

BREAKING THE ICE

WITH
REGULAR EXPRESSIONS

Contents

- 01 The String Story
- 02 Crew Emails
- 03 Confirmative
- 04 Multi-line Strings
- 05 Capture Groups

Why Regular Expressions?

Imagine we run a website that asks people to enter their Orlando-area phone number to receive local restaurant offers.

407-555-1212



First 3 numbers have to be
either 407 or 321 (Orlando
area code)

Why Regular Expressions?

Imagine we run a website that asks people to enter their Orlando-area phone number to receive local restaurant offers.

407-555-1212



Next, there could be a dash,
but it's optional

Why Regular Expressions?

Imagine we run a website that asks people to enter their Orlando-area phone number to receive local restaurant offers.

407-555-1212



Followed by 3 more numbers

Why Regular Expressions?

Imagine we run a website that asks people to enter their Orlando-area phone number to receive local restaurant offers.

407-555-1212



Another optional dash

Why Regular Expressions?

Imagine we run a website that asks people to enter their Orlando-area phone number to receive local restaurant offers.

407-555-1212



And finally, 4 more numbers

What would the code to validate this look like?

Validating Input Without Regular Expressions

```
if  
  ((chars[0] == 4 && chars[1] == 0 && chars[2] == 7)  
   ||  
   (chars[0] == 3 && chars[1] == 2 && chars[2] == 1))  
   &&  
   chars[3] == "-" || chars[3] == "" && ...  
... and that's just the first few numbers!
```

The same validation using a regular expression

```
if (chars.match(/regular expression/))
```



FYI, the syntax can be different per language

Understanding a Regular Expression

```
if (chars.match(/regular expression/))
```

↑
Subject String

The string where we're looking for
a match is called the subject string

↑
aka, regex

A set of characters that represent rules
for searching or matching text in a
string in a concise, predictable way

Our First Regular Expression

Does the phone number have a 407 area code?

```
if (chars.match(/407/))
```



Literally matches the characters 407

number

407-555-1212

subject string



True!

Not Finding a Match

Does the phone number have a 407 area code?

```
if (chars.match(/407/))
```

number

subject string

321-555-1212



False!

Does the phone number have a 407 or 321 area code?

Using the OR / Alternation Operator

Does the phone number have a 407 or 321 area code?

```
if (chars.match(/407|321/))
```

Takes left-most match first

number

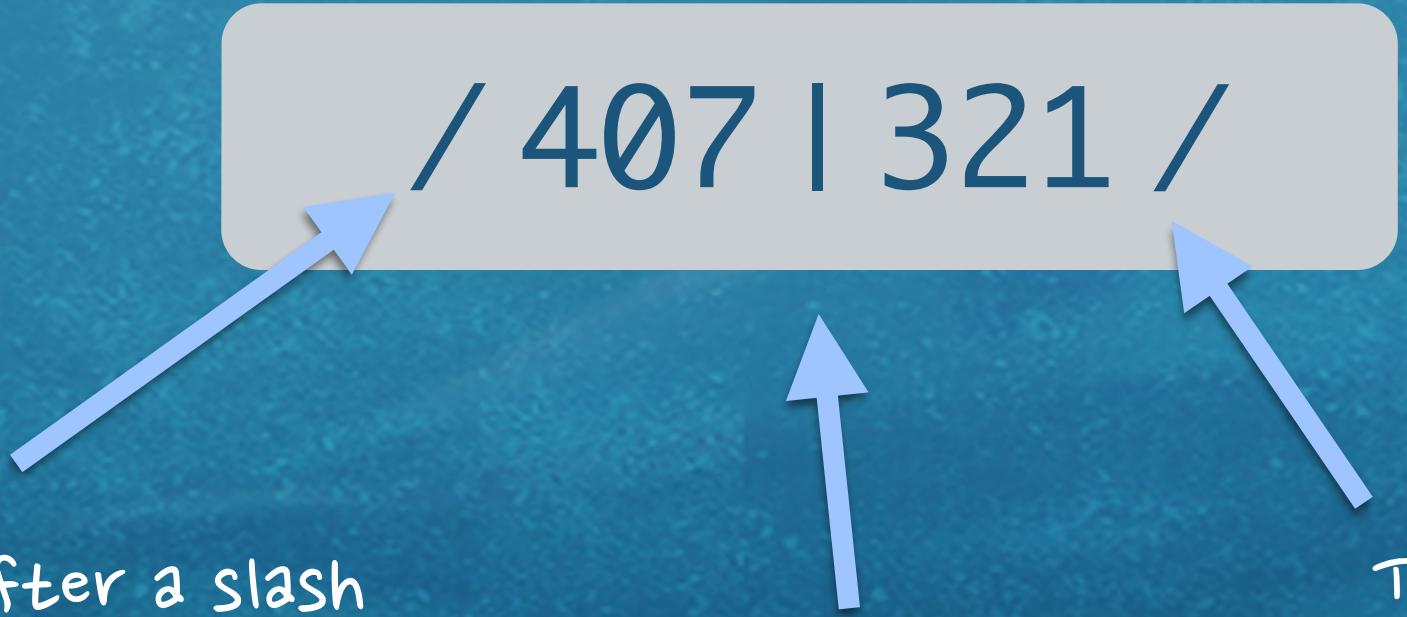
subject string



321-555-1212

True!

