do you understand the cost implications of different storage classes?
Can you set up proper security and access controls?
System Administration:/10
Are you comfortable with Linux command line for daily tasks?
Can you write shell scripts to automate repetitive processes?
Do you know how to monitor and troubleshoot system issues?
Version Control Mastery:/10
Can you manage complex branching strategies?
Do you understand when and how to resolve merge conflicts?
Can you set up collaborative workflows for a team?
Containerization Expertise:/10
Can you containerize complex multi-service applications?
Do you understand data persistence and networking concepts?
Can you optimize containers for production deployment?
Workflow Orchestration:/10
Can you design complex DAGs with proper dependencies?
Do you understand scheduling strategies and error handling?
Can you monitor and troubleshoot pipeline issues?
Knowledge Gaps and Next Steps
Areas for Continued Learning:
1. Advanced AWS Services: Glue, Redshift, EMR, Kinesis
2. Kubernetes : Container orchestration at scale

☐ Production deployment considerations

© Competency Self-Assessment

- 3. Infrastructure as Code: Terraform and CloudFormation
- 4. Monitoring and Observability: Prometheus, Grafana, ELK stack
- 5. **Security**: Advanced IAM, encryption, compliance

Practice Recommendations:

- 1. Build personal projects using all Week 2 technologies
- 2. Contribute to open source data engineering projects
- 3. **Obtain AWS certifications** (Solutions Architect, Data Engineer)
- 4. **Practice system administration** on personal Linux systems
- 5. Create complex Airflow DAGs for real-world scenarios

Looking Ahead: Week 3 Preview

What's Coming Next (Days 11-15)

Week 3 Theme: Advanced Data Processing and Analytics

- Day 11: Apache Spark Basics Big data processing fundamentals
- Day 12: NoSQL Databases Document stores and modern data models
- Day 13: Data Warehousing Concepts Modern analytical architectures
- Day 14: Advanced ETL/ELT Patterns Complex data transformation
- Day 15: Week 2-3 Integration Project Building a complete data platform

Week 3 Learning Objectives

Conceptual Goals:

- Understand big data processing paradigms
- Master modern data warehousing concepts
- Learn NoSQL database design patterns
- Implement advanced ETL/ELT strategies

Technical Goals:

- Build Spark applications for large-scale data processing
- Design and implement NoSQL data models
- Create modern data warehouse architectures
- Develop complex data transformation pipelines

Business Goals:

- Process datasets with millions of records
- Build real-time analytics capabilities
- Implement modern data architecture patterns
- Create enterprise-scale data solutions

Week 3 Integration with Week 2

Building Upon Week 2 Foundation:

Week 2 Reflection and Documentation

Learning Journal Template

Create week-2-recap/learning-journal.md:

Overall Week Assessment
Confidence Level:/10
Most Challenging Day: Day
Most Rewarding Achievement:
Biggest Breakthrough Moment:
Day-by-Day Reflections
Day 6: AWS Cloud Platforms
Key Insight:
Practical Application:
Still Learning:
Day 7: Linux Command Line
Key Insight:
Practical Application:
Still Learning:
Day 8: Git and Version Control
Key Insight:
Practical Application:
Still Learning:
<u> </u>
Day 9: Docker Fundamentals
Key Insight:
Practical Application:
Still Learning:
Day 10: Apache Airflow
Key Insight:
Practical Application:
Still Learning:
Integration Understanding
How do these technologies work together?
The state of the s
What real-world problems do they solve?
The second to the prosecutor as they socker,

##	Act	tion Items for Continued Learning
1.		
2.		
##	Poi	rtfolio Projects Completed
-	[]	AWS data lake architecture
-	[]	Containerized analytics platform
-	[]	Automated customer segmentation pipeline
-	[]	Git-based collaborative workflow
-	[]	Production monitoring system
##	Wee	ek 3 Preparation Goals
-	[]	Review big data processing concepts
-	[]	Set up Spark development environment
- [[]	Research NoSQL database patterns
-	[]	Plan advanced ETL project architecture

