

# Harmonsierung

Datenharmonisierung ist der Prozess, bei dem unterschiedliche Datensätze aus verschiedenen Quellen standardisiert und aufeinander abgestimmt werden, um eine einheitliche und konsistente Darstellung zu gewährleisten. Durch die Harmonisierung von Daten kannst Du die [Datenqualität](#) verbessern und zuverlässige Analysen ermöglichen. Daher ist es besonders wichtig in Bereichen wie Business Intelligence, wo präzise und vergleichbare Daten entscheidend sind.

# Using Regular Expressions for Data Cleaning

Regular expressions (regex) are powerful tools that allow you to search, extract, and manipulate text data efficiently.

Whether you're dealing with messy survey responses, log files, or unstructured text data, mastering regex will help you automate and streamline the data-cleaning process.

## Why Use Regular Expressions for Data Cleaning?

Regular expressions are invaluable when dealing with text data for several reasons:

- **Flexible Pattern Matching:** Regex allows you to define search patterns to match specific sequences of characters.
- **Efficient Text Extraction:** You can quickly extract key information such as dates, phone numbers, email addresses, or specific keywords.
- **Automated Data Cleaning:** Regex can help remove unwanted characters, whitespace, punctuation, or standardize text formats.

For example:

- Extracting emails from a dataset of customer reviews.
- Cleaning and formatting phone numbers in different formats.
- Removing special characters, HTML tags, or punctuation from text data.

## Common Regular Expressions for Data Cleaning

Here are some common regex patterns you'll frequently use:

1. `\d`: Matches any digit (0-9).
2. `\w`: Matches any word character (alphanumeric plus underscore).
3. `\s`: Matches any whitespace character (spaces, tabs).
4. `[a-zA-Z]`: Matches any alphabetic character.
5. `^`: Anchors the pattern to the start of a string.
6. `$`: Anchors the pattern to the end of a string.
7. `+`, `*`, `{}`: Quantifiers that control how many times a character or group of characters can appear.

Tutorial: Applying Regex for Data Extraction and Cleaning

In this tutorial, we'll use Python's built-in re library to demonstrate how to apply regex for data extraction and cleaning tasks. We'll work with text data that contains customer feedback and format it to ensure consistency.

### Step 1: Importing the Necessary Libraries

Python's re module provides a wide range of functions to work with regular expressions. Import the necessary libraries:

```
import re
import pandas as pd
```

### Step 2: Loading a Sample Dataset

We'll use a sample dataset of customer feedback, which contains a mix of structured and unstructured text that needs cleaning.

```
# Sample dataset: Customer feedback
data = pd.DataFrame({
    'FeedbackID': [1, 2, 3],
    'FeedbackText': [
        'Great product! Contact me at 555-123-4567.',
        'Terrible service! My email is example@test.com.',
        'The product was okay... Could be better! Call me @ (555) 987-6543.'
    ]
})
# Displaying the dataset
print(data)
```

This dataset contains FeedbackText with customer comments that include phone numbers, email addresses, and inconsistent punctuation.

### Step 3: Extracting Phone Numbers with Regex

Let's start by extracting phone numbers from the FeedbackText column. We'll use a regex pattern to identify phone numbers in various formats, such as 555-123-4567 or (555) 987-6543.

```
# Function to extract phone numbers using regex
def extract_phone_numbers(text):
    phone_pattern = r'\(?(\d{3})?\)?[-.\s]?(\d{3})[-.\s]?(\d{4})'
    phone_numbers = re.findall(phone_pattern, text)
    return phone_numbers

# Applying the function to extract phone numbers
data['PhoneNumber'] = data['FeedbackText'].apply(extract_phone_numbers)
print("Extracted Phone Numbers:")
print(data[['FeedbackText', 'PhoneNumber']])
```

This regex pattern matches phone numbers in different formats, allowing us to extract them regardless of their formatting in the text. We use the `findall()` function to extract all matches and return them as a list.

#### Step 4: Extracting Email Addresses with Regex

Next, we'll extract email addresses from the `FeedbackText` column using a regex pattern designed to match common email formats.

```
# Function to extract email addresses using regex
def extract_emails(text):
    email_pattern = r'[a-zA-Z0-9_+-.]+@[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+'
    emails = re.findall(email_pattern, text)
    return emails
```

```
# Applying the function to extract email addresses
data['Email'] = data['FeedbackText'].apply(extract_emails)
print("Extracted Emails:")
print(data[['FeedbackText', 'Email']])
```

This regex pattern captures email addresses by matching alphanumeric characters, special characters like underscores and periods, followed by an `@` symbol and a valid domain name. Using `findall()`, we extract all email addresses from the text.

