

Contents

1 PDF Section Extractor - Metadata Enhancement	1
1.1 Overview	1
1.2 New Fields	2
1.2.1 1. <code>empirical_page_number</code> (Integer, Nullable)	2
1.2.2 2. <code>section_name</code> (String, Nullable)	2
1.3 Page Number Line Exclusion	2
1.4 Updated Schema	3
1.5 Usage Examples	3
1.5.1 Basic Extraction	3
1.5.2 Query by Section Name	3
1.5.3 Compare PDF vs Empirical Page Numbers	4
1.5.4 Filter by Empirical Page Number	4
1.5.5 Section Analysis	4
1.5.6 Extract Specific Sections	5
1.6 Implementation Details	5
1.6.1 New Helper Methods	5
1.6.2 Two-Pass Processing	6
1.7 Benefits	6
1.7.1 1. Enhanced Metadata	6
1.7.2 2. Improved Filtering	7
1.7.3 3. Better Analytics	7
1.7.4 4. Cleaner Output	7
1.8 Backward Compatibility	7
1.9 Testing	8
1.10 See Also	8

1 PDF Section Extractor - Metadata Enhancement

1.1 Overview

The PDFSectionExtractor has been enhanced with two new metadata fields: 1. `empirical_page_number`: Page numbers extracted from the document itself 2. `section_name`: Standardized section names for academic papers

Additionally, page number lines are now excluded from the extracted content.

1.2 New Fields

1.2.1 1. empirical_page_number (Integer, Nullable)

Extracts the actual page number from the document by analyzing the first and last lines of each PDF page.

Examples: - "485" → 485 - "486 ... Wang" → 486 - "485-489" → 485 (first number in range)

How it works: 1. For each PDF page, examines the first 5 and last 5 lines 2. Looks for numeric patterns matching page numbers (1-9999) 3. Maps PDF page numbers to empirical page numbers

Example mapping:

```
PDF Page 1 → Empirical Page 108  
PDF Page 2 → Empirical Page 486  
PDF Page 3 → Empirical Page 487  
PDF Page 4 → Empirical Page 488  
PDF Page 5 → Empirical Page 489
```

1.2.2 2. section_name (String, Nullable)

Identifies and standardizes academic section names from headers.

Recognized sections: - Introduction - Abstract - Keywords - Methods / Materials and Methods - Results - Discussion - Conclusion - Taxonomy - Description - Etymology - Holotype / Paratype / Specimen - Acknowledgments - References / Literature Cited - Background / Objectives / Summary - Figures / Tables / Appendix / Supplementary

How it works: 1. When a header is detected, checks if it matches known section keywords 2. Assigns standardized section name to that header 3. Subsequent paragraphs inherit the current section name 4. Section name persists until a new section header is encountered

Example:

```
Row 7: "Abstract – This paper..." → section_name = "Abstract"  
Row 8: "collected from Taiwan..." → section_name = "Abstract" (inherited)  
Row 9: "Key words – discomycetes..." → section_name = "Keywords"  
Row 10: "Introduction" → section_name = "Introduction"  
Row 11: "In Korf's monograph..." → section_name = "Introduction" (inherited)
```

1.3 Page Number Line Exclusion

Page number lines are now automatically detected and excluded from the extracted content.

Patterns detected: - Simple numbers: "485", "486" - With ellipsis: "486 ... Wang" - Page ranges: "485-489" - Formatted: "Page 485", "p. 486"

Benefits: - Cleaner extracted text - No standalone page number paragraphs - Better paragraph continuity

1.4 Updated Schema

```
StructType([
    StructField("value", StringType(), False),
    StructField("doc_id", StringType(), False),
    StructField("attachment_name", StringType(), False),
    StructField("paragraph_number", IntegerType(), False),
    StructField("line_number", IntegerType(), False),
    StructField("page_number", IntegerType(), False),           # PDF page
    StructField("empirical_page_number", IntegerType(), True), # NEW - Nullable
    StructField("section_name", StringType(), True)            # NEW - Nullable
])
```

1.5 Usage Examples

1.5.1 Basic Extraction

```
from pyspark.sql import SparkSession
from pdf_section_extractor import PDFSectionExtractor

spark = SparkSession.builder.appName("PDFExtractor").getOrCreate()
extractor = PDFSectionExtractor(spark=spark)

sections_df = extractor.extract_from_document(
    database='skol_dev',
    doc_id='document-id'
)

# View all fields
sections_df.show()
```

1.5.2 Query by Section Name

```
# Get all Introduction paragraphs
intro = sections_df.filter(sections_df.section_name == "Introduction")
intro.select("paragraph_number", "value").show()

# Get all rows with section names
```

```
sections_only = sections_df.filter(sections_df.section_name.isNotNull())
sections_only.groupBy("section_name").count().show()
```

1.5.3 Compare PDF vs Empirical Page Numbers

```
# See the mapping
sections_df.select("page_number", "empirical_page_number") \
    .distinct() \
    .orderBy("page_number") \
    .show()

# Output:
# +-----+-----+
# |page_number|empirical_page_number|
# +-----+-----+
# | 1 | 108 |
# | 2 | 486 |
# | 3 | 487 |
# +-----+-----+
```

1.5.4 Filter by Empirical Page Number

```
# Get content from document page 486
page_486 = sections_df.filter(sections_df.empirical_page_number == 486)
page_486.select("paragraph_number", "section_name", "value").show()
```