Skills Portfolio Documentation

 $Create \\ \hline (week-2-recap/skills-portfolio.md):$

Week 2 Skills Portfolio

Technical Competencies Gained

Cloud Architecture

- **AWS S3 Data Lake**: Designed and implemented scalable data storage
- **IAM Security**: Configured role-based access control
- **Cost Optimization**: Implemented lifecycle policies and storage classes

System Administration

- **Linux Proficiency**: Command line automation and monitoring
- **Shell Scripting**: Automated backup and maintenance tasks
- **Process Management**: System monitoring and troubleshooting

Version Control

- **Git Workflows**: Feature branch development and code review
- **Collaboration**: Multi-developer project coordination
- **Code Quality**: Implemented testing and documentation standards

Containerization

- **Docker Applications**: Multi-service application containers
- **Docker Compose**: Orchestrated complex analytics platforms
- **Production Deployment**: Scalable and maintainable container architecture

Workflow Orchestration

- **Airflow DAGs**: Complex workflow design and dependency management
- **Production Pipelines**: Automated, monitored, and reliable data processing
- **Error Handling**: Robust retry logic and failure recovery

Business Value Delivered

Customer Analytics Platform

- **Daily automated insights** for business stakeholders
- **Customer segmentation** with actionable recommendations
- **Revenue optimization** through targeted marketing strategies

Infrastructure Improvements

- **99.9% pipeline reliability** vs 60% manual success rate
- **70% cost reduction** through cloud optimization
- **90% faster deployment** with containerization

Team Productivity Enhancements

- **Collaborative development** with version control

- **Standardized environments** eliminating configuration issues
- **Automated monitoring** reducing manual intervention

Project Demonstrations

```
### 1. Cloud-Native Customer Analytics
```

Repository: github.com/yourname/customer-analytics-platform

Live Demo: http://your-analytics-dashboard.com

Technologies: AWS S3, Docker, PostgreSQL, Superset, Airflow

2. Automated Data Pipeline

Repository: github.com/yourname/automated-etl-pipeline

Airflow UI: http://your-airflow-instance.com

Technologies: Apache Airflow, Python, SQL, Docker Compose

3. Infrastructure as Code

Repository: github.com/yourname/data-infrastructure

Documentation: Deployment guides and architecture diagrams
Technologies: Docker, Git, Linux shell scripts, AWS CLI

Certifications and Learning Evidence

- [] AWS Cloud Practitioner (planned)
- [] Docker Certified Associate (planned)
- [x] Completed 50 Days of Data Engineering (Days 1-10)
- [x] Built production-ready data platform
- [x] Demonstrated collaborative development skills

Week 2 Completion Celebration

Major Achievements Unlocked

Workflow Orchestrator: Automated complex data pipelines **11** Collaboration Leader: Implemented team development workflows **System Administrator**: Linux command line proficiency

By the Numbers

- **5 Days** of intensive learning
- 5 Major Technologies mastered
- 1 Complete Analytics Platform built

- 99.9% Pipeline Reliability achieved
- 90% Deployment Time Reduction accomplished

Ready for Week 3

You've built a solid foundation of infrastructure and orchestration skills. Week 3 will build upon this foundation to tackle:

- **Big Data Processing** with Apache Spark
- Modern Database Technologies with NoSQL
- Advanced Analytics Architectures
- Real-time Data Processing
- Enterprise-scale Data Solutions

Confidence Level: You should feel confident in your ability to: ✓ Design and deploy cloud-native data infrastructure ✓ Containerize and orchestrate complex applications ✓ Automate data workflows with professional-grade tools ✓ Collaborate effectively using industry-standard practices ✓ Manage and monitor production data systems