The Parts of a Regular Expression



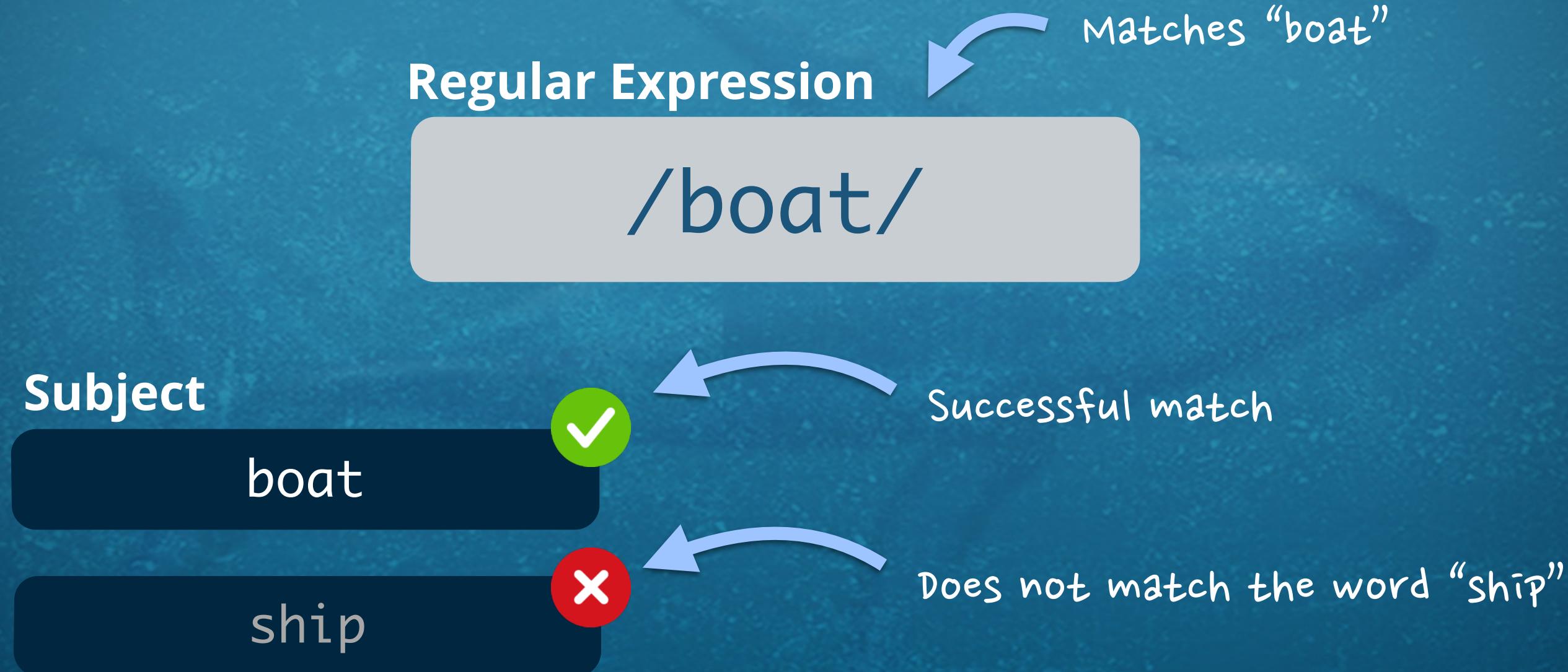
The pattern starts after a slash

The characters in the
middle are called a "pattern"

The pattern ends before this slash

Writing regex is really like asking the question,
“Does a group of characters match a specific pattern?”

Not Just Numbers



Number OR letters

Regular Expression  Matches “boat” or “ship”
`/boat|ship/`

Subject

boat

ship



Successful match



Successful match

`/407 | 321 /`

Just like our number pattern

What Are Regular Expressions Used For?

Common uses include:

Validations

Phone numbers

Email

Passwords

Domain names



Searching

Words in a sentence

Unwanted characters

Extracting sections

Replacing/cleaning/formatting



Section 2

The String Story

BREAKING THE ICE
WITH
REGULAR EXPRESSIONS

Our First Problem: Collecting Names

We are assembling a crew!



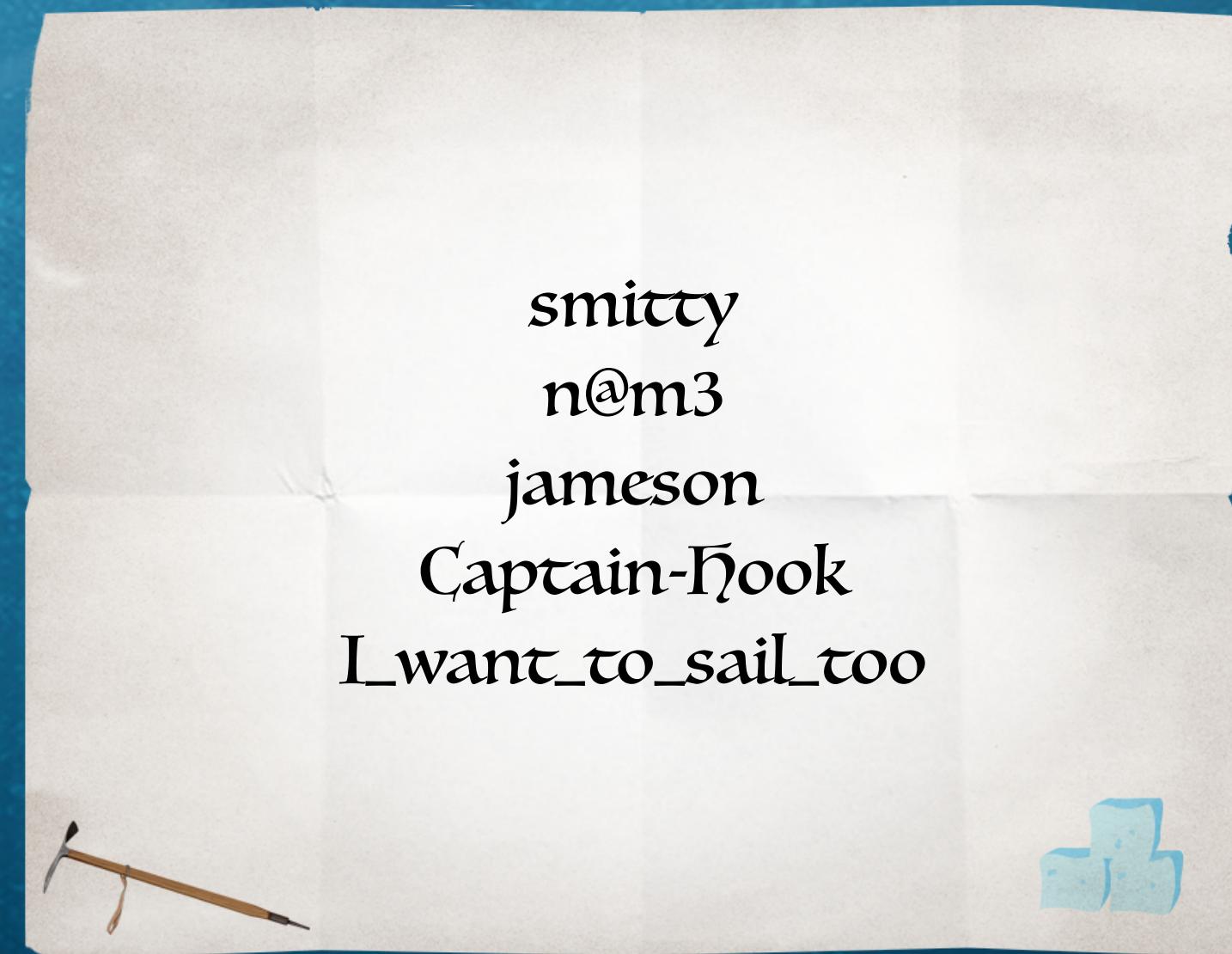
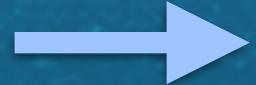
Register

Username

Password

Remember me

Register



We Need to Check Multiple Usernames



Subject

smitty



How Can We Match on Multiple Versions of Names?

