

# Detailed Breakdown: PySpark Preprocessing Pipeline

## 1. Initial Setup and Imports

python

```
from pyspark.ml.feature import OneHotEncoder, StringIndexer, VectorAssembler
```

**What these do:**

- **StringIndexer**: Converts text → numbers
- **OneHotEncoder**: Converts numbers → binary vectors
- **VectorAssembler**: Combines multiple columns → single feature vector

## 2. Define Categorical Columns

python

```
categoricalColumns = ['job', 'marital', 'education', 'default', 'housing', 'loan', 'contact', 'poutcome']
```

**What this is:** List of all text columns that need special processing

**Example data in these columns:**

```
job: "teacher", "student", "manager"  
marital: "single", "married", "divorced"  
education: "primary", "secondary", "tertiary"  
default: "yes", "no"  
housing: "yes", "no"  
loan: "yes", "no"  
contact: "cellular", "telephone", "unknown"  
poutcome: "success", "failure", "unknown"
```

## 3. Initialize Empty Stages List

python

```
stages = []
```

**Purpose:** This list will store all preprocessing steps in order **Think of it as:** A recipe where each step gets added one by one

## 4. The Loop - Processing Each Categorical Column

python

```
for categoricalCol in categoricalColumns:  
    stringIndexer = StringIndexer(inputCol = categoricalCol, outputCol = categoricalCol + "Index")  
    encoder = OneHotEncoder(inputCols=[stringIndexer.getOutputCol()], outputCols=[categoricalCol + "classVec"])  
    stages += [stringIndexer, encoder]
```

**Let's trace through this loop step by step:**

**Iteration 1: categoricalCol = 'job'**

**Step 1a: Create StringIndexer**

python

```
stringIndexer = StringIndexer(inputCol = 'job', outputCol = 'job' + 'Index')  
# Result: StringIndexer(inputCol='job', outputCol='jobIndex')
```

**What it will do:**

Input column 'job': ["teacher", "student", "manager", "teacher"]  
Output column 'jobIndex': [0, 1, 2, 0]

**Step 1b: Create OneHotEncoder**

python

```
encoder = OneHotEncoder(inputCols=['jobIndex'], outputCols=['job' + "classVec"])  
# Result: OneHotEncoder(inputCols=['jobIndex'], outputCols=[jobclassVec])
```

**What it will do:**

Input column 'jobIndex': [0, 1, 2, 0]  
Output column 'jobclassVec': [[1,0,0], [0,1,0], [0,0,1], [1,0,0]]

**Step 1c: Add to stages**

```
python
```

```
stages += [stringIndexer, encoder]
# stages now contains: [StringIndexer(job→jobIndex), OneHotEncoder(jobIndex→jobclassVec)]
```

## Iteration 2: categoricalCol = 'marital'

### Step 2a: Create StringIndexer

```
python
```

```
stringIndexer = StringIndexer(inputCol = 'marital', outputCol = 'marital' + 'Index')
# Result: StringIndexer(inputCol='marital', outputCol='maritalIndex')
```

### Step 2b: Create OneHotEncoder

```
python
```

```
encoder = OneHotEncoder(inputCols=['maritalIndex'], outputCols=['marital' + "classVec"])
# Result: OneHotEncoder(inputCols=['maritalIndex'], outputCols=['maritalclassVec'])
```

### Step 2c: Add to stages

```
python
```

```
stages += [stringIndexer, encoder]
# stages now contains: [
#   StringIndexer(job→jobIndex),
#   OneHotEncoder(jobIndex→jobclassVec),
#   StringIndexer(marital→maritalIndex),
#   OneHotEncoder(maritalIndex→maritalclassVec)
# ]
```

**This continues for all 8 categorical columns...**

**After the complete loop, stages contains 16 items:**

```
python
```

```
stages = [
    StringIndexer(job → jobIndex),
    OneHotEncoder(jobIndex → jobclassVec),
    StringIndexer(marital → maritalIndex),
    OneHotEncoder(maritalIndex → maritalclassVec),
    StringIndexer(education → educationIndex),
    OneHotEncoder(educationIndex → educationclassVec),
    StringIndexer(default → defaultIndex),
    OneHotEncoder(defaultIndex → defaultclassVec),
    StringIndexer(housing → housingIndex),
    OneHotEncoder(housingIndex → housingclassVec),
    StringIndexer(loan → loanIndex),
    OneHotEncoder(loanIndex → loanclassVec),
    StringIndexer(contact → contactIndex),
    OneHotEncoder(contactIndex → contactclassVec),
    StringIndexer(poutcome → poutcomeIndex),
    OneHotEncoder(poutcomeIndex → poutcomeclassVec)
]
```

