

# Regular Expressions

Marion Fischerlehner

Anja Scheuchenstuhl

# Überblick + Motivation

---

- Regular Expressions (auch regex, RE oder regex pattern)
- Zeichenfolgen suchen, extrahieren, ersetzen, validieren
- spezialisierte Programmiersprache in Python integriert
- `import re`

# Funktionen

Funktion	Beschreibung
<code>re.compile(pattern, flags=0)</code>	erstellt aus Regex -Pattern ein Regex -Objekt, das wiederholt aufgerufen werden kann.
<code>re.match(pattern, string, flags=0)</code>	Überprüft, ob der Regex zu Beginn des Strings passt
<code>re.search(pattern, string, flags=0)</code>	Sucht den gesamten String nach einem Muster ab
<code>re.findall(pattern, string, flags=0)</code>	Gibt alle Vorkommen eines Musters im String als Liste zurück
<code>re.sub(pattern, repl, string, count=0, flags=0)</code>	Ersetzt Teile eines Strings, die einem Muster entsprechen
<code>re.split(pattern, string, maxsplit=0, flags=0)</code>	Teilt einen String an allen Stellen, die einem Muster entsprechen

# **Grundlegende REGEX Muster**

# Einfache Zeichen (literal characters)

---

- Meisten Zeichen entsprechen sich selbst
- Achtung: case sensitiv!

```
import re
beispiel = "Einführung in das Programmieren"
muster = "Einführung"

re.findall(muster, beispiel)
```

# Metacharacters

Metacharacter	Bedeutung	Beispiel
<code>.</code> (Punkt)	Passt zu jedem Zeichen (außer Zeilenumbruch)	<code>a.c</code> passt zu <code>abc</code> , <code>aXc</code> , aber nicht zu <code>ac</code>
<code>^</code> (Zirkumflex)	Passt nur am Anfang eines Strings/Zeile	<code>^Hello</code> passt zu <code>Hello World</code> , aber nicht zu <code>Well Hello</code>
<code>\$</code> (Dollar)	Passt nur am Ende eines Strings/Zeile	<code>end\$</code> passt zu <code>the end</code> , aber nicht zu <code>ending soon</code>
<code>*</code> (Sternchen)	0 oder mehr Wiederholungen des vorherigen Zeichens	<code>ba*</code> passt zu <code>b</code> , <code>ba</code> , <code>baa</code> , <code>baaa</code>

# Metacharacters

Metacharacter	Bedeutung	Beispiel
<b>+</b> (Plus)	1 oder mehr Wiederholungen des vorherigen Zeichens	<b>ba+</b> passt zu <b>ba</b> , <b>baa</b> , <b>baaa</b> , aber nicht zu <b>b</b>
<b>?</b> (Fragezeichen)	0 oder 1 Vorkommen des vorherigen Zeichens	<b>ba?</b> passt zu <b>b</b> und <b>ba</b> , aber nicht zu <b>baa</b>
<b>{n}</b>	Genau n Wiederholungen des vorherigen Zeichens	<b>ba{2}</b> passt zu <b>baa</b> , aber nicht zu <b>ba</b> oder <b>baaa</b>
<b>{n,m}</b>	Zwischen n und m Wiederholungen des vorherigen Zeichens	<b>ba{1,3}</b> passt zu <b>ba</b> , <b>baa</b> , <b>baaa</b>

# Metacharacters

Metacharacter	Bedeutung	Beispiel
<code>[]</code> (eckige Klammern)	Definiert eine Zeichenklasse (Menge von Zeichen)	<code>[aeiou]</code> passt zu beliebigen Vokal, <code>[A-Za-z]</code> zu beliebigen Buchstaben
<code> </code> (Pipe)	ODER-Verknüpfung für Auswahl zwischen Optionen	<code>cat dog</code> passt zu <code>cat</code> oder <code>dog</code> .
<code>()</code> (runde Klammern)	Gruppiert Ausdrücke + ermöglicht wiederholte Anwendung	<code>(ab)+</code> passt zu <code>ab</code> , <code>abab</code> , <code>ababab</code>
<code>\</code> (Backslash)	Maskiert Sonderzeichen oder definiert spezielle Sequenzen	<code>\.</code> passt zu einem Punkt, <code>\d</code> zu einer Zahl (0-9)



