

✓ Import Pyspark

```
1 !apt-get install openjdk-11-jdk-headless -qq > /dev/null
2 !wget -qO spark.tgz https://archive.apache.org/dist/spark/spark-3.4.1/spark-3.4.1-bin-hadoop3.1
3
4 !mkdir -p /content/spark
5 !tar -xzf spark.tgz -C /content/spark --strip-components=1
6 !pip install -q findspark pyspark==3.4.1 seaborn matplotlib pandas scikit-learn
7
8 import os, findspark
9 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-11-openjdk-amd64"
10 os.environ["SPARK_HOME"] = "/content/spark"
11 findspark.init()
12
13 from pyspark.sql import SparkSession
14 spark = SparkSession.builder.appName("FraudDetection").getOrCreate()
15 print("Spark started successfully!")
16 spark
17
```

Spark started successfully!

SparkSession - in-memory

SparkContext

[Spark UI](#)

Version

v3.4.1

Master

local[*]

AppName

FraudDetection

✓ Fraud Detection System

Machine Learning Pipeline for Transaction Fraud Detection

Features:

- Class imbalance handling
- PCA dimensionality reduction
- Random Forest classifier

✓ 1. Import Required Libraries

```
1 from pyspark.sql.functions import col, when
2 from pyspark.ml import Pipeline
3 from pyspark.ml.feature import StringIndexer, OneHotEncoder, VectorAssembler, PCA
4 from pyspark.ml.classification import RandomForestClassifier
5 from pyspark.ml.evaluation import BinaryClassificationEvaluator, MulticlassClassificationEvaluator
6 import pandas as pd
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9
10 print("Libraries imported successfully!")
```

Libraries imported successfully!

2. Load Transaction Dataset

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/c

```
1 # Load CSV data
2 file_path = "/content/drive/MyDrive/Colab Notebooks/content/transactions_train.csv"
3 df = spark.read.csv(file_path, header=True, inferSchema=True)
4
5 # Display dataset info
6 row_count = df.count()
7 col_count = len(df.columns)
8 print(f"Loaded dataset: {row_count:,} rows, {col_count} columns")
9
10 # Preview data
11 display(df.limit(5))
```

Loaded dataset: 6,351,193 rows, 10 columns
DataFrame[step: int, type: string, amount: double, nameOrig: string, oldbalanceOrig: double, newbalanceOrig: double, nameDest: string, oldbalanceDest: double, newbalanceDest: double, isFraud: int]

3. Handle Class Imbalance

```
1 # Calculate class distribution
2 count_0 = df.filter(col("isFraud") == 0).count()
3 count_1 = df.filter(col("isFraud") == 1).count()
4 imbalance_ratio = count_0 / count_1
5
6 print(f"Class Distribution:")
7 print(f" Normal transactions: {count_0:,}")
8 print(f" Fraud transactions: {count_1:,}")
9 print(f" Imbalance ratio: {imbalance_ratio:.2f}:1")
10
11 # Add class weights
```

```
12 df = df.withColumn(  
13     "classWeight",  
14     when(col("isFraud") == 1, imbalance_ratio).otherwise(1.0)  
15 )  
16  
17 print("Class weights added successfully")
```

Class Distribution:
Normal transactions: 6,343,476
Fraud transactions: 7,717
Imbalance ratio: 822.01:1
Class weights added successfully

