

# **Address Linkage Key with AnalyticsIQ Milestone 2 Presentation**

## **Project Owners :**

**Warren Smith**

**Jin Wang**

## **Project Team :**

**Kat McElveen, Team Lead**

**Saritha Gudala, Technical Writer**

**Kat Greer, Technical Specialist**

**Byron Smith, Technical Specialist**

**IT7993 Capstone, Spring 2021**

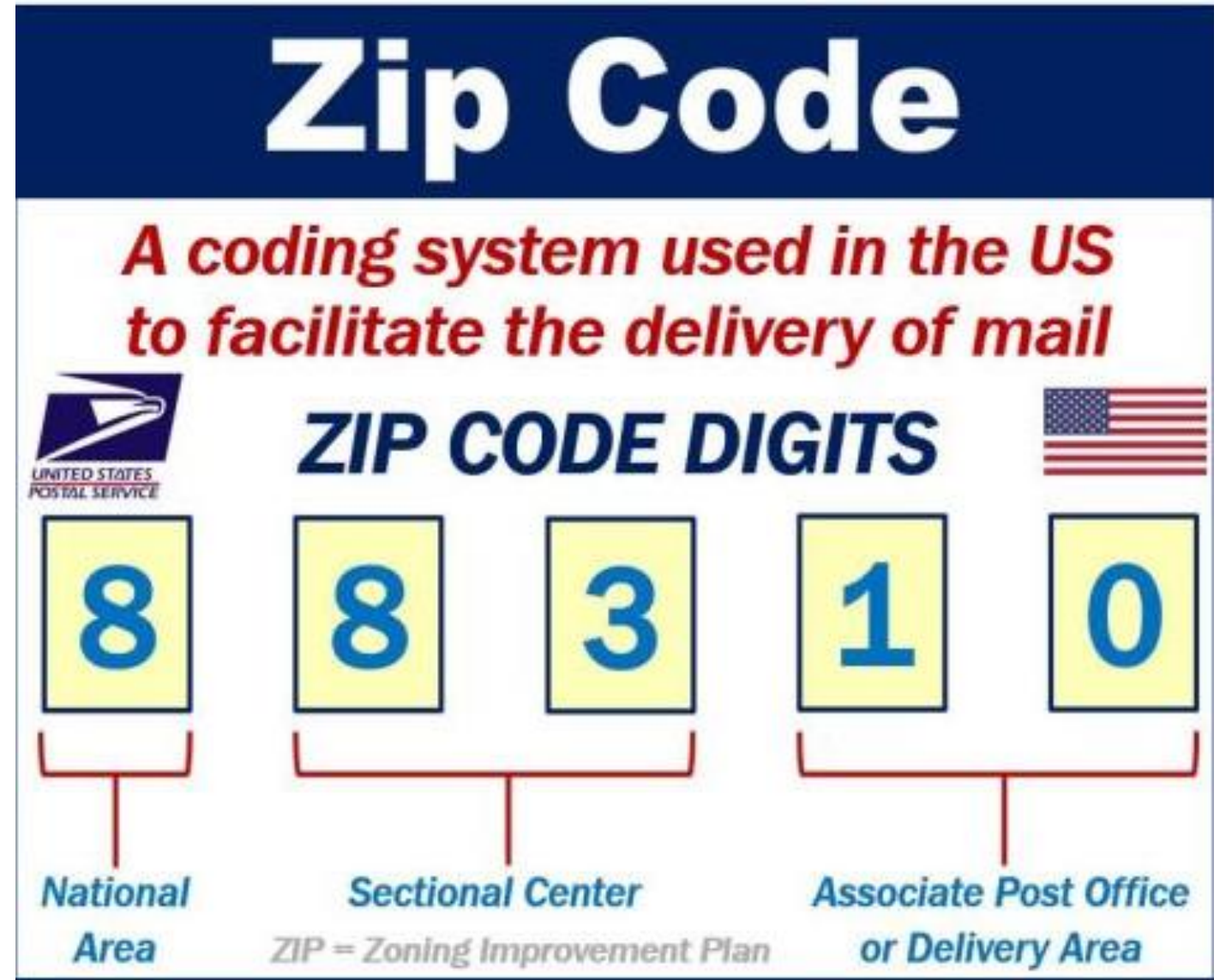
**March 26, 2021**

# Project Overview

- Develop a tool/set of tools to compute the delivery point (DPBC) of a given address following the guidelines of USPS.
- The tool will be initially developed and tested for properly formatted addresses.
- As part of further improvements, it will be tested for improper address format and/or for the presence of certain types of delivery points such as high-rise buildings.