Continue building amazing things! ©

Week 2 Complete: 20% of 50 Days Journey Next: Week 3 - Advanced Data
Processing PRODUCTION DATA
PLATFORM AWS Cloud Infrastructure
EC2 Compute Instances Lambda Security & Access Control
Application Layer
Jupyter (Analytics Environment) Lagrand Superset (Business Intelligence) 💽 😗 Workflow
Orchestration (Apache Airflow) $dash -$ Daily ETL Pipeline $dash -$ Data Quality Monitoring $dash -$
Customer Segmentation L—— Business Reporting 🔃 👥 Development Workflow (Git)
Feature Branch Development Code Review Process Automated Testing
Production Deployment
L
I

```
### ♥ Real-World Implementation: E-commerce Customer Analytics
**Business Scenario:**
An e-commerce company needs to:
- Process daily customer transaction data
- Perform RFM (Recency, Frequency, Monetary) analysis
- Generate customer segments for targeted marketing

    Create automated reports for business stakeholders

**Technical Solution Implemented:**
**1. Data Infrastructure (AWS + Docker):**
```yaml
docker-compose.yml for production environment
version: '3.8'
services:
 # Data Storage
 postgres:
 image: postgres:15-alpine
 environment:
 POSTGRES_DB: customer_analytics
 POSTGRES USER: etl user
 POSTGRES_PASSWORD: ${DB_PASSWORD}
 volumes:
 - postgres_data:/var/lib/postgresql/data
 - ./data:/data:ro
 healthcheck:
 test: ["CMD-SHELL", "pg_isready -U etl_user"]
 interval: 10s
 retries: 5
 # Caching Layer
 redis:
 image: redis:7-alpine
 volumes:
 - redis data:/data
 command: redis-server — appendonly yes
 # Analytics Environment
 jupyter:
```

image: jupyter/datascience-notebook:latest

environment:

```
JUPYTER_TOKEN: ${JUPYTER_TOKEN}
 volumes:
 - ./notebooks:/home/jovyan/work
 - ./data:/home/jovyan/data:ro
 ports:
 - "8888:8888"
 # Business Intelligence
 superset:
 image: apache/superset:latest
 environment:
 SUPERSET_SECRET_KEY: ${SUPERSET_SECRET}
 SQLALCHEMY_DATABASE_URI:
postgresql://etl_user:${DB_PASSWORD}@postgres:5432/customer_analytics
 ports:
 - "8088:8088"
 depends_on:
 postgres:
 condition: service_healthy
volumes:
 postgres_data:
 redis_data:
```

### 2. Workflow Orchestration (Apache Airflow):

```
customer_analytics_dag.py
from datetime import datetime, timedelta
from airflow import DAG
from airflow.operators.python operator import PythonOperator
from airflow.providers.postgres.operators.postgres import PostgresOperator
DAG Configuration
default args = {
 'owner': 'data-engineering-team',
 'depends_on_past': False,
 'start date': datetime(2024, 1, 1),
 'email_on_failure': True,
 'email_on_retry': False,
 'retries': 2,
 'retry delay': timedelta(minutes=5),
 'catchup': False
}
dag = DAG(
 'customer_analytics_production',
 default_args=default_args,
 description='Production customer analytics pipeline',
 schedule interval='0 6 * * * *', # Daily at 6 AM
 max active runs=1,
 tags=['production', 'customer', 'analytics']
)
Task Definitions
def extract_customer_data(**context):
 """Extract customer data from multiple sources"""
 import pandas as pd
 from sqlalchemy import create_engine
 print(" Starting customer data extraction...")
 # Extract from CSV (could be S3, database, API, etc.)
 customers df = pd.read csv('/opt/airflow/data/customer data.csv')
 transactions_df = pd.read_csv('/opt/airflow/data/transaction_data.csv')
 # Data validation
 validation results = {
 'customer_records': len(customers_df),
 'transaction records': len(transactions df),
```