Regular Expression

```
/smitty|james/
```

Match “smitty” - if no match,
then match “james”

Subject

smitty



james



OR / alternation operator allows
match on “smitty” or “james”

Problem: We Need to Repeat the “R”

Regular Expression

/ar/

Subject

ar



This works, but is inefficient

arr



This was a partial match

arrr



This was a partial match

Regular expressions look for matches anywhere in their subject

Problem: Long Expression Repeating

Regular Expression

/ar|arr|arrr|arrrr/

Subject

ar



arr



arrr



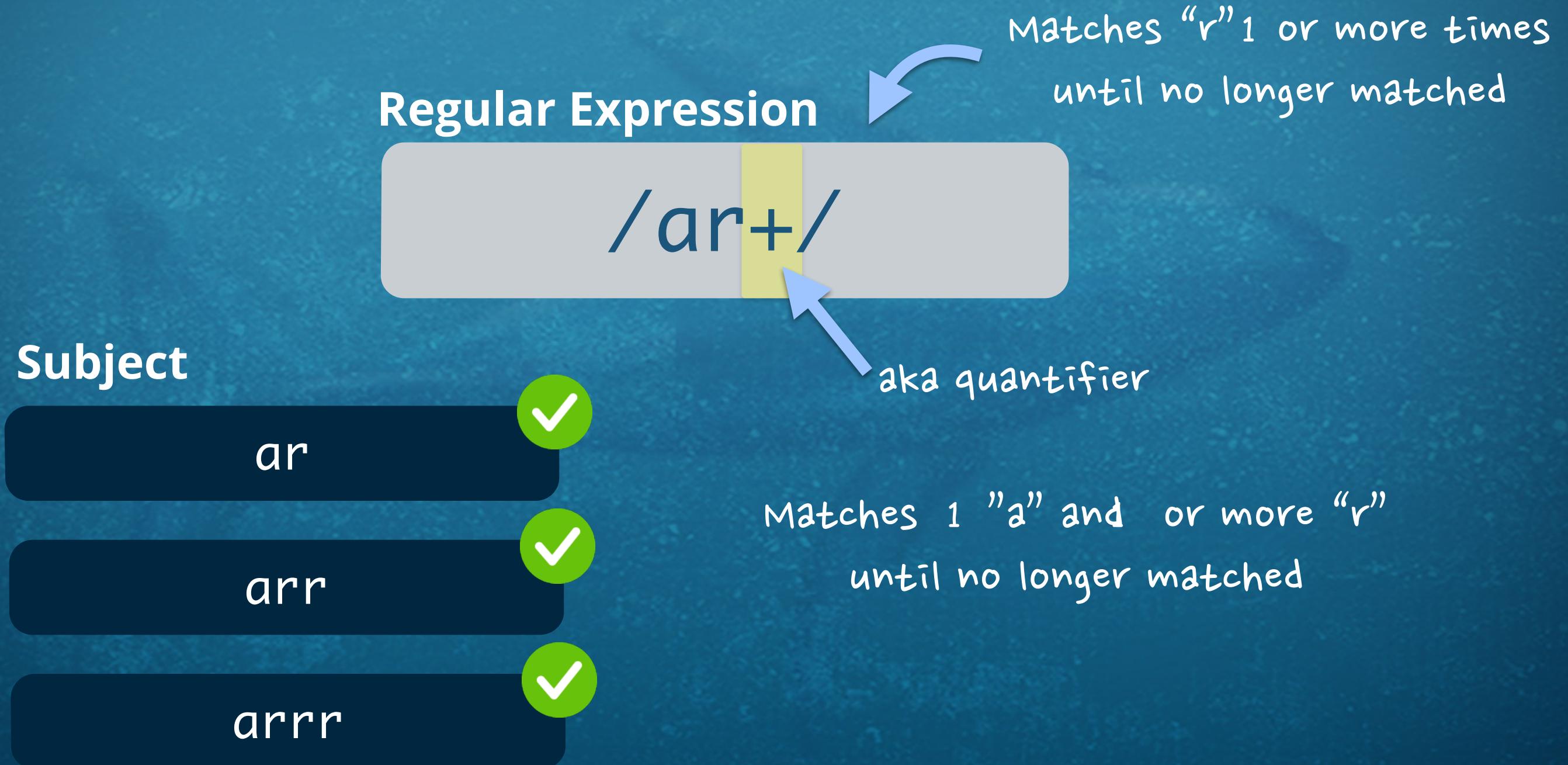
This is going to grow too large



There are just too many
versions to account for

This works, but is inefficient

Solution: The Plus Operator



This means, “Look for the previous character 1 or more times.”

How Can We Match on Multiple Versions of Names?

Regular Expression

```
/smitty|james|ar+/
```

Multiple versions of “arr”

Subject

smitty



james



ar



arrr



Matching on Multiple Versions of Names

Regular Expression

```
/smitty|james|ar+/
```

Subject

james



jameson



Last 2 letters are not matched



This was a partial match

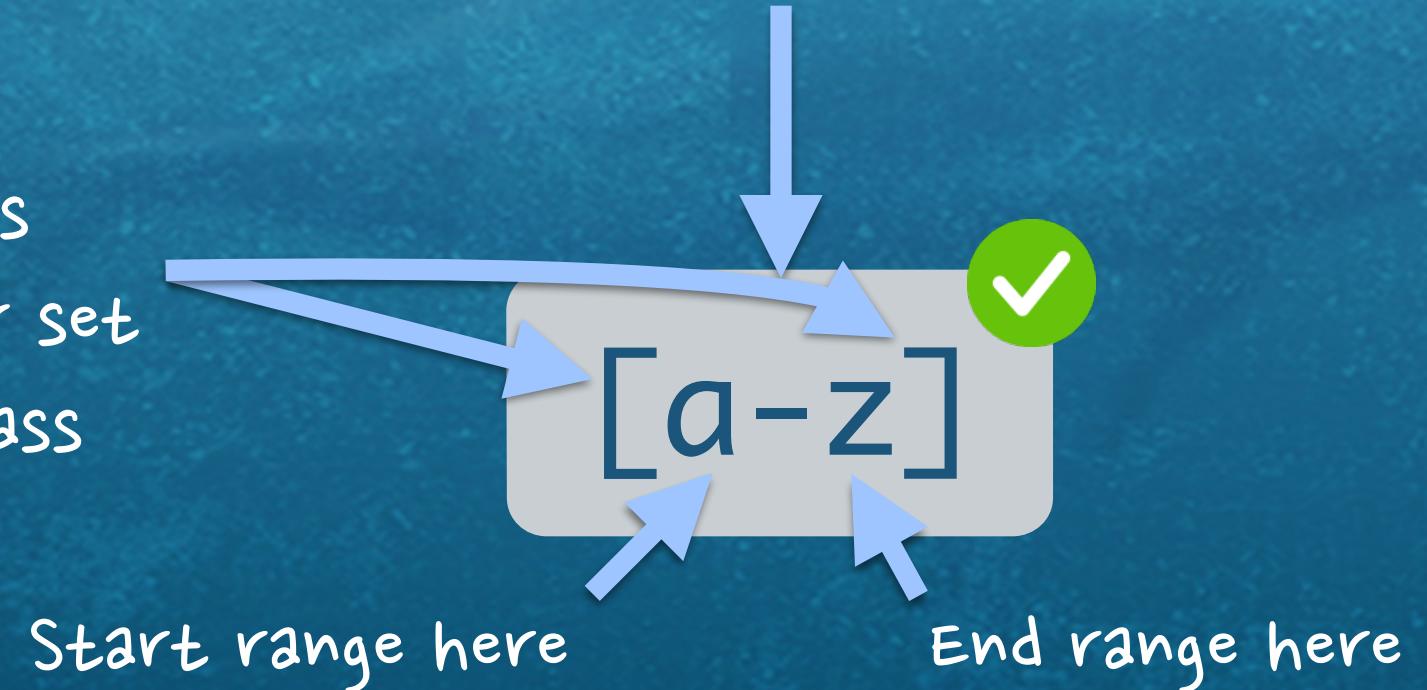
Regular expressions look for matches anywhere in their subject and take the left-most match first