## What is Minimally Viable Product?

- A tool or set of tools which allow me to compute the delivery point code of a given address given:
  - 1) A “clean” address (as defined by publication 28)
  - 2) ZIP and ZIP+4 appended
  - 3) A dataset which contains Highrise addresses



# Project Progress Summary

- Project team continued weekly meetings (Wednesday and Friday) to assess the progress and discuss further plans.
- As a major part of project deals with regular expressions, team members studied and applied regular expressions for implementation.
- Completed the development part of minimally viable product.
  - Team initially self-assigned tasks according to DPBC rules
  - After team discussions, some tasks were combined or discarded because they were redundant, unnecessary, or overly burdensome
  - Team utilized Jupyter Notebooks for local testing
  - Functions and Test Cases were created for the applicable rules
  - Team used Git (via TortoiseGit and command line) to commit and push local edits, which were then merged into the master branch upon approval of Project Owner
- For next milestone, we plan to extend the implementation of tool for secondary address rules (addresses with apartment numbers, PO Boxes, rural route addresses)

# Milestone 2 Deliverables Status Update

## **Complete:**

- Finished the implementation of minimally viable product.
- Currently we can generate a DBPC for all scenarios stated in slide 6 and 7 (excluding 4 and 11)

## **In Progress:**

- CI/CD Integration

## **Changed:**

- Eliminated rules 4 and 11
- Combined rules 2, 7, and 13

## Delivery Point Barcode Rules (Primary)

<p><b>1. General Rule</b></p> <p>Address: 1234 MAIN ST (PO BOX 44, RR 1 BOX 154, HC 1 BOX 1264)  DPBC: 34 (44, 54, 64)</p> <p>Use last two digits. Print code characters in DPBC representing last two digits of primary street number (or post office box, rural route box, or highway contract route number).</p>	<p><b>8. Leading/Embedded Alphas</b></p> <p>Address: 23S41 MAIN ST (23S4 MAIN ST, 2W3S1 MAIN ST, MAINS ST, C8INT)  DPBC: 11 (04, 01, 01)</p> <p>Print code characters in DPBC representing last two digits to right of alphas. If single digit to right of alphas, add leading zero.</p>
<p><b>2. No Numbers</b></p> <p>Address: MAIN St (RR 1, HC 1)  DPBC: 99 (99, 99)</p> <p>Use 99. Print code characters in DPBC representing last two digits of primary street number (or PO Box, rural route, or highway contract route number).</p>	<p><b>9. Slashes (/)</b></p> <p>Address: 123/4 MAIN ST (PO BOX ¼, RR 1 BOX 123/124/125, H 3 BOX 11/13)  DPBC: 23 (03, 23, 07)</p> <p>Print code characters in DPBC representing 99 whenever a slash appears directly next to numeric in the primary street number.</p>
<p><b>3. Single Digits</b></p> <p>Address: 8 MIAN St (PO BOX 1, RR 1 BOX 2, HC 1 BOX 3)  DPBC: 08 (01, 02, 03)</p> <p>Add leading zero. Print code characters in DPBC representing leading zero and single digit.</p>	<p><b>10. Other Embedded Symbols</b></p> <p>Address: 1.23 MAIN ST (PO BOX 1-3, RR 1 BOX 1.23, HC 3 BOX 11*7)  DPBC: 23 (03, 23, 07)</p> <p>Use last two digits to right of the symbol. Print code characters in DPBC representing last two digits to the right of all symbols (except slashes), such as periods and hyphens appearing in primary street numbers. If single digit to right, add leading zero.</p>
<p><b>4. Fractional Number</b></p> <p>Address: 1234 ½ MAIN ST (PO BOX 1 ½, RR 1 BOX 2 ¾, HC 1 BOX 10 ¼)  DPBC: 34 (01, 02, 10)</p> <p>Ignore fraction. Print code characters in DPBC representing two digits to left of fraction. If single digit to left of fraction, add leading zero.</p>	<p><b>11. Embedded Spaces</b></p> <p>Address: 1 23 MAIN ST (PO BOX 1 3, RR 1 BOX 1 7, HC 1 BOX 12 34)  DPBC: 23 (03, 07, 34)</p> <p>Treat embedded spaces like other symbols (Rule 10). Print code characters in DPBC representing last two digits to right of space. If single digit to right, add leading zero.</p>