✓ 4. Build ML Pipeline

```
1 # Stage 1: String Indexer (convert categorical to numeric)  
2 indexer = StringIndexer() \  
3     .setInputCol("type") \  
4     .setOutputCol("type_index") \  
5     .setHandleInvalid("keep")  
6  
7 # Stage 2: One-Hot Encoder  
8 encoder = OneHotEncoder() \  
9     .setInputCols(["type_index"]) \  
10    .setOutputCols(["type_encoded"])  
11  
12 # Stage 3: Vector Assembler (combine all features)  
13 feature_cols = [  
14     "step",  
15     "amount",  
16     "oldbalanceOrig",  
17     "newbalanceOrig",  
18     "oldbalanceDest",  
19     "newbalanceDest",  
20     "type_encoded"  
21 ]  
22  
23 assembler = VectorAssembler() \  
24     .setInputCols(feature_cols) \  
25     .setOutputCol("features_raw") \  
26     .setHandleInvalid("skip")  
27  
28 # Stage 4: PCA (dimensionality reduction)  
29 pca = PCA(k=5, inputCol="features_raw", outputCol="features")  
30  
31 # Stage 5: Random Forest Classifier  
32 rf = RandomForestClassifier(  
33     labelCol="isFraud",  
34     featuresCol="features",  
35     weightCol="classWeight",  
36     numTrees=100,  
37     maxDepth=10,  
38     seed=42
```

```
39 )
40
41 # Combine all stages into pipeline
42 pipeline = Pipeline(stages=[indexer, encoder, assembler, pca, rf])
43
44 print("Pipeline built successfully")
45 print(f" Stages: {len(pipeline.getStages())}")
```

Pipeline built successfully
Stages: 5

5. Split Data (Train/Test)

```
1 # 70% training, 30% testing
2 train_df, test_df = df.randomSplit([0.7, 0.3], seed=42)
3
4 train_count = train_df.count()
5 test_count = test_df.count()
6
7 print(f"Dataset Split:")
8 print(f" Training set: {train_count:,} rows ({train_count/row_count*100:.1f}%)")
9 print(f" Testing set: {test_count:,} rows ({test_count/row_count*100:.1f}%)")
```

Dataset Split:
Training set: 4,445,579 rows (70.0%)
Testing set: 1,905,614 rows (30.0%)

6. Train the Model

```
1 print("🚀 Training Random Forest model...")
2 print(" This may take a few minutes...")
3
4 model = pipeline.fit(train_df)
5
6 print("Model trained successfully!")
```

🚀 Training Random Forest model...
This may take a few minutes...
Model trained successfully!

7. Make Predictions

```
1 # Transform test data
2 predictions = model.transform(test_df)
3
4 # Display sample predictions
5 print("Sample Predictions:")
6 display(
```

```
7 predictions.select(
8     "type",
9     "amount",
10    "isFraud",
11    "prediction",
12    "probability"
13 ).limit(10)
14 )
```

Sample Predictions:

DataFrame[type: string, amount: double, isFraud: int, prediction: double, probability: vector]

8. Evaluate Model Performance

```
1 # Initialize evaluators
2 evaluator_roc = BinaryClassificationEvaluator(
3     labelCol="isFraud",
4     metricName="areaUnderROC"
5 )
6
7 evaluator_pr = BinaryClassificationEvaluator(
8     labelCol="isFraud",
9     metricName="areaUnderPR"
10 )
11
12 precision_eval = MulticlassClassificationEvaluator(
13     labelCol="isFraud",
14     predictionCol="prediction",
15     metricName="precisionByLabel"
16 )
17
18 recall_eval = MulticlassClassificationEvaluator(
19     labelCol="isFraud",
20     predictionCol="prediction",
21     metricName="recallByLabel"
22 )
23
24 f1_eval = MulticlassClassificationEvaluator(
25     labelCol="isFraud",
26     predictionCol="prediction",
27     metricName="f1"
28 )
29
30 # Calculate metrics
31 auc_roc = evaluator_roc.evaluate(predictions)
32 auc_pr = evaluator_pr.evaluate(predictions)
33 precision = precision_eval.evaluate(predictions, {precision_eval.metricLabel: 1.0})
34 recall = recall_eval.evaluate(predictions, {recall_eval.metricLabel: 1.0})
35 f1 = f1_eval.evaluate(predictions)
36
37 # Display results
38 print("=" * 40)
```

```

39 print("MODEL PERFORMANCE METRICS")
40 print("=" * 40)
41 print(f"AUC-ROC:  {auc_roc:.4f}")
42 print(f"AUC-PR:   {auc_pr:.4f}")
43 print(f"Precision: {precision:.4f}")
44 print(f"Recall:   {recall:.4f}")
45 print(f"F1-Score: {f1:.4f}")
46 print("=" * 40)