Matching All Characters in the Alphabet Using Ranges

Regular Expression

/a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z/

Square brackets
create a character set
aka character class



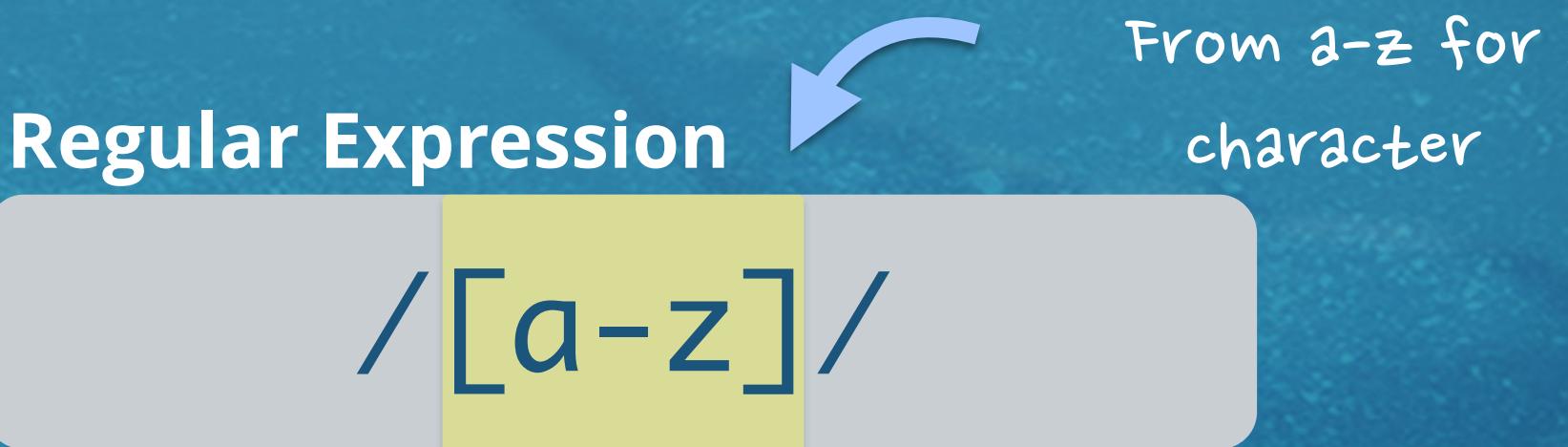
This works, but is inefficient



NOTES

- A range only works in a character set
- This character set represents 1 character

Using Our Character Set in a Pattern



Subject

smitty



Matches the first character represented by our character set.

How can we change this to match all 6 characters?

Matching Multiple Characters With Multiple Character Sets

Regular Expression

/[a-z][a-z][a-z][a-z][a-z][a-z]/

Subject

smitty



james



Does not match



From a-z for 6 characters

How can we change this pattern
to accept any amount of letters?

Checking for 1 Character Set Multiple Times

Regular Expression

A diagram of the regular expression `/[a-z]+`. The character set `[a-z]` is highlighted in yellow, and the plus operator `+` is highlighted in green. A blue arrow points from the text "Plus operator" to the green-highlighted plus sign.

`/[a-z]+"`

Plus operator



This means, “Look for the previous character 1 or more times.”

Subject

smitty



james



jameson



Now they all match!

ar



arr



arrr



Problem: Capital Letters

Regular Expression

```
/[a-z]+/
```

Subject

Blackbeard



!

Problem: Our pattern is ignoring capital letters

Matching Uppercase and Lowercase Characters

Regular Expression

```
/[a-zA-Z]+/
```

From a-z, A-Z for 1 or more times

Subject

Blackbeard



Now we're matching capital letters!

Simplifying Patterns With a Casing Modifier

Regular Expression

```
/[a-z]+/i
```

Letters after final slash
are called “modifiers”

The **i** modifier means “case insensitive,” which will
match uppercase and lowercase characters.

Subject

Blackbeard



Problem: Names With 2 Words Not Matching

Regular Expression

```
/[a-z]+/i
```

Subject

Captain hook



Problem: Not matching whitespace

How Can We Match Whitespace?

The diagram illustrates two regex patterns for matching whitespace:

- Pattern 1:** `/Captain hook/` (Incorrect)
A yellow box highlights the space between "Captain" and "hook". A blue arrow points from this space to the handwritten note: "using a literal space is hard to read - you might miss it!". An orange circle with a exclamation mark (!) is placed next to the highlighted space.
- Pattern 2:** `/Captain\shook/` (Correct)
A yellow box highlights the backslash character "\s". A blue arrow points from this character to the handwritten note: "\s is the same as saying 'a whitespace character'".

Whitespace can include:

- Spaces
- Tabs
- New lines

Matching Spaces in a Character Set

Regular Expression

```
/[a-z\s]+/i
```

\s

is the same as saying
“a whitespace character”

Subject

Captain hook



NOTE: The order of characters doesn't matter here!

[\sa-zA-Z] +

is the same thing!

Problem: How Can We Match Numbers?

Regular Expression

```
/[a-zA-Z\s]+/i
```

Does not account for numbers

Subject

Captain hook



Long John the 3rd



Problem: we do not match any numbers

Matching Numbers in a Range

Regular Expression

```
/[a-zA-Z0-9\s]+/i
```

From a-z, A-Z, 0-9, whitespace
for 1 or more times

using a range to match numbers

Subject

Captain hook

Long John the 3rd



We now match on numbers 0-9

Refactoring With the Word Metacharacter

\w

is the same as

[a-zA-Z0-9]

Regular Expression

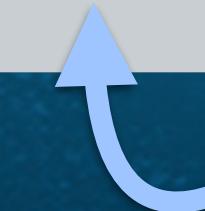
/[a-zA-Z0-9\s]+/i

!



These have the same result,
but the shortcut is easier to read.