#### Step 5: Cleaning and Formatting Text Data

Now let's clean the `FeedbackText` column by removing unwanted characters such as punctuation, special symbols, or multiple spaces.

```
# Function to clean text by removing special characters and punctuation
def clean_text(text):
    cleaned_text = re.sub(r'[^\w\s]', '', text) # Removes all non-word characters
    cleaned_text = re.sub(r'\s+', ' ', cleaned_text) # Replaces multiple spaces with a single space
    return cleaned_text.strip()
```

```
# Applying the function to clean text data
data['CleanedFeedback'] = data['FeedbackText'].apply(clean_text)
print("Cleaned Feedback Text:")
print(data[['FeedbackText', 'CleanedFeedback']])
```

The regex pattern `[^\w\s]` matches any character that is not a word character or whitespace, effectively removing punctuation and special symbols. Additionally, we use `\s+` to replace multiple spaces with a single space.

#### Step 6: Standardizing Date Formats

If your dataset contains dates in multiple formats, you can use regex to extract and standardize them. For instance, you might encounter dates like `12/31/2023`, `2023-12-31`, or `December 31, 2023`. You can define regex patterns to capture and standardize these formats.

```
# Example regex pattern to standardize dates (not in this dataset, but for demonstration)
def standardize_dates(text):
    date_pattern = r'\b(\d{1,2})[-/](\d{1,2})[-/](\d{2,4})\b'
    standardized_text = re.sub(date_pattern, r'\3-\1-\2', text) # Reformat dates as YYYY-MM-DD
    return standardized_text
```

```
# Applying the function to a sample text (for demonstration purposes)
sample_text = "The event is on 12/31/2023."
standardized_text = standardize_dates(sample_text)
print(f"Standardized Date Text: {standardized_text}")
```

This pattern matches dates in formats like MM/DD/YYYY or MM-DD-YYYY and reorders them to the YYYY-MM-DD format using `re.sub()`.

## Conclusion

In today's post, we explored how to use regular expressions (regex) to clean and extract data from unstructured text. Regex is a powerful tool for automating data cleaning tasks, making it easier to prepare text data for analysis or modeling.

## Key Takeaways:

- Regex allows you to search, extract, and manipulate text data efficiently by defining flexible patterns.
- You can use regex to extract specific elements like phone numbers, email addresses, or dates from text.
- Regex is also useful for cleaning text data by removing special characters, punctuation, or unwanted whitespace.

# Regex 101 For Data Scientists

Understanding the basics of regular expressions is essential for any data scientist..

Regular expressions, also known as regex, are a powerful and flexible tool for working with text data. They are useful for a wide range of tasks, including:

1. **Searching:** Regular expressions allow you to search for specific patterns in text data. You can use them to find all occurrences of a pattern in a string, or to search for a pattern in a larger text file or dataset.
2. **Matching:** Regular expressions allow you to match specific patterns in text data. You can use them to check if a string conforms to a certain pattern, or to extract specific pieces of information from a larger string.
3. **Manipulating:** Regular expressions allow you to manipulate text data in various ways. You can use them to replace certain patterns with other values, to split strings into smaller substrings, or to remove unwanted characters or patterns.
4. **Validating:** Regular expressions can be used to validate the format of text data, such as email addresses, phone numbers, and other types of information.

Overall, regular expressions are a valuable tool for working with text data in data science, and they can save you time and effort when dealing with large datasets or unstructured data. In this blog post, we'll cover the fundamentals of regex as well as some little-known functions.

## Basic syntax

The basic syntax of regular expressions consists of a pattern to match, surrounded by forward slashes (/) or other delimiters. For example, the regex pattern `/cat/` would match the string "cat" in the text "the cat in the hat".

You can use various special characters and metacharacters to define more complex patterns. For example, the caret (^) symbol matches the start of a string, the dollar sign (\$) matches the end of a string, and the dot (.) matches any single character. You can use square brackets ([]) to match any one of a set of characters, and the asterisk (\*) to match zero or more of the preceding character or pattern.