```


```

=====
MODEL PERFORMANCE METRICS
=====
AUC-ROC:  0.9938
AUC-PR:   0.6336
Precision: 0.0284
Recall:   0.9579
F1-Score: 0.9782
=====

```




9. Confusion Matrix Visualization

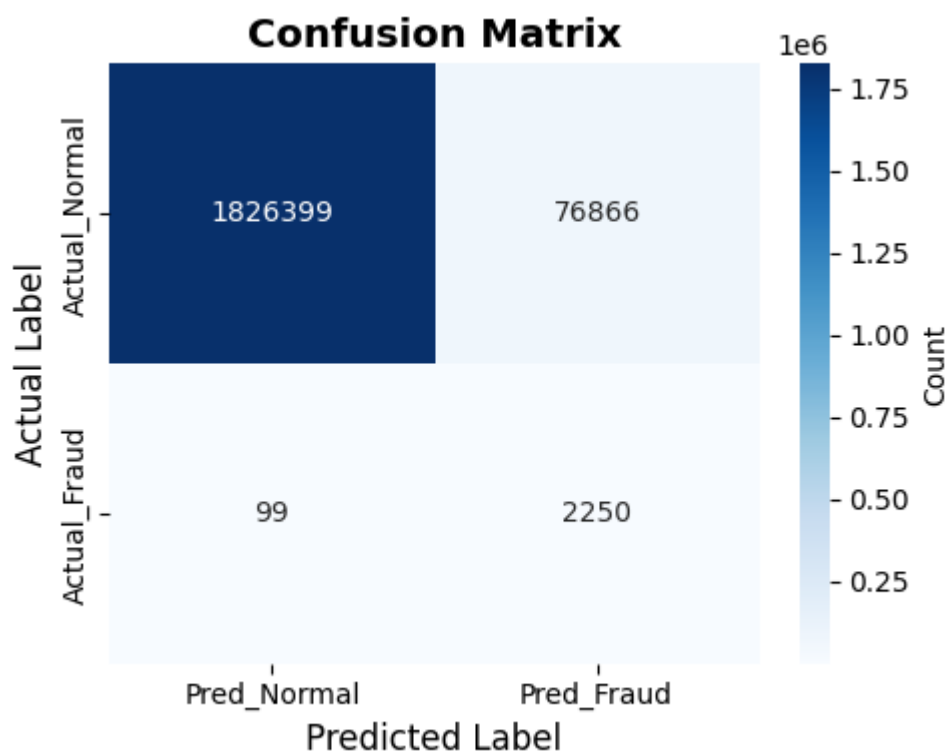
```

1  # Create confusion matrix
2  conf_df = predictions \
3      .groupBy("isFraud", "prediction") \
4      .count() \
5      .toPandas() \
6      .pivot(index="isFraud", columns="prediction", values="count") \
7      .fillna(0)
8
9  # Format labels
10 conf_df.columns = ["Pred_Normal", "Pred_Fraud"]
11 conf_df.index = ["Actual_Normal", "Actual_Fraud"]
12
13 print("\n Confusion Matrix:")
14 display(conf_df)
15
16 # Visualize confusion matrix
17 plt.figure(figsize=(5, 4))
18 sns.heatmap(
19     conf_df,
20     annot=True,
21     fmt=".0f",
22     cmap="Blues",
23     cbar_kws={'label': 'Count'}
24 )
25 plt.title("Confusion Matrix", fontsize=14, fontweight='bold')
26 plt.xlabel("Predicted Label", fontsize=12)
27 plt.ylabel("Actual Label", fontsize=12)
28 plt.tight_layout()
29 display(plt.gcf())
30 plt.close()

