


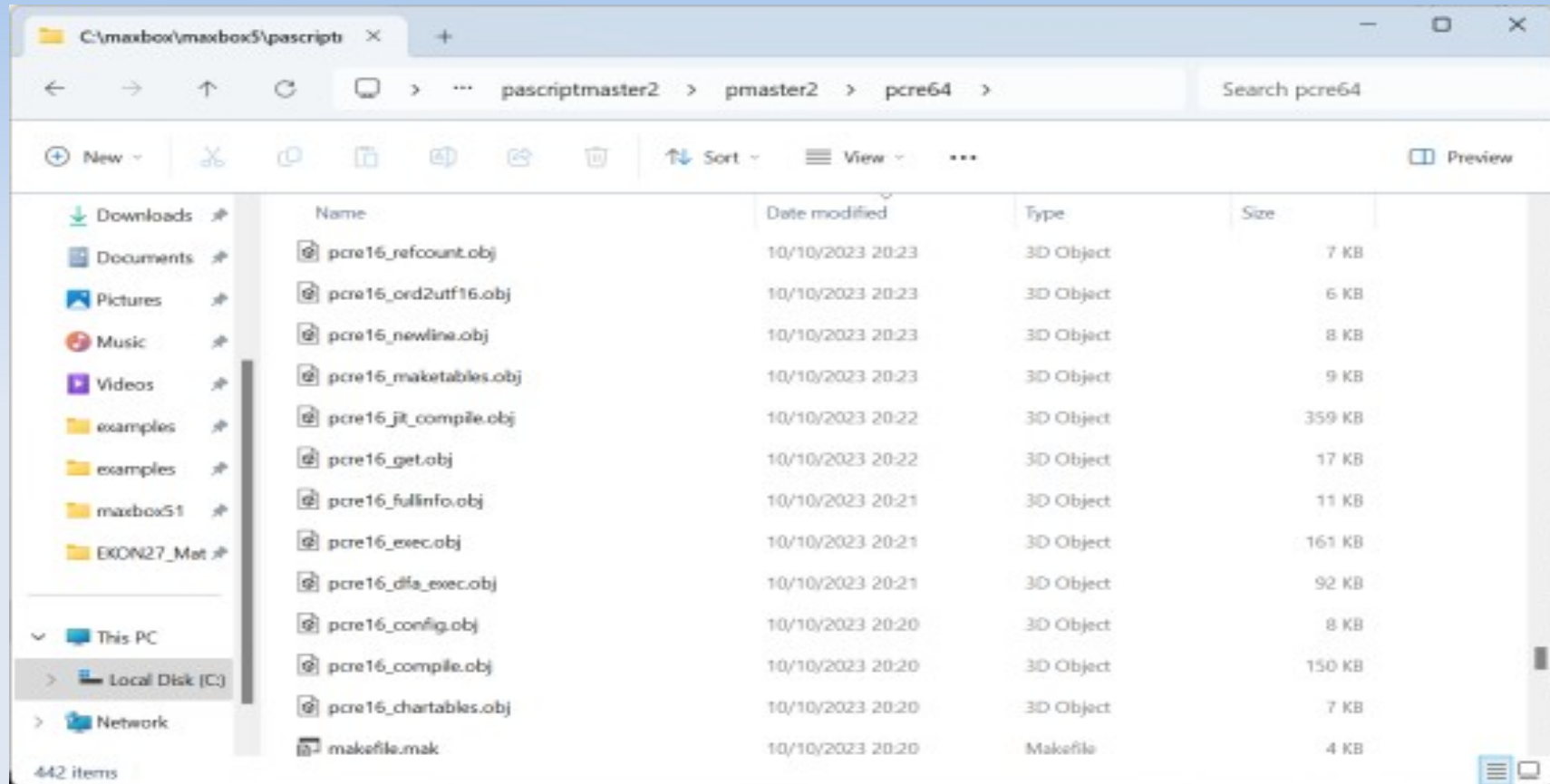


Modern Regex

Nov. 2024 V Max Kleiner

- TPerlRegEx with out of the box demos
- PCRE library
- TRegEx with TMatch/TMatchEvaluator 
- <https://www.pcre.org/original/doc/html/pcre16.html>
- This session shows you various ways of using modern regex in your application.
- Python re lib
- <https://regex101.com/>