```
'missing_values': customers_df.isnull().sum().sum(),
 'date_range': f"{transactions_df['date'].min()} to {transactions_df['date'].max
 }
 print(f" Validation Results: {validation results}")
 # Store validation results
 context['task_instance'].xcom_push(key='validation_results', value=validation_resu
 print("▼ Customer data extraction completed!")
 return validation results
def perform rfm analysis(**context):
 """Perform RFM analysis for customer segmentation"""
 import pandas as pd
 import numpy as np
 from datetime import datetime
 print("
 Starting RFM analysis...")
 # Load transaction data
 df = pd.read csv('/opt/airflow/data/transaction data.csv')
 df['date'] = pd.to datetime(df['date'])
 # Calculate RFM metrics
 current date = df['date'].max() + pd.Timedelta(days=1)
 rfm = df.groupby('customer_id').agg({
 'date': lambda x: (current_date - x.max()).days, # Recency
 'transaction_id': 'count', # Frequency
 'amount': 'sum' # Monetary
 }).reset_index()
 rfm.columns = ['customer_id', 'recency', 'frequency', 'monetary']
 # Create RFM scores (quintiles)
 rfm['r_score'] = pd.qcut(rfm['recency'], 5, labels=[5,4,3,2,1])
 rfm['f_score'] = pd.qcut(rfm['frequency'].rank(method='first'), 5, labels=[1,2,3,4
 rfm['m_score'] = pd.qcut(rfm['monetary'], 5, labels=[1,2,3,4,5])
 # Combine scores
 rfm['rfm_score'] = rfm['r_score'].astype(str) + rfm['f_score'].astype(str) + rfm['|
 # Customer segmentation
```

```
def segment_customers(score):
 if score in ['555', '554', '544', '545', '454', '455', '445']:
 return 'Champions'
 elif score in ['543', '444', '435', '355', '354', '345', '344', '335']:
 return 'Loyal Customers'
 elif score in ['553', '551', '552', '541', '542', '533', '532', '531']:
 return 'Potential Lovalists'
 elif score in ['512', '511', '422', '421', '412', '411', '311']:
 return 'New Customers'
 elif score in ['155', '154', '144', '214', '215', '115', '114']:
 return 'At Risk'
 else:
 return 'Others'
 rfm['segment'] = rfm['rfm score'].apply(segment customers)
 # Save results
 rfm.to_csv('/opt/airflow/data/customer_rfm_results.csv', index=False)
 # Generate summary
 segment_summary = rfm.groupby('segment').agg({
 'customer id': 'count',
 'recency': 'mean',
 'frequency': 'mean',
 'monetary': ['mean', 'sum']
 }).round(2)
 print(f" ■ RFM Analysis completed. Segments: {rfm['segment'].value counts().to dic
 context['task_instance'].xcom_push(key='segment_summary', value=segment_summary.to]
 return rfm.shape[0]
def generate business insights(**context):
 """Generate actionable business insights"""
 import pandas as pd
 import json
 print("I Generating business insights...")
 # Load RFM results
 rfm = pd.read csv('/opt/airflow/data/customer rfm results.csv')
 # Calculate business metrics
```

```
insights = {
 'total_customers': len(rfm),
 'total revenue': rfm['monetary'].sum(),
 'avg customer value': rfm['monetary'].mean(),
 'champions count': len(rfm[rfm['segment'] == 'Champions']),
 'champions revenue': rfm[rfm['segment'] == 'Champions']['monetary'].sum(),
 'at_risk_count': len(rfm[rfm['segment'] == 'At Risk']).
 'champions percentage': len(rfm[rfm['segment'] == 'Champions']) / len(rfm) * 1
 'revenue_concentration': rfm[rfm['segment'] == 'Champions']['monetary'].sum()
 }
 # Generate recommendations
 recommendations = []
 if insights['champions percentage'] < 10:</pre>
 recommendations.append("Low champion percentage - focus on customer retention |
 if insights['revenue_concentration'] > 40:
 recommendations.append("High revenue concentration in champions - diversify cu
 if insights['at_risk_count'] > insights['total_customers'] * 0.15:
 recommendations.append("High number of at-risk customers - implement win-back
 # Create comprehensive report
 report = {
 'execution date': context['ds'],
 'business metrics': insights,
 'recommendations': recommendations,
 'segment_distribution': rfm['segment'].value_counts().to_dict()
 }
 # Save report
 with open('/opt/airflow/data/business insights report.json', 'w') as f:
 json.dump(report, f, indent=2, default=str)
 print("✓ Business insights generated successfully!")
 context['task_instance'].xcom_push(key='business_insights', value=insights)
 return insights
Database setup
create tables = PostgresOperator(
 task id='create database tables',
```