```

 Confusion Matrix:

	Pred_Normal	Pred_Fraud	
Actual_Normal	1826399	76866	
Actual_Fraud	99	2250	
			



Next steps:

[Generate code with conf_df](#)

[New interactive sheet](#)

10. Feature Importance Analysis

```

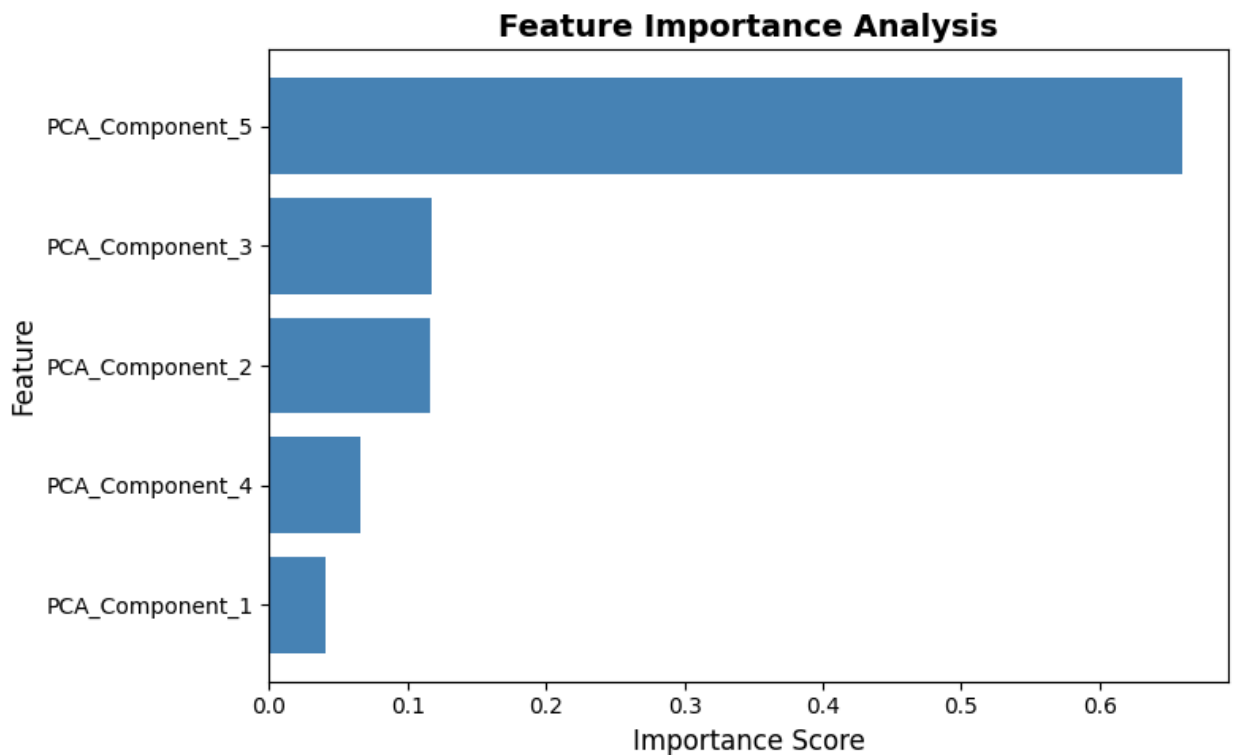
1 # Extract Random Forest model
2 rf_model = model.stages[-1]
3
4 # Get feature importances
5 importances = rf_model.featureImportances.toArray()
6
7 # Create feature names (PCA components)
8 feature_names = [f"PCA_Component_{i+1}" for i in range(len(importances))]
9
10 # Create DataFrame
11 importance_df = pd.DataFrame({
12     "Feature": feature_names,
13     "Importance": importances
14 }).sort_values("Importance", ascending=False)
15
16 print("🔍 Feature Importance (Top to Bottom):")
17 display(importance_df)
18