/[\w\s]+/



Also includes underscore

01 The String Story

→ 02 Crew Emails

03 Confirmative

04 Multi-line Strings

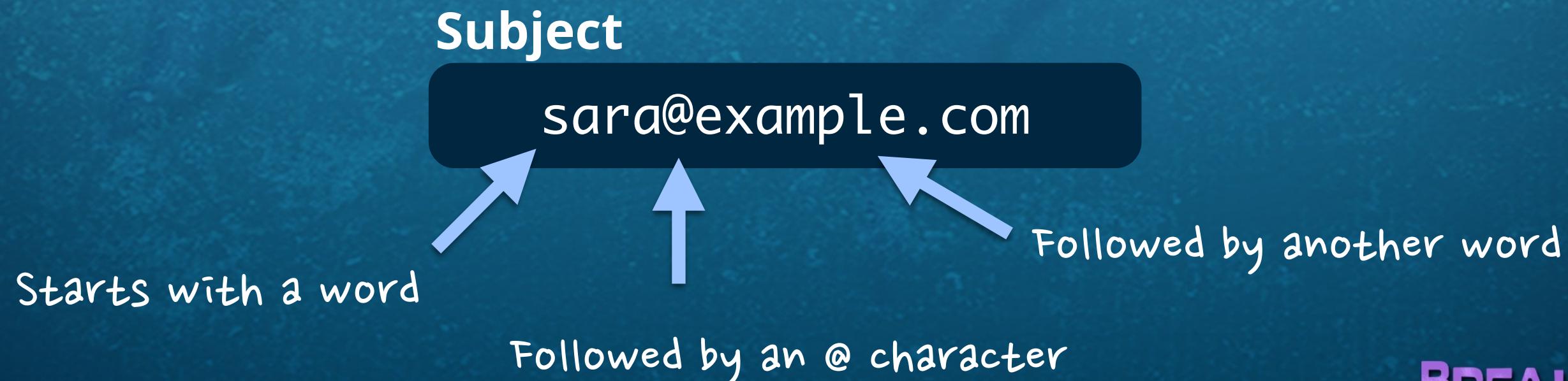
05 Capture Groups

Analyze Our New Subject String

We now have our list of names for our voyage.



Let's work on
verifying their
emails!



Starting the Regular Expression

Regular Expression

`\w@\w\w`



Remember: `\w` only searches for 1 word-like character.

Subject

`sara@example.com`

Pattern Goals

1 word character repeated 1 or more times

✓ @ symbol

1 word character repeated 1 or more times

Matching Full Words

Regular Expression

```
/\w+@\w+/
```

Subject

```
sara@example.com
```



Because the “.” is not a word, it is not matched

Pattern Goals

- ✓ 1 word character repeated 1 or more times
- ✓ @ symbol
- ✓ 1 word character repeated 1 or more times
- dot literal
- com

A Successful Match

Regular Expression

```
\w+@\w+. \w+/
```

Subject

```
sara@example.com
```



We successfully matched!

Pattern Goals

- ✓ 1 word character repeated 1 or more times
- ✓ @ symbol
- ✓ 1 word character repeated 1 or more times
- ✓ dot literal
- ✓ com

Matching the .com

Regular Expression

```
\w+@\w+\.\w+
```

Subject

sara@example.com

sara@example!com

Pattern Goals

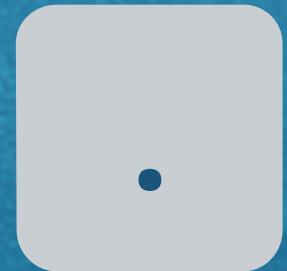
- ✓ 1 word character repeated 1 or more times
- ✓ @ symbol
- ✓ 1 word character repeated 1 or more times
- ✗ dot literal
- ✓ com

!

Hold up - an exclamation point
shouldn't be matching here

The . Wildcard

The “.” is a wildcard metacharacter



The “.” wildcard matches any character except newline



Escaping a Character in a Pattern

using a backslash, we escape special characters

Regular Expression

```
\w+@\w+\.\w+/
```

Subject

sara@example.com

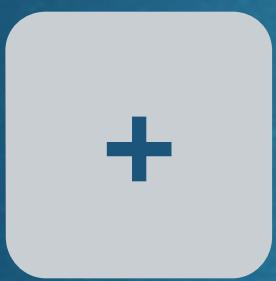
sara@example!com

Pattern Goals

- ✓ 1 word character repeated 1 or more times
- ✓ @ symbol
- ✓ 1 word character repeated 1 or more times
- ✓ dot literal
- ✓ com

Great! The incorrect email is not matched.

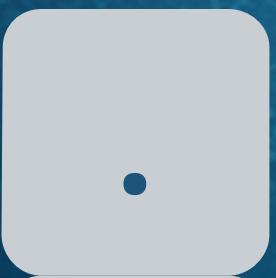
Escaping Characters



→ Matches a character 1 or more times



→ Matches literal “+” character



→ Matches any character except newline



→ Matches literal “.” character

Being More Specific About the Top-level Domain (TLD)

Regular Expression

```
\w+@\w+\.\com|net|org|edu/i
```

4 common TLDs

Subject

sara@example.net



OR / alternation
operators not
evaluating properly

Pattern Goals

- ✗ 1 word character repeated 1 or more times
- ✗ @ symbol
- ✗ 1 word character repeated 1 or more times
- ✗ dot literal
- ✓ com or net or org or edu

Grouping Together TLDs

Regular Expression

The parentheses group our TLDs together

```
^\w+@\w+\.(com|net|org|edu)/i
```

Subject

sara@example.net



Pattern Goals

- ✓ 1 word character repeated 1 or more times
- ✓ @ symbol
- ✓ 1 word character repeated 1 or more times
- ✓ dot literal
- ✓ com or net or org or edu