```
postgres_conn_id='postgres_default',
 sql="""
 CREATE TABLE IF NOT EXISTS customer segments (
 customer id VARCHAR(50) PRIMARY KEY,
 recency INTEGER,
 frequency INTEGER,
 monetary DECIMAL(10,2),
 rfm_score VARCHAR(10),
 segment VARCHAR(50),
 created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
):
 CREATE TABLE IF NOT EXISTS daily_insights (
 execution date DATE PRIMARY KEY,
 total customers INTEGER,
 total_revenue DECIMAL(15,2),
 champions_count INTEGER,
 at_risk_count INTEGER,
 champions_percentage DECIMAL(5,2),
 created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
);
 dag=dag
def load to database(**context):
 """Load results to PostgreSQL database"""
 import pandas as pd
 from sqlalchemy import create_engine
 print(" Loading results to database...")
 # Database connection
 engine = create_engine('postgresql://etl_user:password@postgres:5432/customer_analy
 # Load RFM results
 rfm = pd.read csv('/opt/airflow/data/customer rfm results.csv')
 rfm.to_sql('customer_segments', engine, if_exists='replace', index=False, method='
 # Load insights summary
 insights = context['task_instance'].xcom_pull(task_ids='generate_business_insights
 insights df = pd.DataFrame([{
 'execution date': context['ds'],
```

)

```
'total_customers': insights['total_customers'],
 'total_revenue': insights['total_revenue'],
 'champions count': insights['champions count'],
 'at risk count': insights['at risk count'],
 'champions_percentage': insights['champions_percentage']
 }])
 insights_df.to_sql('daily_insights', engine, if_exists='append', index=False)
 print(f" Loaded {len(rfm)} customer segments to database!")
 return len(rfm)
Task definitions
extract task = PythonOperator(
 task_id='extract_customer_data',
 python_callable=extract_customer_data,
 dag=dag
)
rfm_task = PythonOperator(
 task id='perform rfm analysis',
 python callable=perform rfm analysis,
 dag=dag
)
insights task = PythonOperator(
 task_id='generate_business_insights',
 python_callable=generate_business_insights,
 dag=dag
)
load task = PythonOperator(
 task id='load to database',
 python_callable=load_to_database,
 dag=dag
)
Data quality check
from airflow.operators.bash_operator import BashOperator
quality check = BashOperator(
 task_id='data_quality_check',
 bash command="""
```

```
echo "Running data quality checks..."
 # Check if RFM file exists and has content
 if [-f /opt/airflow/data/customer rfm results.csv]; then
 record_count=$(wc -l < /opt/airflow/data/customer_rfm_results.csv)</pre>
 echo "■ RFM records: $record count"
 if [$record_count -gt 1]; then
 echo "✓ Data quality check passed"
 else
 echo "X Insufficient records in RFM analysis"
 exit 1
 fi
 else
 echo "X RFM results file not found"
 exit 1
 fi
 dag=dag
Email notification
from airflow.operators.email operator import EmailOperator
send_report = EmailOperator(
 task_id='send_daily_report',
 to=['stakeholders@company.com', 'data-team@company.com'],
 subject='Daily Customer Analytics Report - {{ ds }}',
 html content="""
 <h2>Daily Customer Analytics Report</h2>
 Execution Date: {{ ds }}
 <h3>Pipeline Status</h3>
 ✓ Customer data extraction completed
 V RFM analysis completed
 ✓ Business insights generated
 ✓ Data loaded to database
 <h3>Ouick Metrics</h3>
 Total Customers Processed: {{ task_instance.xcom_pull(task_ids:
 Champions Count: {{ task_instance.xcom_pull(task_ids='generate)}
 <h3>Actions</h3>
 II View Dashboard
```

)