For example, the regex pattern `/^The/` would match the string "The" at the beginning of a sentence, and the pattern `/hat$/` would match the string "hat" at the end of a sentence. The pattern `/cat./` would match the strings "cat," "cat!" or "cat?", and the pattern `/[Tt]he/` would match the strings "The" and "the". The pattern `/ca*/` would match the strings "c", "ca", or "caa", and so on.

## Common regex functions and Python's re library

There are many functions and libraries available for working with regular expressions in different programming languages. In Python, for example, you can use the re library to search and manipulate strings using regex patterns.

The re.search() function allows you to search for a pattern in a string and return a match object if the pattern is found. The re.findall() function returns a list of all matches in a string.

The re.sub() function allows you to replace all occurrences of a pattern with a different string.

For example, the following code uses the re.search() function to find the first occurrence of the regex pattern /cat/ in a string:

Copy code

```
import re

string = "the cat in the hat"

match = re.search("cat", string)

if match:

    print("Match found:", match.group())

else:

    print("Match not found")
```

The output of this code would be "Match found: cat".

re.findall()

The re.findall() function returns a list of all matches in a string. It is often used to extract specific pieces of information from a larger string.

For example, the following code uses the re.findall() function to find all occurrences of the regex pattern /\d+/ (which matches one or more digits) in a string:

```
import re

string = "There are 3 cats and 2 dogs."

matches = re.findall("\d+", string)

print(matches)
```

The output of this code would be ['3', '2'].

re.sub()

The re.sub() function allows you to replace all occurrences of a pattern with a different string. It takes three arguments: the pattern to search for, the replacement string, and the string to perform the search in.

For example, the following code uses the re.sub() function to replace all occurrences of the regex pattern /cat/ with the string "dog" in a string:

```
import re
string = "the cat in the hat"
new_string = re.sub("cat", "dog", string)
print(new_string)
```

The output of this code would be "the dog in the hat".

#### Lesser-known Regex Functions for Data Scientists

1. `re.split()` - This function allows you to split a string into a list of substrings based on a regex pattern. It is similar to the `str.split()` function, but it allows you to use a regex pattern as the delimiter.
2. `re.finditer()` - This function returns an iterator over all matches in a string. It is similar to `re.findall()`, but it returns match objects instead of just the matched strings.
3. `re.subn()` - This function is similar to `re.sub()`, but it returns a tuple containing the modified string and the number of substitutions made.
4. `re.escape()` - This function returns a string with all non-alphanumeric characters escaped, making it safe to use in a regex pattern.
5. `re.purge()` - This function clears the regex cache, which can be useful if you are working with a large number of regex patterns and want to free up memory.
6. `re.fullmatch()` - This function tries to match a regex pattern to the entire string. It returns a match object if the pattern matches the entire string, or `None` if it does not.

It's always a good idea to familiarize yourself with the functions and options available in your programming language's regex library, as they can save you time and make your code more efficient.

In conclusion, regular expressions (regex) are a powerful and essential tool for data scientists. They allow you to search, match, and manipulate text data in a flexible and efficient way. Whether you are cleaning and preprocessing data, extracting specific pieces of information, or performing other text-based tasks, regular expressions can help you get the job done.

It's important to understand the basic syntax of regular expressions and how to use various special characters and metacharacters to define patterns. You should also be familiar with the functions and libraries available in your programming language of choice for working with regular expressions. With a little practice and some trial and error, you can master the art of regex and use it to your advantage in your data science projects.



# Working with Data

What you need to know for Data Management and Data Wrangling

Regular Expressions are fancy wildcards. Typically abbreviated "regex", they allow you to find / match as well as replace inexact patterns and even special characters (like tabs and line breaks) in text. This is useful in many programming languages, but also for finding-and-replacing in documents. Most programming languages and many text editors, including Microsoft Word (2010+) allow the use of standard regex symbols and syntax. There are some differences in how regular expressions work in different software, but much of the basics are identical. The term "grep" which is used by programmers synonymous with search, stands for "global regular expression print".

Regular Expressions and Wildcards can be used for finding and replacing or removing text in a file, database, or filename. It is very powerful and can even assist in cleaning and reformatting data.