Problem: Need to Ignore Characters Before and After a Subject

Regular Expression

```
\w+@\w+\.(com|net|org|edu)/i
```

Subject

sara@example.comlksh

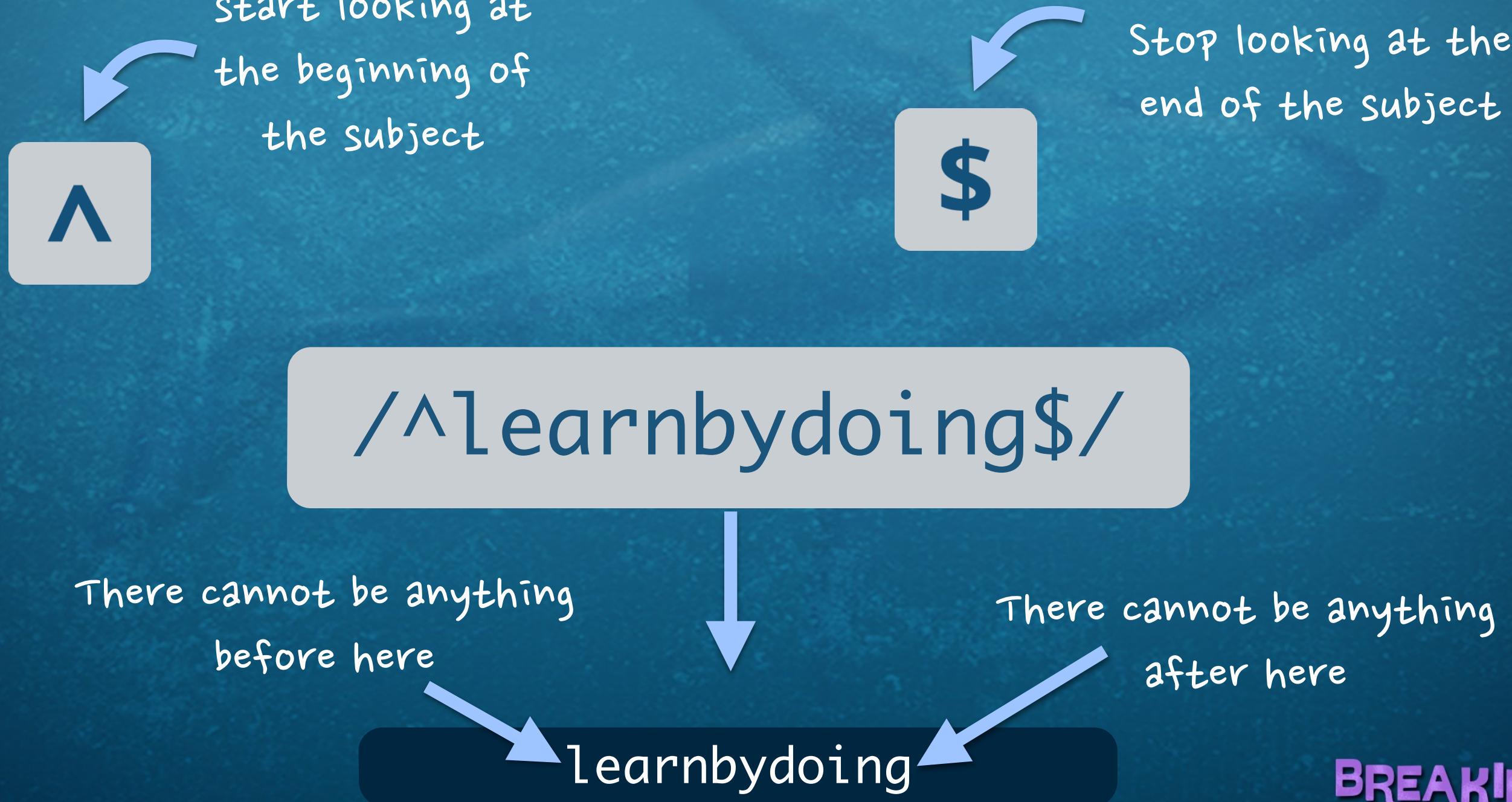


We get a partial match but shouldn't!
These characters should cause the match to fail.

Pattern Goals

- ✓ 1 word character repeated 1 or more times
- ✓ @ symbol
- ✓ 1 word character repeated 1 or more times
- ✓ dot literal
- ✓ matches com or net or org or edu
- ✓ no text before or after email

Introducing Anchors



Adding Anchors to Our Pattern

Regular Expression

```
/^\w+@\w+\.(com|net|org|edu)$/i
```

Start matching at the beginning of the line

Stop matching at the end of the line

Subject

sara@example.com



|sara@example.com|ksh

Great - our invalid email is no longer matched!

- ### Pattern Goals
- ✓ 1 word character repeated 1 or more times
 - ✓ @ symbol
 - ✓ 1 word character repeated 1 or more times
 - ✓ dot literal
 - ✓ matches com or net or org or edu
 - ✓ no text before or after email

Contents

01 The String Story

02 Crew Emails

→ 03 Confirmative

04 Multi-line Strings

05 Capture Groups

Joining the Voyage

Find each shipmate's answer as to whether they are willing to voyage!

Subject Strings

ok, i will do it

okie dokie

Ahooy, Okay!!!

why sure, i can go

arrrr, yes matey

My answer me mate, is yes

y



Accepted Answer Keywords

ok

Okay

sure

yes

y

Starting the Expression by Matching the First Confirmation

using ok literal to directly match “ok”

Regular Expression

/ok/



Possible Subjects

ok, i will do it



okie dokie



- Pattern Goals
- ✓ ok
 - ✓ Okay
 - ✓ sure
 - ✓ yes
 - y
 - ✗ answer by itself
 - ✓ ignore other text



We want the answer to be by itself and not part of another word

Detecting Word Boundaries With the Boundary Metacharacter

Boundary metacharacter



“whole words only”

Match words surrounded by boundary

Regular Expression

$\^b\w+\b/g$

Possible Subjects



word surrounded by
boundary

Using Boundaries in a Pattern

Regular Expression

```
/\bok\b/
```

Ensures our match is a single word

\b

= boundary metacharacter

Possible Subjects

ok, i will do it



okie|dokie



Great - we don't get a match because "ok" does not have a boundary on the right, it only has letters

- Pattern Goals
- ✓ ok
 - ✓ Okay
 - ✓ sure
 - ✓ yes
 - ✓ y
 - ✓ answer by itself
 - ✓ ignore other text

Problem: We Need to Match Multiple Versions of a Word

Regular Expression

```
\bok\b/
```

Possible Subjects

Ahooy, Okay!!!

Pattern Goals

- ✓ ok
- ✓ 0kay
- ✓ sure
- ✓ yes
- ✓ y
- ✓ answer by itself
- ✓ ignore other text

!

wait - we want "Okay" to match too

Inefficient Matching With Multiple OR Statements

Regular Expression

```
/\bok\b | \bokay\b/
```

“ok”

OR

“okay”



But what if there
were 50 different
options?

If we had a way to make “ay” optional,
we could match with a single pattern.

```
/\bok(ay)\b/
```



How can we make
this optional?

- Pattern Goals
- ✓ ok
 - ✓ Okay
 - ✓ sure
 - ✓ yes
 - ✓ y
 - ✓ answer by itself
 - ✓ ignore other text

Checking for a Pattern 0 or More Times

Regular Expression

/ship/



All of these characters are required

/ship?/



characters in a pattern that are followed
by a question mark are optional

Match the letters “ship” 1 time

ship



shirt



Match the letters “shi” 1 time
followed by 0 or 1 “p” characters

ship



shirt



Marking a Group of Characters as Optional

Regular Expression

```
/pirate\s(ship)?/
```



Group characters in parentheses followed
by "?" to make them all optional

pirate ship

pirate boat

Matching Different Versions of the Same Word

Regular Expression

```
\Abok(ay)?\b/i
```

Match 0 or 1 times

"i" modifier for upper-
and lowercase

Pattern Goals

- ✓ ok
- ✓ Okay
- sure
- yes
- y
- ✓ answer by itself
- ✓ ignore other text

Possible Subjects

Ahooy, 0kay!!!