<p><b>5. Trailing Alphas</b></p> <p>Address: 1234A MAIN ST (PO BOX 4A, RR 1 BOX 154A, HC 1 BOX 12644AA) DPBC: 34 (04, 54, 44)</p> <p>Ignoring trailing alphas. Print code characters in DPBC representing last two digits to left of space and alphas. If single digit to left of space and alphas, add leading zero.</p>	<p><b>12. Numeric Street Names</b></p> <p>Address: 8 33 ST (123 7<sup>th</sup> ST) DPBC: 08 (23)</p> <p>Ignore numeric street name. Print code characters in DPBC representing last two digits of primary street number (Rule 1).</p>
<p><b>6. Spaces and Alphas</b></p> <p>Address: 1234 A MAIN ST (PO BOX 4A, RR 1 BOX 154A, HC1 BOX 12644AA) DPBC: 34 (04, 54, 44)</p> <p>Ignoring space and alphas. Print code characters in DPBC representing two digits to left of space and alphas. If single digit to left of space and alphas, add leading zero.</p>	<p><b>13. All other Anomalies</b></p> <p>Use 99. Print code characters in DPBC representing 99 for conditions not covered by Rules 1 – 12.</p>
<p><b>7. Alphas Only</b></p> <p>Address : A Main St (PO Box AA, RR 1 Box X, HC 1 Box AB) DPBC: 99 (99, 99, 99)</p> <p>Ignore alphas and se 99. Print code characters in DPBC representing 99 when alphas appear as the only primary street number.</p>	

# Task Completion Summary

## Completed

- **General Rule**
- **No Numbers Rule**
- **Single Digits Rule**
- **Trailing Alphas Rule**
- **Spaces and Alphas Rule**
- **Alphas Only Rule**
- **Leading/Embedded Alphas**
- **All other Anomalies**

## In Progress

- **Other Embedded Symbols Rule- Testing Phase**
- **Slashes Rule (/)-Testing Phase**



# Task Status Update from GitLab

The screenshot displays the GitLab web interface for the 'Address Linkage Key' project. The top navigation bar includes the GitLab logo, 'Projects', 'Groups', and 'More' dropdowns, along with a search bar and user profile icons. The left sidebar contains a navigation menu with options: Project overview, Repository, Issues (11), Boards (selected), Labels, Service Desk, Milestones, Merge Requests (0), CI / CD, Operations, Packages & Registries, Analytics, Wiki, Snippets, and Members.

The main content area shows the 'Issue Boards' view for the 'Address Linkage Key' project. The board is organized into four columns: 'To-do' (2 items), 'Doing' (2 items), 'Done' (7 items), and 'Closed' (7 items). Each column has a header with a status label, a count, and icons for adding, deleting, and settings. The 'To-do' column contains 'Secondary DPC Rules (page 37)' (#17) and 'Compute Delivery Point Check Digit (page 45)' (#16). The 'Doing' column contains 'Single Digits Rule (#3)' (#5) and 'Leading/Embedded Alphas Rule (#8)' (#10). The 'Done' column contains 'Create Controller Script' (#1), 'All other Anomalies Rule (#13)' (#15), 'House Number (#0)' (#18), 'No Numbers Rule (#2)' (#4), 'Trailing Alphas Rule (#5)' (#7), 'General Rule (#1)' (#3), and 'Other Embedded Symbols Rule (#10)' (#12). The 'Closed' column contains 'Slashes Rule (#9)' (#11) with 'Doing' and 'Won't Fix' labels, 'Numeric Street Names Rule (#12)' (#14), 'Alphas Only Rule (#7)' (#9), 'Spaces and Alphas Rule (#6)' (#8), 'Embedded Spaces Rule (#11)' (#13), 'Fractional Number Rule (#4)' (#6) with a 'To-do' label, and 'Create "Source of Truth" Data Samples' (#2) with a 'Doing' label.