[xkcd cartoon about Regular Expressions "saving the day"](#)

## Background

Wildcards (the concept behind RegEx)

- [Searching Effectively Using Truncation](#) - Animated Java tutorial explains why wildcards are useful
- [Wildcard Definition](#) (Computer Hope) - Examples of \*, ?, and [] from DOS Windows, and Linux ([more examples](#))
- [How to Use Wildcards When Searching in Word 2013](#) - Examples of using Word's Special Wildcards

Showcase for what RegEx can do

- Writing: [Why writers want to learn regular expressions](#) by Jakob Persson - Several useful examples in context
- Business: [Everyday Text Patterns](#) by Dan Nguyen -- A narrative showcases several uses for a business
- History: [Understanding Regular Expressions](#) by the Programming Historian - Complex follow-along example with LibreOffice Calc.
  - See also [Cleaning OCR'd text](#), a complete example utilizing Python.

## Getting Started

Follow-Along Exercises

- [RegexOne](#) - Step-by-step instruction with real-time practice
- [Learn regular expressions in about 55 minutes](#) by Sam Hughes - Guided practice
- Video: [Using Regular Expressions](#) from Lynda.com

## Build and Test

- [RegExr](#) - Lots of help, quick reference and a pleasing interface; see capture groups by replacing \$& with \$1, \$2, etc
- [Regular Expressions 101](#) - Good for programmers; can mimic php, javascript, and python; has a library of scripts

## Reference

- [What is Regex](#) and [Quick Reference](#) from Computer Hope
- [Regular Expressions Cheat Sheet](#) by Dave Child
- [Regular Expression Cheat Sheet](#) by Nicola Pietrolungo

## Advanced

### Test Yourself

- [The regular expression game](#) (Machine Learning Lab) - State what you know and get 12 timed exercises to see if you can capture the correct strings.
- [Regex Crossword](#)
- [Regex Golf](#)

### Tips

- [How to Build a Good RegEx](#) (Sitepoint) - Examples illustrate the process of designing a string match
- [Regex tips](#) (David Birnbaum) - Shows and explains some commonly needed expressions, especially for XSLT

## Books

All books are available ONLINE through Mason.

- [Mastering Regular Expressions](#) by Jeffrey E. F. Fried
- [Regular Expressions Cookbook](#) by Steven Levithan; Jan Goyvaerts
- [Beginning Regular Expressions](#) by Andrew Watt

## Comprehensive Websites

- [Regular-Expressions.info](#) - From Jan Goyvaerts, creator of the software Regex Buddy -- See also the [Quick Start](#)
- [Rex Egg](#) - from Rex, a regex consultant -- Including a [Cheat Sheet](#)

## Implementation in Specific Programs

The vast majority of RegEx symbols are exactly the same across software, so these resources are broadly useful.

Stata	<a href="#">What are regular expressions and how can I use them in Stata?</a> (STATA) More examples: <a href="#">How can I extract a portion of a string variable using regular expressions?</a> (UCLA)
R	Finding: <a href="#">Regular Expression in R</a> (Gloria Li and Jenny Bryan) Extracting <a href="#">Extract information from texts with regular expressions in R</a> (Kun Ren)
Python	<a href="#">Regular Expression HOW TO</a> and <a href="#">Operations</a> (Python Documentation) <a href="#">Text Processing in Python</a> by David Mertz -- related: <a href="#">Matching Patterns in Text: The Basics</a>
Notepad++	<a href="#">How to use regular expressions in Notepad++</a> (NP++ Wiki) <a href="#">A guide to using regular expressions and extended search mode</a> (Mark Antoniou)
Atom	<a href="#">Find and Replace</a> (Atom Manual) <a href="#">Regular expressions in Atom for code consistency</a> (Modelit)
Office	Microsoft: <a href="#">Finding and replacing characters using wildcards</a> (Graham Mayor, MVP) OpenOffice: <a href="#">Regular Expressions in Writer</a> (OpenOffice Documentation)
Perl	<a href="#">Perl Regular Expression</a> (Perl Documentation) <a href="#">Teach Yourself Perl in 21 Days</a> by Laura Lemay - excerpt: <a href="#">Pattern Matching</a>
SQL	<a href="#">SQL Wildcards</a> (W3 Schools)

## Additional Links:

<https://www.integrate.io/blog/using-regular-expressions-in-big-data/>

<https://medium.com/analytics-vidhya/how-regular-expression-is-used-in-getting-insightful-data-e558ab21744a>