ok, i will do it



Problem: The OR Operator Is Splitting the Boundary

Regular Expression

```
\Abok(ay)?|sure\b/i
```

Pattern Goals

- ✓ ok
- ✓ Okay
- ✓ sure
- yes
- y
- ✗ answer by itself
- ✓ ignore other text

Possible Subjects

why |sure|, i can go

|ensure| code is good



we do not want a partial match of the answer

Solution: Group the Pattern Within the Boundary

Group ensures boundary is applied to all answers

Regular Expression

```
\Ab(ok(ay)?|sure)\b/i
```

Pattern Goals

- ✓ ok
- ✓ Okay
- ✓ sure
- yes
- y
- ✓ answer by itself
- ✓ ignore other text

Possible Subjects

why sure, i can go



ensure code is good

Great - “ensure” is no longer matched

Continuing to Match More Valid Subjects

Regular Expression

```
\Ab(ok(ay)?|sure|yes)\b/i
```

Pattern Goals

- ✓ ok
- ✓ Okay
- ✓ sure
- ✓ yes

y

- ✓ answer by itself
- ✓ ignore other text

Possible Subjects

arrr, yes matey



y



! Single "y" is not matching

Another Optional Part of the Answer

Regular Expression

```
\Ab(ok(ay)?|surely(es)?)\b/i
```

Pattern Goals

- ✓ ok
- ✓ Okay
- ✓ sure
- ✓ yes
- ✓ y
- ✓ answer by itself
- ✓ ignore other text

Possible Subjects

arrr, yes matey



y



Our Finished Pattern

Regular Expression

```
\b(ok(ay)?|surely(es)?)\b/i
```

Pattern Goals

- ✓ ok
- ✓ 0kay
- ✓ sure
- ✓ yes
- ✓ y
- ✓ answer by itself
- ✓ ignore other text

Possible Subjects

ok, i will do it



Ahooy, 0kay!!!



why sure, i can go



arr, yes matey



y



Section 2

Confirmative

Sailor Taglines

All sailors have entered their taglines!



Pattern Goals

- does not contain numbers
- 40 characters or less
- 20 characters or more

Work like a captain, play like a pirate.

Keep calm and say Arr.

Shiver me timbers matey.

Why are pirates pirates? cuz they arr.



Start With a Base Pattern for Matching the Tagline

Regular Expression

```
/[a-z]+/i
```

Subject

Work like a captain, play
like a pirate.



Not matching whitespace

Pattern Goals

- ✓ does not contain numbers
- 40 characters or less
- 20 characters or more

Expanding the Pattern to Include Whitespace

Regular Expression

```
/[a-z\s]+/i
```

Subject

Work like a captain, play
like a pirate.



We also need to match the comma

Pattern Goals

- ✓ does not contain numbers
- 40 characters or less
- 20 characters or more

Including A Comma

Regular Expression

```
/[a-z\s,]+/i
```

Subject

Work like a captain, play
like a pirate.



Pattern Goals

- ✓ does not contain numbers
- 40 characters or less
- 20 characters or more

Writing a Shorter Pattern With the NOT Operator

Everything matched
with this pattern...

```
/[a-z\s,]+/i
```

...is also matched
by this one.

```
/[^d]+/i
```

Subject

Work like a captain, play
like a pirate.



The “^” means “not”
when placed within a
character set

This means “any number”

so, this pattern means
“anything that’s not a number”

1 Caret Symbol — 2 Different Meanings

^ used to anchor
beginning of subject

Regular Expression

/**[^^\d]+\$/**

Subject

Work like a captain, play
like a pirate.



Pattern Goals

- ✓ does not contain numbers
- 40 characters or less
- 20 characters or more

\d

= digit character (number)

^

= not the following character(s)

Even Shorter With Negated Shorthand Characters

[^\d]

is the same as

\D

match every character
except numbers

[^\s]

is the same as

\S

match every character
except whitespace

[^\w]

is the same as

\W

match every character
except words



NOTE: You don't need a character set when using negated shorthand characters

Using the Negated Shorthand Syntax in a Pattern

Regular Expression

```
/^\D+$/
```

Subject

Work like a captain, play
like a pirate.



Pattern Goals

- ✓ does not contain numbers
- 40 characters or less
- 20 characters or more

Problem: We Want to Set Min and Max Characters



we want a minimum
of 20 characters...

Subject

Work like a



Pattern Goals

- ✓ does not contain numbers
- 40 characters or less
- 20 characters or more

Regular Expression

```
/^\D+$/
```



Subject

Work like a captain, play
like a pirate. Work like a
captain, play like a
pirate. Work like a
captain, play like a
pirate.



...and a maximum of 40 characters

Matching a Specific Number of Times With Interval Expressions

/[a-z]{2}/

Subject

test



Matches any character from a-z
exactly 2 times

Matches at least this
amount

/[a-z]{1,3}/

Subject

t

te

test



Matches the character in question a minimum of
1 time and a maximum of 3 times

Matches at most this
amount

First, Set the Minimum Amount of Characters

Regular Expression

```
\D{20,}/i
```

Ensure length is greater than or equal to 20

Subject

```
Work like a captain, play  
like a pirate.
```

Subject with at least 20 characters is matched

Subject

```
Work like a captain.
```



Great! subject with length less than 20 characters is not matched.

Pattern Goals

- ✓ does not contain numbers
- 40 characters or less
- ✓ 20 characters or more

Next, Set a Maximum Number of Characters

Regular Expression

```
/^\D{20,40}$/
```

Subject between 20 and 40 characters matched

Subject

Work like a captain, play
like a pirate.



Great! Long subject no longer matched.



Work like a captain, play
like a pirate. Work like a
captain, play like a
pirate. Work like a
captain, play like a
pirate.