At the bottom left, there is a 'Collapse sidebar' button. The bottom right corner shows a scroll bar.

## Handler Function

```
15  
16 def handler(df):  
17     df = house_number_extract(df)  
18     df = apply_rule_1(df)  
19     df = apply_rule_5(df)  
20     df = apply_rule_8(df)  
21     df = apply_rule_10(df)  
22     df = apply_rule_13(df)  
23     #apply rule 3 (single digit rule) last  
24     df = apply_rule_3(df)  
25  
26     return df  
27
```

## Rule Functions

```
32 def apply_rule_1(df):
33     # Use last two digits. Print code characters in DPBC representing last two digits of primary street number
34     df = df.withColumn('dpc',
35         f.when(
36             f.col('address_line_2').isNull() &
37             f.col(' housenumber').isNotNull() &
38             f.col(' housenumber').rlike('^[0-9]*$'),
39             f.col(' housenumber').substr(-2,2))
40     return df
41
42 def apply_rule_2(df):
43     # Use 99 when address contains no house number
44     df = df.withColumn('dpc',f.when((f.col('dpc').isNull()) & (f.col(' housenumber').isNull()),f.lit('99'))
45         .otherwise(f.col('dpc')))
46     return df
47
48 def apply_rule_3(df):
49     # Add a leading 0 when address contains a single digit house number
50
51     # left pads dpc column with 0 up to length 2
52     df = df.withColumn('dpc', f.lpad('dpc', 2, '0'))
53     return df
54
55 def apply_rule_5(df):
56     #Ignoring trailing alphas. Print code characters in DPBC representing last two digits to left of space and alphas
57     df = df.withColumn('dpc',
58         f.when(
59             f.col('address_line_2').isNull() &
60             f.col(' housenumber').isNotNull() &
61             f.col('dpc').isNull() &
62             #f.col(' housenumber').rlike('^[a-zA-Z0-9]*$'),
63             f.col(' housenumber').rlike('^[0-9]+[A-Z]+$'),
64             f.regexp_extract(f.col(' housenumber'), '(\d+)', 1).substr(-2,2)).otherwise(f.col('dpc'))
```

## Test Case Classes

```
126 class TestRule1(SparkTestCase):
127     def test(self):
128         # Test Case: Rule 1 (General Rule) should provide the last two digits of a primary street number
129         # run a query that gives me the input that the function expects and the expected value from that
130         # parse operation
131         testdf = self.spark.sql("SELECT '1234 MAIN ST' AS address_line_1, CAST(NULL AS string) AS address_line_2, '1234' AS
132         housenumber, CAST(NULL AS string) AS dpc, '34' AS expected")
133         #run the function
134         testdf = apply_rule_1(testdf)
135         #gather the results [row 0 only] and compare the returned dpc to the expected dpc
136         row = testdf.collect()[0]
137         self.assertEqual(row.expected, row.dpc, "rule 1 created match")
138
139 class TestRule2(SparkTestCase):
140     def test(self):
141         # Test Case: Rule 2 should use 99 when the address contains no house number.
142         # run a query that gives me the input that the function expects and the expected value from that
143         # parse operation
144         testdf = self.spark.sql("SELECT 'MAIN ST' AS address_line_1, CAST(NULL AS string) AS housenumber, CAST(NULL AS string)
145         AS dpc, '99' AS expected")
146         #run the function
147         testdf = apply_rule_2(testdf)
148         #gather the results [row 0 only] and compare the returned dpc to the expected dpc
149         row = testdf.collect()[0]
150         self.assertEqual(row.expected, row.dpc, "rule 2 fail; expected does not match dpc")
151
152 class TestRule3(SparkTestCase):
153     def test(self):
154         # Test Case: Rule 3 should add a leading 0 when dpc is a single digit
155         # run a query that gives me the input that the function expects and the expected value from that
156         # parse operation
157         testdf = self.spark.sql("SELECT '8 MAIN ST' AS address_line_1, '8' AS dpc, '08' AS expected")
158         #run the function
159         testdf = apply_rule_3(testdf)
160         #gather the results [row 0 only] and compare the returned dpc to the expected dpc
161         row = testdf.collect()[0]
162         self.assertEqual(row.expected, row.dpc, "rule 3 fail; expected does not match dpc")
```

# Running Test Cases

To run test cases, run the module containing the test classes (address\_link/\_\_init\_\_.py file) in a terminal. The script will connect to a local spark instance and run the test cases.