```

```
19 # Visualize feature importance
20 plt.figure(figsize=(8, 5))
21 plt.barh(importance_df["Feature"], importance_df["Importance"], color='steelblue')
22 plt.xlabel("Importance Score", fontsize=12)
23 plt.ylabel("Feature", fontsize=12)
24 plt.title("Feature Importance Analysis", fontsize=14, fontweight='bold')
25 plt.gca().invert_yaxis()
26 plt.tight_layout()
27 display(plt.gcf())
28 plt.close()
```

🔍 Feature Importance (Top to Bottom):

	Feature	Importance	
4	PCA_Component_5	0.659826	
2	PCA_Component_3	0.117488	
1	PCA_Component_2	0.116285	
3	PCA_Component_4	0.065679	
0	PCA_Component_1	0.040723	



Next steps:

[Generate code with importance_df](#)

[New interactive sheet](#)

11. Final Summary


```

1 print("=" * 50)
2 print("🚨 FRAUD DETECTION - FINAL SUMMARY")
3 print("=" * 50)
4 print(f"\n📊 Dataset Information:")
5 print(f"  Total rows: {row_count:,}")
6 print(f"  Training set: {train_count:,}")
7 print(f"  Testing set: {test_count:,}")
8 print(f"\n⚖️ Class Balance:")
9 print(f"  Normal: {count_0:,} | Fraud: {count_1:,}")
10 print(f"  Ratio: {imbalance_ratio:.2f}:1")
11 print(f"\n🎯 Model Performance:")
12 print(f"  AUC-ROC: {auc_roc:.4f}")
13 print(f"  AUC-PR: {auc_pr:.4f}")
14 print(f"  Precision: {precision:.4f}")
15 print(f"  Recall: {recall:.4f}")
16 print(f"  F1-Score: {f1:.4f}")
17 print("\n✅ Status: Completed successfully!")
18 print("✅ Serverless: Compatible")
19 print("✅ Whitelist: Safe")
20 print("=" * 50)

```

```

=====
🚨 FRAUD DETECTION - FINAL SUMMARY
=====

📊 Dataset Information:
  Total rows: 6,351,193
  Training set: 4,445,579
  Testing set: 1,905,614

⚖️ Class Balance:
  Normal: 6,343,476 | Fraud: 7,717
  Ratio: 822.01:1

🎯 Model Performance:
  AUC-ROC: 0.9938
  AUC-PR: 0.6336
  Precision: 0.0284
  Recall: 0.9579
  F1-Score: 0.9782

✅ Status: Completed successfully!
✅ Serverless: Compatible
✅ Whitelist: Safe
=====

```

```

1 # รวม HTML Visualization กับผลโมเดล
2 from IPython.display import display, HTML
3
4 # แสดงผลโมเดล
5 display(HTML(f"""
6 <h2>Model Summary</h2>
7 <ul>
8 <li>AUC-ROC: {auc_roc:.4f}</li>
9 <li>AUC-PR: {auc_pr:.4f}</li>
10 <li>Precision: {precision:.4f}</li>

```

```
11 <li>Recall: {recall:.4f}</li>
12 <li>F1: {f1:.4f}</li>
13 </ul>
14 """))
15
16 # แทรก visualization จาก under.html
17 with open("/content/drive/MyDrive/Colab Notebooks/content/under.html", "r") as f:
18     html_vis = f.read()
19
20 display(HTML(html_vis))
```



```

In [6]: import pandas as pd
from imblearn.under_sampling import RandomUnderSampler
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, confusion_matrix, accu

# 1. โหลดข้อมูล
df = pd.read_csv(r'C:\Users\Asus F15\Downloads\transactions_train.csv', n

# 2. ตัวอย่างลดขนาดข้อมูล (ถ้าข้อมูลใหญ่มาก)
df_sample = df.sample(n=600000, random_state=42)

# 3. แยก X กับ y
X = df_sample.drop(columns=['isFraud'])
y = df_sample['isFraud']

# 4. Label Encoding (แปลง categorical เป็นตัวเลข)
categorical_cols = X.select_dtypes(include=['object']).columns
for col in categorical_cols:
    le = LabelEncoder()
    X[col] = le.fit_transform(X[col])

# 5. ทำ Undersampling
rus = RandomUnderSampler(random_state=42)
X_resampled, y_resampled = rus.fit_resample(X, y)

print("ข้อมูลก่อน undersampling:\n", y.value_counts())
print("ข้อมูลหลัง undersampling:\n", pd.Series(y_resampled).value_counts())

# 6. แบ่งข้อมูล train/test
X_train, X_test, y_train, y_test = train_test_split(
    X_resampled, y_resampled, test_size=0.2, random_state=42, stratify=y_

# 7. สร้างโมเดล Random Forest
model = RandomForestClassifier(n_estimators=100, random_state=42)

# 8. ฝึกโมเดล
model.fit(X_train, y_train)

# 9. ทำนายข้อมูล test
y_pred = model.predict(X_test)

# 10. ประเมินผลโมเดล
print("Accuracy:", accuracy_score(y_test, y_pred))
print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))

```

```

ข้อมูลก่อน undersampling:
isFraud
0      599639
1         361
Name: count, dtype: int64

```

```
ข้อมูลหลัง undersampling:
isFraud
0      361
1      361
Name: count, dtype: int64
Accuracy: 0.9448275862068966
```

```
Confusion Matrix:
[[68  5]
 [ 3 69]]
```

```
Classification Report:
              precision    recall  f1-score   support

     0       0.96      0.93      0.94         73
     1       0.93      0.96      0.95         72

 accuracy          0.94         145
 macro avg         0.95         145
weighted avg         0.95         145
```

```
In [7]: # 5. ทำ Undersampling
rus = RandomUnderSampler(random_state=42)
X_resampled, y_resampled = rus.fit_resample(X, y)

print("ขนาดข้อมูลก่อน Undersampling:")
print(y.value_counts())

print("ขนาดข้อมูลหลัง Undersampling:")
print(pd.Series(y_resampled).value_counts())
```

```
ขนาดข้อมูลก่อน Undersampling:
isFraud
0      599639
1         361
Name: count, dtype: int64
ขนาดข้อมูลหลัง Undersampling:
isFraud
0      361
1      361
Name: count, dtype: int64
```

```
In [11]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
    X_resampled, y_resampled, test_size=0.2, random_state=42
)
```

```
In [13]: !pip install xgboost
```

```
Collecting xgboost
  Downloading xgboost-3.0.5-py3-none-win_amd64.whl.metadata (2.1 kB)
Requirement already satisfied: numpy in c:\users\asus f15\anaconda3\lib\site-packages (from xgboost) (2.1.3)
Requirement already satisfied: scipy in c:\users\asus f15\anaconda3\lib\site-packages (from xgboost) (1.15.3)
Downloading xgboost-3.0.5-py3-none-win_amd64.whl (56.8 MB)
----- 0.0/56.8 MB ? eta -:--:--
----- 0.3/56.8 MB ? eta -:--:--
----- 0.8/56.8 MB 2.6 MB/s eta 0:00:
```

22

```
----- 1.0/56.8 MB 1.6 MB/s eta 0:00:
34 ----- 1.3/56.8 MB 1.6 MB/s eta 0:00:
35 ----- 1.8/56.8 MB 1.8 MB/s eta 0:00:
32 ----- 2.6/56.8 MB 2.1 MB/s eta 0:00:
26 ----- 2.6/56.8 MB 2.1 MB/s eta 0:00:
26 ----- 2.9/56.8 MB 1.7 MB/s eta 0:00:
32 ----- 3.4/56.8 MB 1.8 MB/s eta 0:00:
30 ----- 3.9/56.8 MB 1.9 MB/s eta 0:00:
28 ----- 4.5/56.8 MB 2.0 MB/s eta 0:00:
27 ----- 4.7/56.8 MB 2.0 MB/s eta 0:00:
26 ----- 5.0/56.8 MB 1.8 MB/s eta 0:00:
29 ----- 5.5/56.8 MB 1.9 MB/s eta 0:00:
28 ----- 6.3/56.8 MB 2.0 MB/s eta 0:00:
26 ----- 6.6/56.8 MB 2.0 MB/s eta 0:00:
25 ----- 6.6/56.8 MB 2.0 MB/s eta 0:00:
25 ----- 6.8/56.8 MB 1.8 MB/s eta 0:00:
28 ----- 7.3/56.8 MB 1.9 MB/s eta 0:00:
27 ----- 7.9/56.8 MB 1.9 MB/s eta 0:00:
27 ----- 7.9/56.8 MB 1.9 MB/s eta 0:00:
27 ----- 8.1/56.8 MB 1.8 MB/s eta 0:00:
27 ----- 8.4/56.8 MB 1.8 MB/s eta 0:00:
28 ----- 8.9/56.8 MB 1.8 MB/s eta 0:00:
28 ----- 8.9/56.8 MB 1.8 MB/s eta 0:00:
28 ----- 9.2/56.8 MB 1.7 MB/s eta 0:00:
29 ----- 9.4/56.8 MB 1.7 MB/s eta 0:00:
29 ----- 9.4/56.8 MB 1.7 MB/s eta 0:00:
29 ----- 10.0/56.8 MB 1.6 MB/s eta 0:0
0:29 ----- 10.5/56.8 MB 1.7 MB/s eta 0:0
0:28 ----- 10.7/56.8 MB 1.6 MB/s eta 0:0
0:29 ----- 11.3/56.8 MB 1.7 MB/s eta 0:0
0:28 ----- 11.5/56.8 MB 1.7 MB/s eta 0:0
0:27
```

```
----- 11.8/56.8 MB 1.7 MB/s eta 0:0
0:28
----- 11.8/56.8 MB 1.7 MB/s eta 0:0
0:28
----- 12.1/56.8 MB 1.6 MB/s eta 0:0
0:29
----- 12.3/56.8 MB 1.6 MB/s eta 0:0
0:29
----- 12.8/56.8 MB 1.6 MB/s eta 0:0
0:28
----- 13.6/56.8 MB 1.6 MB/s eta 0:0
0:27
----- 13.9/56.8 MB 1.7 MB/s eta 0:0
0:26
----- 13.9/56.8 MB 1.7 MB/s eta 0:0
0:26
----- 14.2/56.8 MB 1.6 MB/s eta 0:0
0:27
----- 14.9/56.8 MB 1.6 MB/s eta 0:0
0:26
----- 15.2/56.8 MB 1.7 MB/s eta 0:0
0:26
----- 15.2/56.8 MB 1.7 MB/s eta 0:0
0:26
----- 16.0/56.8 MB 1.6 MB/s eta 0:0
0:25
----- 16.8/56.8 MB 1.7 MB/s eta 0:0
0:24
----- 17.6/56.8 MB 1.7 MB/s eta 0:0
0:23
----- 18.4/56.8 MB 1.8 MB/s eta 0:0
0:22
----- 19.1/56.8 MB 1.8 MB/s eta 0:0
0:21
----- 19.4/56.8 MB 1.8 MB/s eta 0:0
0:21
----- 19.9/56.8 MB 1.8 MB/s eta 0:0
0:21
----- 20.4/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 20.4/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 21.0/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 21.2/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 22.0/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 22.0/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 22.3/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 22.5/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 22.8/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 23.1/56.8 MB 1.8 MB/s eta 0:0
0:20
----- 23.6/56.8 MB 1.8 MB/s eta 0:0
0:19
----- 24.4/56.8 MB 1.8 MB/s eta 0:0
0:19
```

U. L. J.