



Lab 6 Text

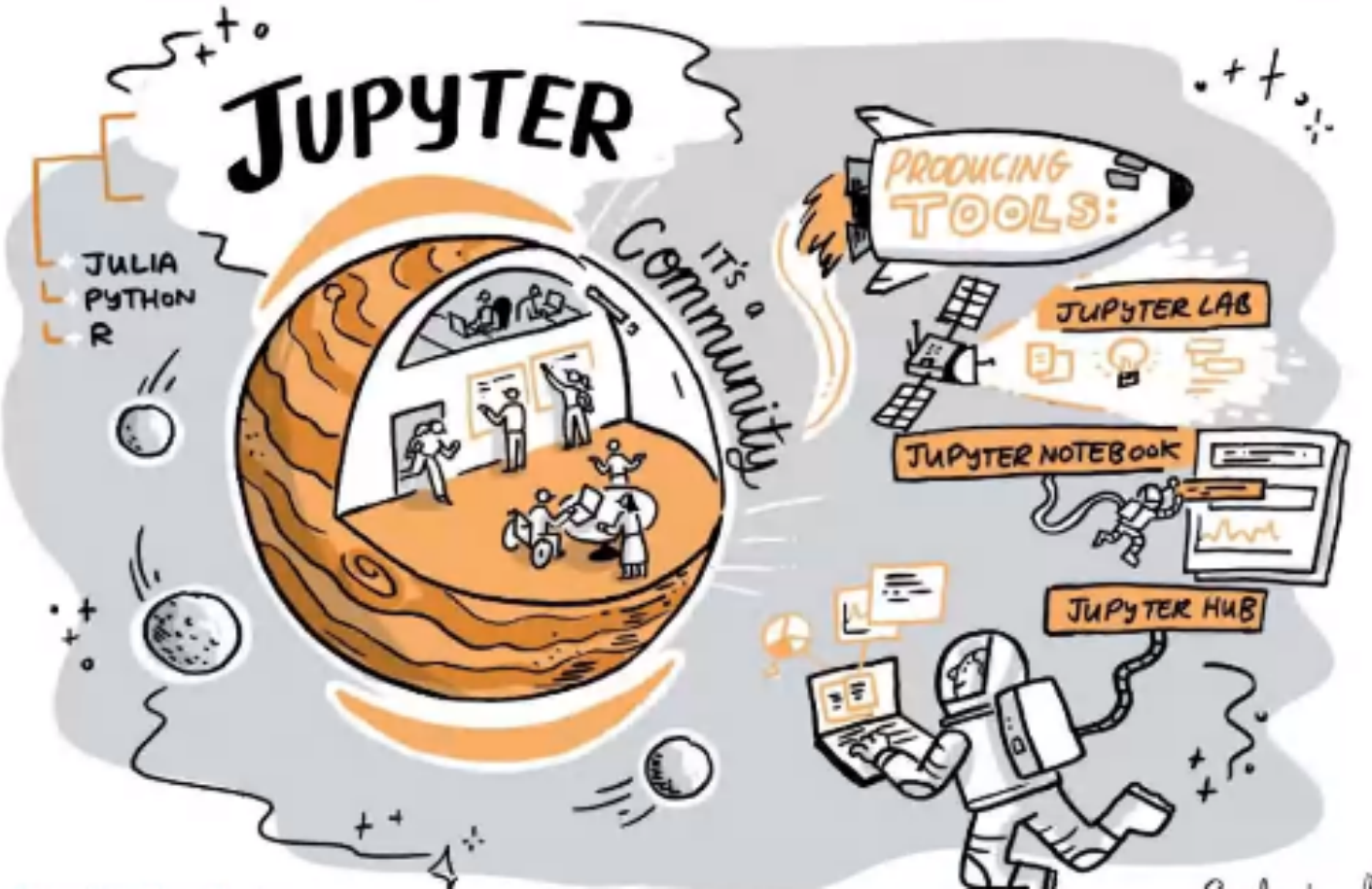
Programming for Business Analytics
MG-GY 8401