1.5.5 Section Analysis

```
# Count paragraphs per section
sections_df.groupBy("section_name") \
    .count() \
    .orderBy("count", ascending=False) \
    .show()

# Output:
# +-----+-----+
# | section_name|count|
# +-----+-----+
# | NULL | 6 | (not in a named section)
# | Etymology | 4 |
# | Holotype | 3 |
# | Abstract | 2 |
# +-----+-----+
```

Get average paragraph length per section

```

from pyspark.sql.functions import length, avg

sections_df.groupBy("section_name") \
    .agg(avg(length("value")).alias("avg_length")) \
    .orderBy("avg_length", ascending=False) \
    .show()

```

1.5.6 Extract Specific Sections

```

# Extract just the abstract
abstract = sections_df.filter(sections_df.section_name == "Abstract") \
    .orderBy("paragraph_number") \
    .select("value") \
    .collect()

abstract_text = " ".join([row.value for row in abstract])
print(f"Abstract: {abstract_text}")

# Extract methods section
methods = sections_df.filter(
    (sections_df.section_name == "Methods") | 
    (sections_df.section_name == "Materials and Methods")
)

```

1.6 Implementation Details

1.6.1 New Helper Methods

1.6.1.1 _extract_empirical_page_number(lines_buffer: List[str]) -> Optional[int] Searches for page numbers in a list of lines.

```

first_lines = ["MYCOTAXON", "Volume 108, pp. 485–489", ...]
last_lines = [...continued", "486", ""]

```

```

page_num = extractor._extract_empirical_page_number(first_lines + last_lines)
# Returns: 486

```

1.6.1.2 _is_page_number_line(line: str) -> bool Checks if a line appears to be a standalone page number.

```

extractor._is_page_number_line("486")                      # True
extractor._is_page_number_line("486 ... Wang")            # True
extractor._is_page_number_line("Page 486")                 # True
extractor._is_page_number_line("This is text")             # False

```

1.6.1.3 `_get_section_name(header_text: str) -> Optional[str]`

Extracts standardized section name from header text.

```
extractor._get_section_name("INTRODUCTION")           # "Introduction"
extractor._get_section_name("Materials and Methods") # "Materials and Methods"
extractor._get_section_name("Key words")             # "Keywords"
extractor._get_section_name("Some Other Header")      # None
```

1.6.2 Two-Pass Processing

The parser now uses a two-pass approach:

Pass 1: Extract empirical page numbers

```
# Identify page boundaries from PDF markers
page_boundaries = [(1, 0, 150), (2, 151, 300), ...]

# Extract empirical page numbers for each PDF page
empirical_page_map = {1: 108, 2: 486, 3: 487, ...}
```

Pass 2: Parse sections with metadata

```
# Parse sections while tracking:
# - Current section name
# - Empirical page number from map
# - Skip page number lines

for line in lines:
    if is_page_marker(line):
        update_current_page()
    elif is_page_number_line(line):
        skip() # Don't include in output
    elif is_header(line):
        update_section_name()
        add_record_with_metadata()
    # ...
```

1.7 Benefits

1.7.1 1. Enhanced Metadata

- Track actual document page numbers vs PDF page numbers
- Know which section each paragraph belongs to
- Better document structure understanding

1.7.2 2. Improved Filtering

```
# Get only Introduction and Methods sections
key_sections = sections_df.filter(
    sections_df.section_name.isin(["Introduction", "Methods"])
)

# Get content from specific document pages
pages_2_3 = sections_df.filter(
    sections_df.empirical_page_number.between(486, 487)
)
```

1.7.3 3. Better Analytics

```
# Analyze section distribution
sections_df.groupBy("section_name", "empirical_page_number") \
    .count() \
    .show()

# Compare content across sections
sections_df.groupBy("section_name") \
    .agg(
        count("*").alias("paragraphs"),
        avg(length("value")).alias("avg_paragraph_length")
    ) \
    .show()
```

1.7.4 4. Cleaner Output

- No standalone page number lines in extracted content
- More coherent paragraph flow
- Better text quality for downstream processing

1.8 Backward Compatibility

Fully compatible with existing code: - All previous fields still present
- New fields are nullable (won't break existing queries) - API remains unchanged

Schema evolution:

```
# Old schema (6 fields)
# New schema (8 fields - added 2)

# Existing code still works
sections_df.select("value", "doc_id", "page_number").show()
```

```
# New functionality available
sections_df.select("section_name", "empirical_page_number").show()
```

1.9 Testing

All tests pass with new functionality:

```
$ python pdf_section_extractor.py
✓ DataFrame created successfully
✓ Schema correct (8 columns)
✓ 27 sections extracted
✓ Empirical page numbers extracted: {1: 108, 2: 486, 3: 487, 4: 488, 5: 489}
✓ Section names identified: Abstract, Keywords, Introduction, etc.
✓ Page number lines excluded from output

$ python test_new_fields.py
✓ All fields populated correctly
✓ Section name tracking works
✓ Empirical page numbers mapped correctly
✓ Page number lines excluded
```

1.10 See Also

- PDF_DATAFRAME_OUTPUT.md - DataFrame features
- PDF_DATAFRAME_MIGRATION_SUMMARY.md - Migration guide
- pdf_section_extractor.py - Implementation
- test_new_fields.py - Test examples

Update Date: 2025-12-22 **Status:**  Complete and tested **Breaking Changes:** None (backward compatible) **New Features:** Section names, empirical page numbers, page number exclusion