Compile first



- *The supplied pcrelib.dll contains PCRE 7.9, compiled with Unicode support which works with FreePascal, Lazarus, Delphi, Jupyter and maXbox.*



RegEx Research

- TPerlRegEx is a Delphi VCL wrapper around the open-source PCRE (Perl-Compatible Regular Expressions) library. It provides powerful regexp capabilities similar to those found in the Perl programming language.
- This version of TPerlRegEx is compatible with the TPerlRegEx class in the RegularExpressionsCore unit in Delphi XE.
- <https://maxbox4.wordpress.com/2024/05/10/modern-regexp/>
- <https://entwickler-konferenz.de/blog/machine-learning-mit-cai/>

Cross-platform RegEx

- Use the TRegEx class from the System.RegularExpressions unit.
- This class provides methods and properties for working with regular expressions, such as Match () and Replace () for matching and replacing strings, and Captures () and Groups () for accessing matched groups.

Let's practice: maxbox51\examples\1313_regex_db12.pas

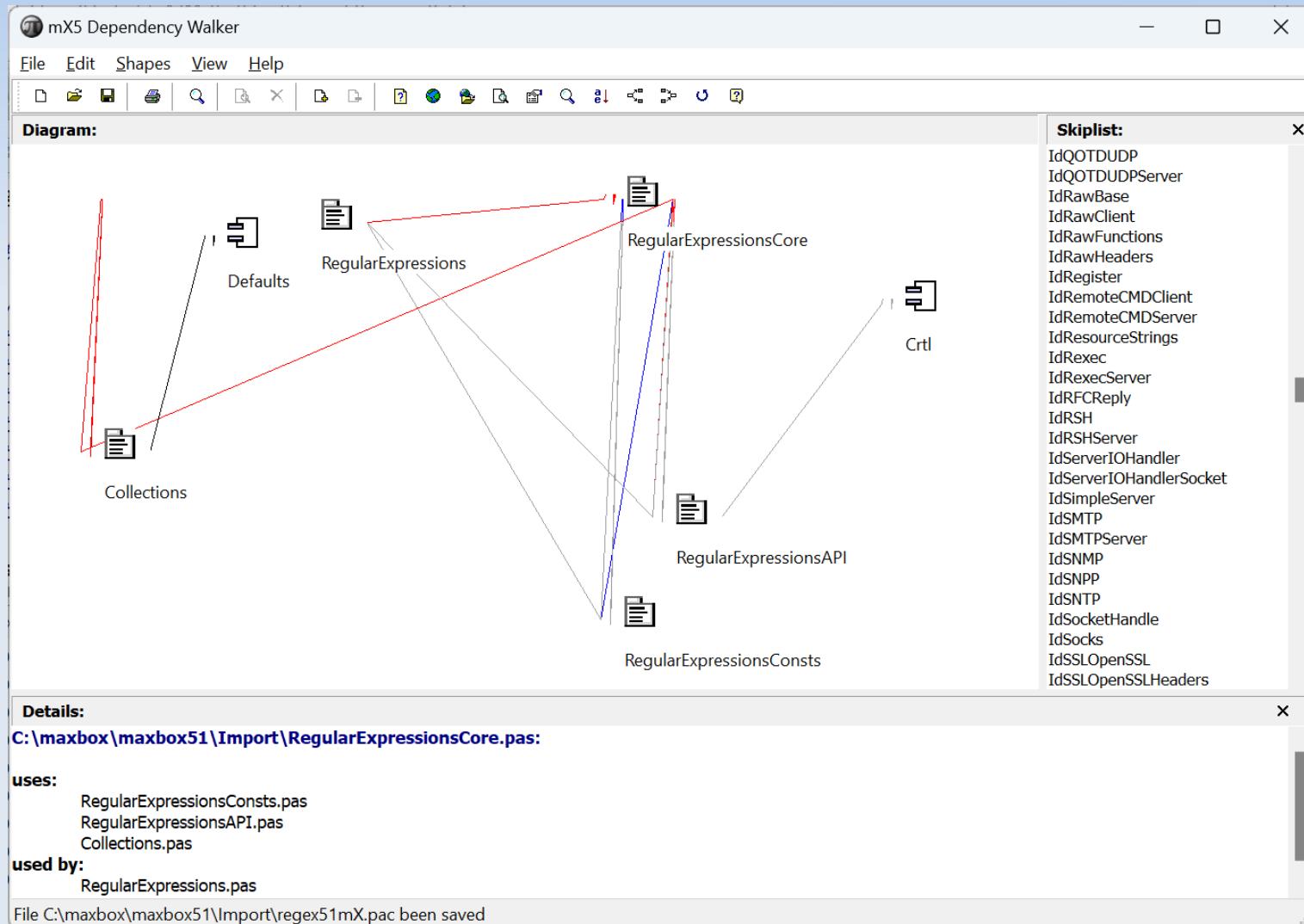
TestRegExMultiMatcher .\1317_regex_matchevaluator1.pas