# Metacharacters

---

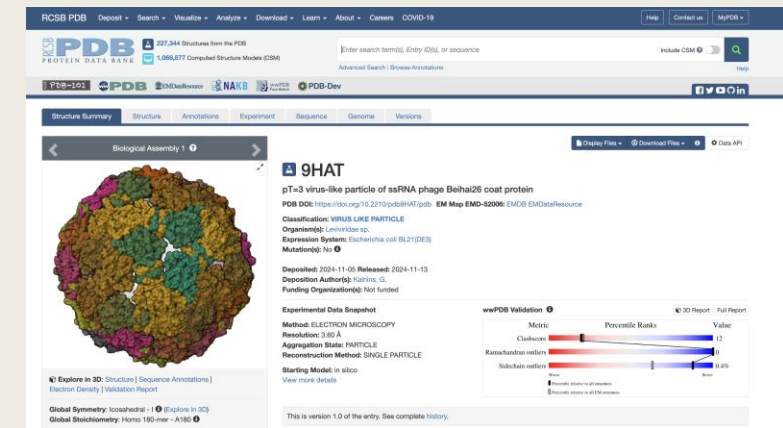
```
sequence = "AUGCGUUACGGAUGUGCCGCGUAAGCUA"  
pattern = re.compile(r"^AUG")  
match = pattern.findall(sequence)
```

```
sequence = "AUGCGAAAUACGGAUGUGCCGCGUAAGCUA"  
pattern = re.compile(r"(UAA)$")  
match = pattern.findall(sequence)
```

```
sequence = "AUGCGAAAUACGGAUGUGCCGCGUAAGCUA"  
pattern = re.compile(r"(UAA|UAG|UGA)$")  
match = pattern.findall(sequence)
```

# **Praktische Regex- Anwendungen**

# Regex in der Praxis



Erstelle ein Pattern, damit nur Proteine aus der die nach 4 Stellen mit „.pdb“ enden ausgegeben werden. Tipp: 1. Character ist immer numerisch

Lösung: `re.compile(r"[1-9][A-Z0-9]{3} \.pdb\b")`

2MLB.pdb  
2MMP.pdb  
2MP8.pdb  
2MPK.pdb  
2MR9.pdb  
2MRL.pdb  
2MSG.pdb  
2MSJ.pdb  
2MUL.pdb  
2MW9.pdb  
2MWA.pdb  
2MWB.pdb  
2MWF.pdb  
2MWR.pdb  
2MXD.pdb  
2MYX.pdb  
2N2U.pdb  
2N35.pdb  
2N3S.pdb

EHEE\_rd4\_0172.pdb  
EHEE\_rd4\_0195.pdb  
EHEE\_rd4\_0300.pdb  
EHEE\_rd4\_0325.pdb  
EHEE\_rd4\_0340.pdb  
EHEE\_rd4\_0394.pdb  
EHEE\_rd4\_0463.pdb  
EHEE\_rd4\_0499.pdb  
EHEE\_rd4\_0502.pdb  
EHEE\_rd4\_0510.pdb  
EHEE\_rd4\_0625.pdb  
EHEE\_rd4\_0726.pdb  
EHEE\_rd4\_0840.pdb  
EHEE\_rd4\_0864.pdb  
EHEE\_rd4\_0877.pdb  
EHEE\_rd4\_0929.pdb  
GG:run1\_0874\_0002.pdb  
GG:run1\_1113\_0003.pdb  
GG:run1\_1390\_0006.pdb

r9\_1163\_TrROS\_Hall.pdb  
r9\_340\_TrROS\_Hall.pdb  
r9\_366\_TrROS\_Hall.pdb  
r9\_733\_TrROS\_Hall.pdb  
r9\_814\_TrROS\_Hall.pdb  
v2K28S:R29S\_2M7K.pdb  
v2K43S\_2KVV.pdb  
v2R14S:R16S\_2L3X.pdb  
v2R31S:R32S\_2N5D.pdb  
v2\_10GW.pdb  
v2\_1PGY.pdb  
v2\_1PV0.pdb  
v2\_1WR7.pdb  
v2\_2BN8.pdb  
v2\_2HDZ.pdb  
v2\_2J6K.pdb  
v2\_2KIS.pdb  
v2\_2L15.pdb  
v2\_2LC2.pdb

# Zusammenfassung

---

- wichtiges Tool beim Arbeiten mit komplexen Texten und Daten
- Zeichenfolgen suchen, extrahieren, ersetzen oder validieren
- wichtig, sich mit den grundlegenden Metacharacters und der Regex-Syntax auseinanderzusetzen