## 5. Handle the Target Variable

```
python
```

```
label_stringIdx = StringIndexer(inputCol = 'deposit', outputCol = 'label')
stages += [label_stringIdx]
```

**Purpose:** Convert target variable 'deposit' (yes/no) to numbers (0/1) **Why no OneHotEncoder?**: Target variables don't need one-hot encoding

### Transformation:

```
deposit: ["yes", "no", "yes", "no"]
label: [0, 1, 0, 1] # or [1, 0, 1, 0] depending on frequency
```

**stages now has 17 items** (16 + 1)

## 6. Define Numeric Columns

```
python
```

```
numericCols = ['age', 'balance', 'duration', 'campaign', 'pdays', 'previous']
```

**What this is:** Columns that are already numbers, no conversion needed

### Example data:

```
age: [25, 35, 45]  
balance: [1200, 2500, 5000]  
duration: [180, 300, 420]
```

## 7. Prepare Inputs for VectorAssembler

```
python
```

```
assemblerInputs = [c + "classVec" for c in categoricalColumns] + numericCols
```

Let's break this down:

### Part 1: [c + "classVec" for c in categoricalColumns]

This is a **list comprehension** that creates:

```
python
```

```
categoricalColumns = ['job', 'marital', 'education', 'default', 'housing', 'loan', 'contact', 'poutcome']
```

# The comprehension does:

```
[c + "classVec" for c in categoricalColumns]
```

# Which expands to:

```
['job' + 'classVec', 'marital' + 'classVec', 'education' + 'classVec', ...]
```

# Final result:

```
['jobclassVec', 'maritalclassVec', 'educationclassVec', 'defaultclassVec',  
'housingclassVec', 'loanclassVec', 'contactclassVec', 'poutcomeclassVec']
```

### Part 2: + numericCols

This concatenates the lists:

```

python

assemblerInputs = ['jobclassVec', 'maritalclassVec', 'educationclassVec', 'defaultclassVec',
                   'housingclassVec', 'loanclassVec', 'contactclassVec', 'poutcomeclassVec'] +
                  ['age', 'balance', 'duration', 'campaign', 'pdays', 'previous']

# Final result:
assemblerInputs = ['jobclassVec', 'maritalclassVec', 'educationclassVec', 'defaultclassVec',
                   'housingclassVec', 'loanclassVec', 'contactclassVec', 'poutcomeclassVec',
                   'age', 'balance', 'duration', 'campaign', 'pdays', 'previous']

```

## 8. Create VectorAssembler

```

python

assembler = VectorAssembler(inputCols=assemblerInputs, outputCol="features")
stages += [assembler]

```

**Purpose:** Combine all processed features into a single vector column

**What it will do:**

```

python

# Input columns:
jobclassVec: [[1,0,0], [0,1,0], [1,0,0]]
maritalclassVec: [[1,0], [0,1], [1,0]]
age: [25, 35, 28]
balance: [1200, 2500, 1800]

# Output column 'features':
[[1,0,0, 1,0, 25, 1200], [0,1,0, 0,1, 35, 2500], [1,0,0, 1,0, 28, 1800]]

```

**stages now has 18 items** (17 + 1)

## 9. Complete Pipeline Summary

```

python

stages = [
    # Process categorical columns (16 stages)
    StringIndexer(job → jobIndex),
    OneHotEncoder(jobIndex → jobclassVec),
    StringIndexer(marital → maritalIndex),
    OneHotEncoder(maritalIndex → maritalclassVec),
    # ... 6 more pairs ...

    # Process target variable (1 stage)
    StringIndexer(deposit → label),

    # Combine all features (1 stage)
    VectorAssembler([all_encoded_features + numeric_cols] → features)
]

```

## 10. Data Transformation Flow

### Original Data:

age	job	marital	deposit
25	teacher	single	yes
35	student	married	no

### After StringIndexers:

age	jobIndex	maritalIndex	label
25	0	0	1
35	1	1	0

### After OneHotEncoders:

age	jobclassVec	maritalclassVec	label
25	[1,0,0]	[1,0,0]	1
35	[0,1,0]	[0,1,0]	0

### After VectorAssembler:

features	label
[25, 1,0,0, 1,0,0]	1
[35, 0,1,0, 0,1,0]	0

## Key Insights

1. **Loop Efficiency:** Instead of writing 16 separate transformers, one loop handles all categorical columns
2. **Consistent Naming:** All encoded columns end with "classVec" for easy identification
3. **Order Matters:** StringIndexer must come before OneHotEncoder
4. **Target Separate:** Target variable ('deposit') gets different treatment than features
5. **Vector Assembly:** All features must be combined into one column for ML algorithms

This preprocessing pipeline converts messy real-world data into clean numerical input that machine learning algorithms can understand!