<https://github.com/maxkleiner/maXbox/blob/master/logisticregression2.ipynb>



Packages View

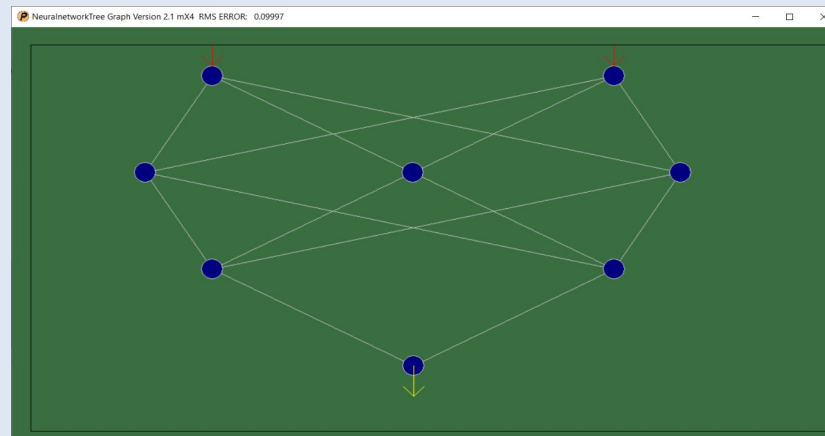
- RegularExpression as Client uses Core (uses System.SysUtils, System.RegularExpressionsCore;)



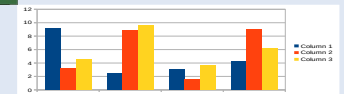


Web & RegEx = WebEx

- `writeln(RegexMatch(IdWhois1.Wholes('domain ibm.com'), '.*Registry Expiry Date.*', false));`
- > Registry Expiry Date: 2025-03-20T04:00:00Z
- `writeln(RegexMatch(IdWhois1.Wholes('domain wordpress.com'), '.*Registry Expiry Date.*', false));`
- > Registry Expiry Date: 2033-03-03T12:13:23Z



.\examples\1302_restcountries_API_24_mcJSON1regexEKON28.txt





Be aware of static record

Performance Object versus Static

The source for `TRegEx.IsMatch(const Input, Pattern: string; Options: TRegExOptions)` shows that a `TRegEx` is created at every invocation (which is a costly operation):

```
1  procedure TForm1.Button1ClickRegExPerformance1(Sender: TObject);
2  var
3      i: integer;
4      t1, t2: cardinal; arex: TRegEx;
5  begin
6      t1 := GetTickCount;
7      arex:= TRegEx.Create1('')
8      for i:= 1 to 100000 do
9          arex.IsMatch3('abcdef', 'cd', [rrIgnoreCase]);
10     t2:= GetTickCount;
11     writeln('rextime1: '+IntToStr(t2-t1)+' ms');
12 end;
```

Conclusion: Use explicit object instance, especially in loops

Demo: <https://maxbox4.wordpress.com/2024/05/10/modern-regex/>



TMatch

A more modern implementation is to code with a TMatch and TMatchCollection class.

This example demonstrates the use of TMatchCollection and TGroupCollection. This example assumes that you have placed a TButton, a TEdit and a TMemo on a form.

```
10 // Creates and lists the match collection, the matches in that
11 // collection and the groups in those matches.
12 procedure TForm1.Button1Click(Sender: TObject);
13 const
14     bigString = 'Look for a the strings in this strang of strungs.';
15     littlestring = '(str)([iau]ng)';
16 var
17     regex: TRegex;
18     i, j: integer;
19     mygrps: TGroupCollection;
20 begin
21     regex:= TRegex.Create(littlestring);
22     mycoll:= regex.Matches(bigString);
23     Edit1.Text:= 'Count: ' + IntToStr(mycoll.Count);
24     memo1.Lines.Add('First Collection: ');
25     for i:= 0 to mycoll.Count-1 do begin
26         memo1.Lines.Add('Match #' + IntToStr(i) + ': ' + mycoll.Item[i].Value);
27         memo1.Lines.Add('Group: ' + IntToStr(i));
28         mygrps:= mycoll.Item[i].Groups;
29         for j:= 0 to mygrps.Count-1 do
30             memo1.Lines.Add('Value: ' + mygrps.Item[j].Value);
31     end;
32 end;
```




TMatch II

The item of a TMatchCollection returns the Match identified by index from the collection (ex. `tmatches[it-1].value` below).

[https://docwiki.embarcadero.com/CodeExamples/Alexandria/en/TMatchCollectionCount_\(Delphi\)](https://docwiki.embarcadero.com/CodeExamples/Alexandria/en/TMatchCollectionCount_(Delphi))

In general matches from a TRegex returns all the matches present in the input string and is useful to iterate through a group or captured group (next slide):

Matches returns all the matches present in the Input string in the form of a TMatchCollection instance. If the Pattern parameter is not present the regular expression used is specified in the TRegex constructor.

StartPos specifies the starting position to start the search. TMatchCollection has no public constructor. It is created as the return value of the Matches method. The collection is populated with one TMatch instance for each match found in the input string. The Count property is a length of the TMatchCollection set. Length specifies the substring, starting at StartPos to match with the regular expressions.

Code as script:

https://sourceforge.net/projects/maxbox/files/Examples/13_General/646_pi_evil2_64_12.TXT/download



From PI Package

■ Numeric Analysis of PI Explore

```
1  function getMatchString2(arex, atext: string): string;
2  var Match: TMatch; tMatches: TMatchCollection;
3      myenum: TMatchCollectionEnumerator;
4  begin
5      with TRegEx.Create1(arex) do
6      try
7          it:= 0;
8          { Match format search...}
9          result:= result+CRLF;
10         if ismatch(atext) then
11             tMatches:=Matches(aText);
12         writeln('captured groups: '+itoa(tmatches.count ));
13         repeat
14             Inc(it);
15             result:= result+Format(#09'%d: %-12s',[it, tmatches[it-1].value])
16             if it mod 5=0 then
17                 result:= result+#13#10;
18             //until match(atext).success; //MatchNext < 0;
19             until it = tmatches.count;
20         finally
21             Free;
22         end;
23         WriteLn('Done REX2 - Hit NOthing to exit');
24     end;
```



Big Iterator as Collection

- TMatchCollection has no public constructor. It is created as the return value of the Matches method. A collection is populated with one TMatch instance for each match found in input string.

```
regEx := TRegEx.create('common':"[\w]', [rrNotEmpty]);  
// Execute search of TMatch  
for it := 0 to envlist.count-1 do  
  if regEx.match((envlist[it])).success then begin  
    writeln(itoa(cnt)+' ':'+envlist[it]);  
    inc(cnt)  
  end;
```

```
langitem: 22 Schweiziska edsförbundet, swe  
langitem: 23 İsviçre Konfederasyonu, tur  
langitem: 24 سوئیس متحده, urd  
langitem: 25 瑞士联邦, zho
```

Using Groups

- `System.RegularExpressions.TMatch.Groups`
- *Contains a collection of groups from the most recent match with a regular expression.*
- A regular expression pattern can include subpatterns, which are defined by enclosing a portion of regex pattern in parentheses. Every such subpattern captures a subexpression or group. For ex., the regex pattern `(\d{3})-(\d{2})-(\d{4})`, which matches social security numbers.
The first group consists of the first three digits and is captured by the first portion of the regular expression, `(\d{3})`.

```
writeln(regx1.match3('2-2321 55-99878 456-545','(\d{2})-(\d{4})').value);
```

1311_RestClientLibrary_httprequestC_EKON28.txt



Let's compile

- Contains a collection of groups from recent match with a reg expression.

The screenshot shows the maXbox5 64-bit ScriptStudio IDE. The main editor displays a Pascal script for a procedure `TForm1Button1Click`. The script defines a regular expression `littlestring = '(str)([iau]ng)'` and uses `regex.Matches(bigString)` to find matches in the string `'Look for a the strings in this strang of strungs.'`. The results pane at the bottom shows the execution output:

```
Count: 3
First Collection:
Match #0: string
Group: 0
Value: string
Value: str
Value: ing
Match #1: strang
Group: 1
Value: strang
```

The right-hand pane shows the 'Interface List' for the script, listing various functions and procedures defined in the script, including `ReadUntil`, `StripTags2`, `getMatchString`, `TimesTable2`, `ParseAttributes`, `getMatchString2`, `getMatchStringSortGroup2Re`, `rex_tester_dual`, `getMatchStringSortGroup2`, `getMatchStringSortGroup3`, `RexusageExample`, and `Locs`.



Collection Matches

- A collection of groups as the result of a match with a single regular expression. A regex pattern can include subpatterns, which are defined by enclosing a portion of the regex.

```
1  var mygrps: TGroupCollection;
2
3  regex:=TRegEx.create('"common": "[\w]*(su)*', [rroNotEmpty, rrosingleline]);
4  mycoll:= regex.Matches(envlist.text);
5  writeln('Count: ' + IntToStr(mycoll.Count));
6  // Execute search3 of TMatchCollection
7  for it:= 0 to mycoll.count-1 do begin
8      writeln(itoa(it)+':'+(mycoll.Item[it].Value));
9      mygrps:= mycoll.Item[it].Groups;
10     for j := 0 to mygrps.Count-1 do
11         writeln('Value: ' + mygrps.Item[j].Value);
12     end;
13     regex.Free;
14     envlist.Free;    (//*)
15
```



The Main Python

```
1  procedure PyCodeREgEx(tex: string);
2  begin
3      with TPythonEngine.Create(nil) do begin
4          try
5              loadDLL;
6              //autofinalize:= false;
7              ExecString('import re, json');
8              execstr('txt = "The rain in Explain"');
9              //if it starts with "The" and ends with "plain":
10             Println('regex test: '+EvalStr('re.search("^The.*plain$", txt)'));
11             execstr('x = re.search("^The.*plain$", txt)');
12             Println('regex test: '+EvalStr('type(x)'));
13             Println('regex test: '+EvalStr('bool(x)'));
14             // Println(EvalStr('decode_and_print_json('+JSONDATA+')'));
15         except
16             raiseError;
17         finally
18             unloadDLL;
19             free;
20         end;
21     end;
22 end;
```



re module

Python has a built-in package called `re`, which can be used to work with Regular Expressions as `re.match()` or `re.search()`.