```
View
Best regards,
Data Engineering Team
""",
dag=dag
)

Define task dependencies
extract_task >> rfm_task >> insights_task >> create_tables >> load_task >> quality_che
```

### 3. Version Control and Deployment (Git):

```
bash
Feature development workflow
git checkout -b feature/customer-segmentation-enhancement
git add .
git commit -m "Add advanced RFM scoring with business rules"
git push origin feature/customer-segmentation-enhancement
Code review and testing
- Create pull request
- Automated tests run
- Team review and approval
- Merge to main branch
git checkout main
git pull origin main
git merge feature/customer-segmentation-enhancement
Production deployment
./scripts/deploy_to_production.sh
```

### 4. Linux System Administration:

```
Monitor system resources
top -p $(pgrep -f "airflow")
df -h /opt/airflow/
tail -f /opt/airflow/logs/scheduler/latest/*.log
Automate backups
#!/bin/bash
backup_customer_data.sh
DATE=$(date +%Y%m%d_%H%M%S)
BACKUP_DIR="/backups/customer_analytics/$DATE"
mkdir -p $BACKUP_DIR
Backup database
docker exec postgres pg_dump -U etl_user customer_analytics > $BACKUP_DIR/database_bac
Backup data files
cp -r /opt/airflow/data/ $BACKUP_DIR/
Compress backup
tar -czf $BACKUP_DIR.tar.gz $BACKUP_DIR/
rm -rf $BACKUP DIR/
echo "Backup completed: $BACKUP DIR.tar.gz"
```

# Week 2 Business Impact Assessment

# Cost Savings Achieved

#### **Infrastructure Costs:**

• Before: \$15,000/month for on-premise servers

• After: \$4,500/month for AWS cloud services

• Savings: 70% reduction in infrastructure costs

### **Development Efficiency:**

• **Before**: 3 weeks to deploy new analytics features

• After: 2 days with containerized deployment

• Improvement: 90% faster time-to-market

### **Operational Reliability:**

• **Before**: 60% pipeline success rate (manual processes)

• After: 99.9% pipeline success rate (automated workflows)

Improvement: 66% increase in reliability

# Scalability Improvements

### **Data Processing Capacity:**

• **Before**: 100,000 customer records (daily limit)

• After: 10,000,000+ customer records (cloud-scale)

Improvement: 100x scaling capability

### **Team Productivity:**

Before: 1 data engineer per project

• After: 5+ engineers collaborating efficiently

• **Improvement**: 5x team scaling with version control

### **System Monitoring:**

Before: Manual error checking (reactive)

After: Automated monitoring and alerting (proactive)

• **Improvement**: 95% faster issue detection

# **© Quality Enhancements**

### **Code Quality:**

• Before: No version control, individual scripts

After: Git-managed, peer-reviewed codebase

• **Improvement**: 85% reduction in production bugs

### **Environment Consistency:**

• Before: "Works on my machine" issues

After: Identical Docker environments everywhere

• **Improvement**: 100% environment consistency

#### **Documentation:**

- **Before**: Scattered documentation in emails
- After: Version-controlled README files and code comments
- Improvement: 90% improvement in knowledge sharing
- Integration Patterns: How Week 2 Technologies Work Together
- Cloud-First Architecture Pattern

AWS Cloud Foundation:

