

Lesson 32 Summary: ML Pipeline

Data Science with Python – Key Concepts

Data Science Program

ML Pipeline

Pipeline Concept

Chain transformers
+ estimator
Reproducible workflow

Common Steps

Impute -> Scale ->
Encode -> Model
Sequential transforms

Benefits

No data leakage
Clean code
Easy deployment

ColumnTransformer

Different transforms for different cols
Numeric vs categorical handling

GridSearchCV Integration

Tune all params together
pipeline_step_param syntax

```
Pipeline([("scaler", StandardScaler()), ("model", LogisticRegression())])
```

Why Pipelines?

Benefits:

- **No data leakage:** Proper train/test separation
- **Reproducibility:** Single object captures workflow
- **Deployment:** Easy to save and load

Pipelines prevent common ML mistakes

Chain of steps:

- **Transformers:** Scale, encode, impute
- **Estimator:** Final model (must be last)
- **Sequence:** Each step transforms data for next

All steps except last must have transform method

Different transforms per column type:

- **Numeric:** Scale, impute with mean
- **Categorical:** One-hot encode, impute with mode
- **Parallel:** Apply transforms simultaneously

ColumnTransformer handles mixed data types

Common Pipeline Pattern

Typical workflow:

- **Step 1:** Handle missing values
- **Step 2:** Scale numeric features
- **Step 3:** Encode categoricals
- **Step 4:** Fit model

This pattern covers most ML preprocessing needs

Tune all parameters together:

- **Syntax:** step__parameter
- **Example:** model__C for regularization
- **Benefit:** Finds optimal preprocessing + model

Double underscore separates step name and parameter

Deployment-ready:

- **joblib:** Save entire pipeline
- **Load:** Full preprocessing + model
- **Predict:** Works on raw data

Saved pipeline includes all fitted transformers

Pipeline tips:

- **Name steps:** Use descriptive names
- **Test early:** Verify with small data
- **Document:** Record preprocessing choices

Clear naming makes debugging easier

Essential Commands:

Task	Code
Create pipeline	<code>Pipeline([('scaler', ...), ('model', ...)])</code>
Fit	<code>pipe.fit(X_train, y_train)</code>
Predict	<code>pipe.predict(X_test)</code>
Grid search	<code>GridSearchCV(pipe, param_grid)</code>
Save	<code>joblib.dump(pipe, 'model.pkl')</code>

Pipelines are essential for production ML