RegEx Functions

The `re` module offers a set of functions that allows us to search a string for a match:

Function	Description
<u>findall</u>	Returns a list containing all matches
<u>search</u>	Returns a <u>Match object</u> if there is a match anywhere in the string
<u>split</u>	Returns a list where the string has been split at each match
<u>sub</u>	Replaces one or many matches with a string

- <https://docs.python.org/3/library/re.html>



Demo FindFiles()

```
procedure TForm1FormCreate(Sender: TObject);
var k,t: integer;
    items: TStringList;
begin
    items:= TStringList.create;
    for k:= 0 to 9 do
        StringGrid1.Cells[0, k+1]:= cs10Labels[k];

    //FindAllFiles(ComboBox1.Items, 'csdata');
    FindFiles(exepath+'data', '*.bmp', items);
    writeln(items.text);
    for t:= 1 to items.count-1 do
        ComboBox1.Items.add(items[t]);
    if ComboBox1.Items.Count > 0 then begin
        ComboBox1.text:= ComboBox1.Items[0];
        if FileExists(ComboBox1.text) then begin
            Image1.Picture.LoadFromFile(ComboBox1.text);
            Image2.Picture.LoadFromFile(ComboBox1.text);
            label1.Caption:= extractfilename(ComboBox1.text);
        end;
    end;
end;
```



Compare Delphi Python

```
1  procedure PyCodeREgEx(tex: string);
2  begin
3      with TPythonEngine.Create(nil) do begin
4          try
5              loadDLL;
6              //autofinalize:= false;
7              ExecString('import re, json');
8              execstr('targetstr = "Emma''s luck numbers are 251, 761, and 231"');
9              execstr('pattern = re.compile("\d{3}")');
10             execstr('matches = pattern.findall(targetstr)');
11             Println('py regex test: '+EvalStr('matches'));
12         except
13             raiseError;
14         finally
15             unloadDLL;
16             free;
17         end;
18     end;
19 end;
```

```
1  writeln('show matchall/findall in python:');
2  with TregEx.create1('\d{3}') do begin
3      matchg:= Match('Emma''s luck numbers are 251, 761, and 231 ');
4      repeat
5          write([' '+ matchg.value+ ' ']);
6          matchg:= matchg.NextMatch;
7      until not matchg.success
8      free;
9  end;
```



Unicode Group Samples



// {\$APPTYPE CONSOLE} ☺ † π 📌

- `writeln(regx.ReplaceAll('\u0418\u0443, \u0427\u0436\u044d\u0446\u0437\u044f\u043d',
'\u([0-9a-f]{0,4})', '$$+', [rrIgnoreCase, rroSingleLine]));`

UC_teststr:= 'Düsseldorf, Köln, 北京市, دوسلدورف, إسرائيل, Αλφα !@#\$',

myEval: TMatchEvaluator;

myeval:= @EvaluatorU;

`Writeln(regx.Replace7('\u0418\u0443, \u0427\u0436\u044d\u0446\u0437\u044f\u043d',
'\u([0-9a-f]{4})', myeval, [rrIgnoreCase]));`

`mycoll:= regx1.matches2('match non-english words like können or
móc zu Çin', '(?s)(.[^\x00-\x7F]\b)+')`

<https://www.regexpal.com/>

<https://github.com/maxkleiner/maXbox/blob/master/objectdetector3.ipynb>



Conclusion

- Internally Delphi uses class `TPerlRegEx` and has such methods for groups and collections.
- Number of matched groups stored in the `Groups` array. E.g. when the regex `"(a)|(b)"` matches "a", `GroupCount` will be 1. When the same regex matches "b", `GroupCount` will be 2.
- The static `TMatch` record or class as instance provides several properties with details about the match. `Success` indicates if a match was found.
- You can use a numeric index to `Item[]` for numbered capturing groups, or a string index for named capturing groups thanks to variants!
> `whatGotMatched:= Match.Groups['MatchName'].Value;`

Method: Design Regex with a Online Site like regex101.com

Model: Object **Regex Pattern** + Subject Data

Metric: Test with generic data and community pattern

<https://raw.githubusercontent.com/breitsch2/maXbox4/master/assets/graph3.html>



Modern Regex

Thanks for coming!

Materials:

<https://maxbox4.wordpress.com/2024/05/10/modern-regex/>

<https://maxbox4.wordpress.com/2024/06/20/ekon-28/>

<https://medium.com/@maxkleiner1/modern-regex-d9d3450fbd36>

<https://maxbox5.wordpress.com/2024/07/22/ekon-28/>