Agenda

- Text
- Regular Expressions

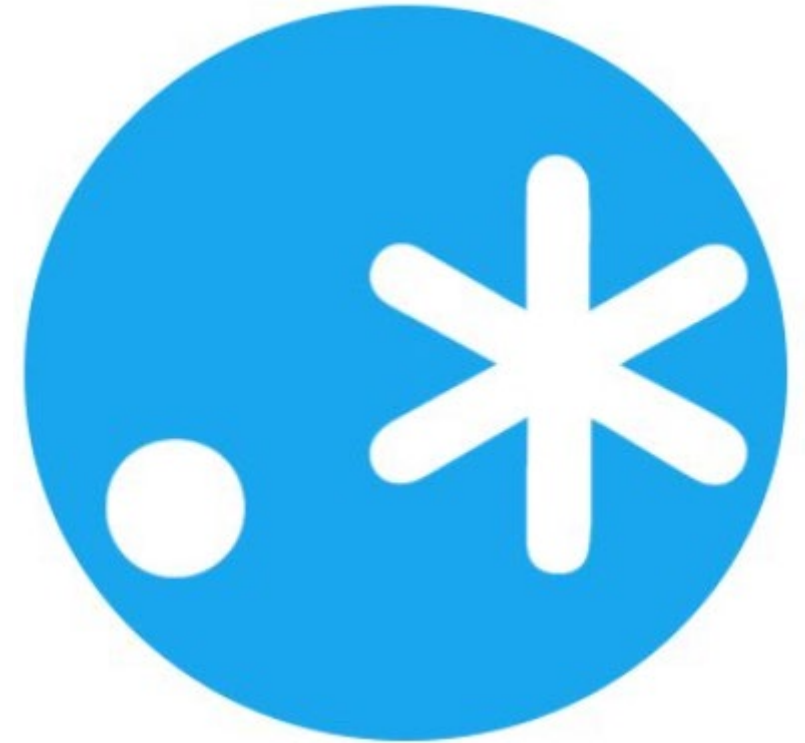


- Describe the learning objectives.
- Summarize the relevant take-aways.
- Ask about unclear information.



The ability to extract patterns of characters in text leads to many applications

- ▶ Access information in digital libraries
- ▶ Filter text such as spam emails
- ▶ Validate data-entry fields such as date or URL





operation	order	example	matches	does not match
concatenation	3	AABAAB	AABAAB	every other string
or	4	AA BAAB	AA BAAB	every other string
closure (zero or more)	2	AB*A	AA ABBBBBBA	AB ABABA
parenthesis	1	A(A B)AAB	AAAAB ABAAB	every other string
		(AB)*A	A ABABABABA	AA ABBA



operation	example	matches	does not match
any character (except newline)	<code>.U.U.U.</code>	CUMULUS JUGULUM	SUCCUBUS TUMULTUOUS
character class	<code>[A-Za-z][a-z]*</code>	word Capitalized	camelCase 4illegal
at least one	<code>jo+hn</code>	john joooooooohn	jhn jjohn
zero or one	<code>joh?n</code>	jon john	any other string
repeated exactly {a} times	<code>j[aeiou]{3}hn</code>	jaoehn jooohn	jhn jaeiouhn
repeated from a to b times: {a,b}	<code>j[ou]{1,2}hn</code>	john juohn	jhn jooohn



regex	matches	does not match
<code>. *SPB. *</code>	RASPBERRY CRISPBREAD	SUBSPACE SUBSPECIES
<code>[0-9]{3}-[0-9]{2}-[0-9]{4}</code>	231-41-5121 573-57-1821	231415121 57-3571821
<code>[a-z]+@([a-z]+\.)+(edu com)</code>	horse@pizza.com horse@pizza.food.com	frank_99@yahoo.com hug@cs



operation	example	matches	does not match
built-in character classes	<code>\w+</code> <code>\d+</code>	fawef 231231	this person 423 people
character class negation	<code>[^a-z]+</code>	PEPPERS3982 17211!↑å	porch CLAmS
escape character	<code>cow\.com</code>	cow.com	cow\$com



operation	example	matches	does not match
beginning of line	<code>^ark</code>	ark two ark o ark	dark
end of line	<code>ark\$</code>	dark ark o ark	ark two
lazy version of zero or more <code>*?</code>	<code>5.*?5</code>	5005 55	5005005

5*5 would
match this!

Write a regular expression to match the pattern

- XYYYYY!

at the end of strings. Here

- X is a digit
- Y is a letter
- ! is the exclamation point