- Open a terminal in Jupyter Notebook
- cd to work/address-linkage-key directory
- Run all tests with either of these commands
  - *python3 address\_link/\_\_init\_\_.py*
  - *make test (looks for all .py files and runs unit tests)*
- Run an individual test with this command
  - *python3 address\_link/\_\_init\_\_.py TestRule1.test*

## Running Test Cases

```
(base) jovyan@f49c99065829:~$ cd work/address-linkage-key
(base) jovyan@f49c99065829:~/work/address-linkage-key$ python3 address_link/__init__.py
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.spark.unsafe.Platform (file:/usr/local/spark-3.0.1-bin-hadoop3.2/jars/
-unsafe_2.12-3.0.1.jar) to constructor java.nio.DirectByteBuffer(long,int)
WARNING: Please consider reporting this to the maintainers of org.apache.spark.unsafe.Platform
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
21/03/24 01:17:56 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java c
s where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
.....F
=====
FAIL: test (__main__.TestRule8)
-----
Traceback (most recent call last):
  File "address_link/__init__.py", line 184, in test
    self.assertEqual(row.expected, row.dpc, "Rule 8 fails; expected does not match dpc")
AssertionError: '41' != None : Rule 8 fails; expected does not match dpc
-----
Ran 7 tests in 12.576s

FAILED (failures=1)
(base) jovyan@f49c99065829:~/work/address-linkage-key$
```

# Learnings

- As our whole project deals with generating DPBC (delivery point barcode) in multiple scenarios for different format of addresses. Regex functions (regular expressions) are best way to deal with different address formats of and generate a DPBC.
- As a team we all got chance to learn and implement regular expressions.
- Exposure to the wonderful world of DevOps.
  - Example of Continuous Integration (CI): configure GitLab to run unit tests upon commit.
  - Example of Continuous Delivery (CD): if all tests are okay, ship to production.

# Regular Expressions

- A regular is a special text string for describing a search pattern. Regular expressions commonly referred to as regex, regexp, or re.
- The re module is used to write regular expressions (regex) in Python. To load this module, we need to use the import statement. The following line of code is necessary to include at the top of your code

```
import re
```

- Example Regex: ('^[0-9]+[A-Z]+\$').The sample regular expression extracts only numbers from the specified column.

Example Address: 1234A MAIN ST

Extracts New Address: 1234 (Ignores Alpha)



# Common Syntax(Regex)

Table 1. Common Regular Expression Syntax

Syntax	Description
.	Matches any one character
^	Anchor; matches from the start of a string
\$	Anchor; matches at the end of a string
\	Escape character
	Pipe Character OR; C T will match C or T
*	Matches zero or more repetitions of the previous character
+	Matches one or more repetitions of the previous character
?	Matches zero or one repetitions of the previous character
{n}	Quantifier; matches n repetitions of the previous character
{n, x}	Quantifier; matches from n to x repetitions of the previous character
[ ]	Character group; e.g. [AGCT] will match the characters AGCT
[^ ]	Negated character group e.g. [^AGCT] will match any characters not in this group
( )	Matches the pattern specified in the parentheses exactly

# Challenges

- As part of initial project kick off, we have aimed to develop a tool for all 13 rules (more details about rules are included in slide 6 & 7).
- The rule numbers 4 (Fractional number) and 11(Embedded Spaces) are excluded from the scope of project development due to limitations in real time scenarios and current pyspark framework as advised by project owners.

# Gantt Chart(Screenshot)

[illegible]

# Milestone 3 Plans

- Design an enhanced and optimized product.
- Planning to work on secondary address rule, where the tool or tools will determine the delivery point for secondary address lines (e.g. apartment or suite numbers) as well as PO Boxes and Rural Route addresses.
- The tool or tools will determine if specific types of delivery points, such as high-rises, are included in the dataset.

# REFERENCES

1. [USPS Publication 28 \(Postal Addressing Standards\)](#)
2. [USPS CASS Technical Guide](#)
3. <https://regexr.com/>



***THANK YOU***