Pattern Goals

- ✓ does not contain numbers
- ✓ 40 characters or less
- ✓ 20 characters or more

Contents

01 The String Story

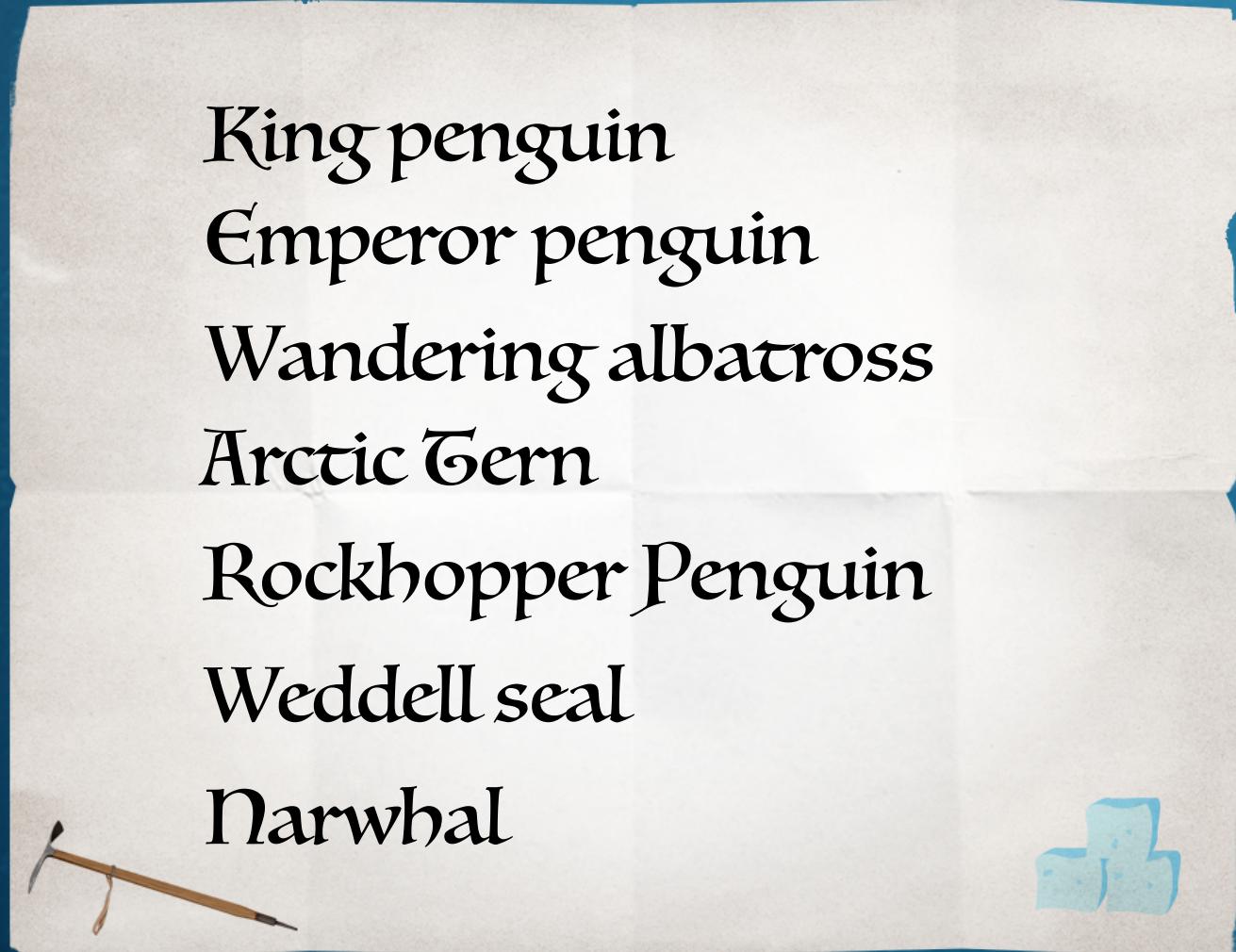
02 Crew Emails

03 Confirmative

→ 04 Multi-line Strings

05 Capture Groups

Finding More Than Just 1 Match

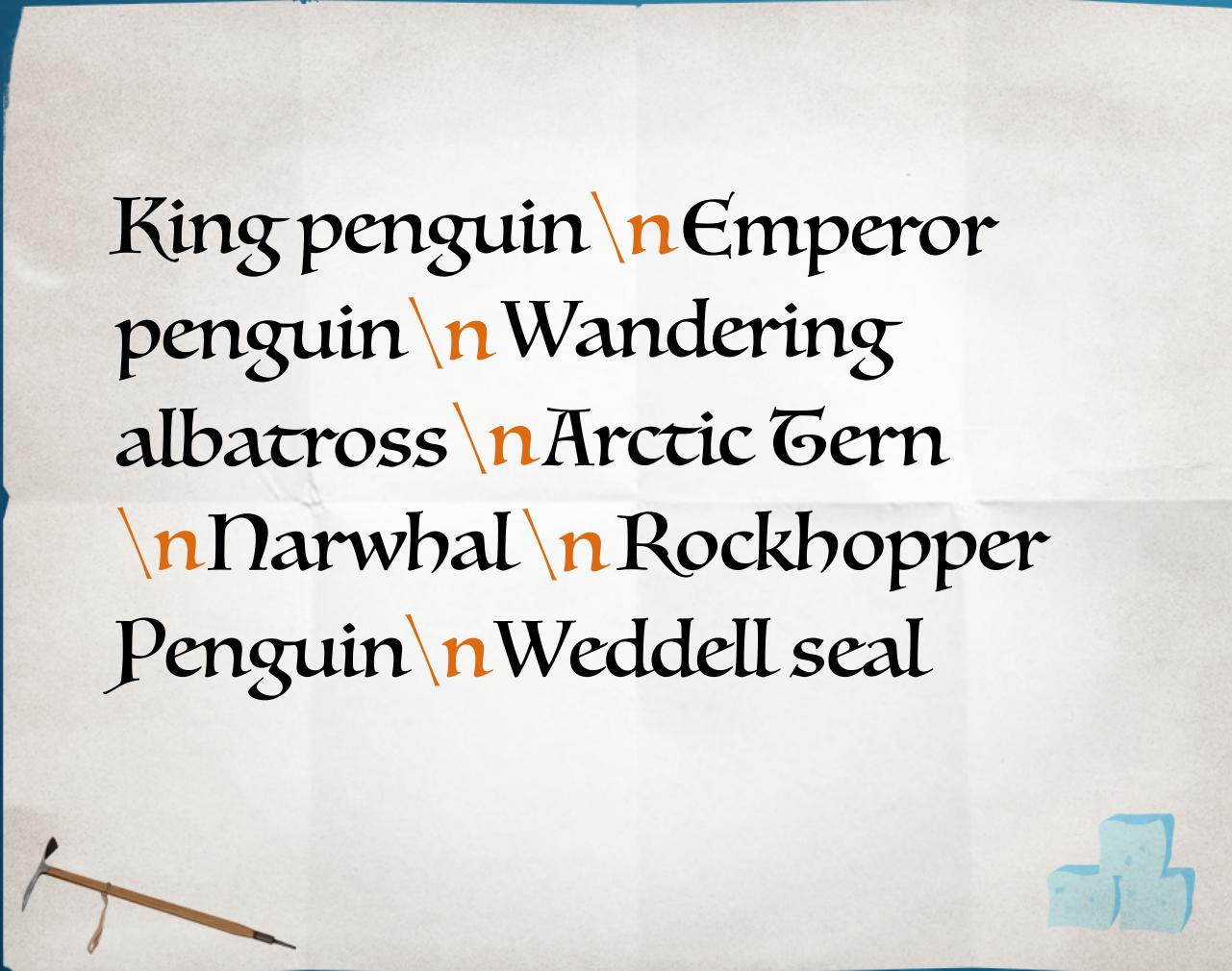


Single text file

We want to find all the birds



Looking at the List of Birds as a Single, Long String



\n

Multiline strings are delimited by the newline character



WARNING! Different operating systems might use other newline characters.

Analyze the Subject String

Subject

```
King penguin
Emperor penguin
Wandering albatross
Arctic Tern
Narwhal
Rockhopper Penguin
Weddell seal
```

Has multiple lines separated by newline characters

Each animal name is 1-2 words

All animal names have mixed casing

Matching the First Animal Name

Regular Expression

/penguin/i

Subject

King penguin

Emperor penguin

Wandering albatross

Arctic Tern

Narwhal

Rockhopper Penguin

Weddell seal

Match the first

occurrence of "penguin"

But we want to match
all penguins, not just
the first

Pattern Goals

- ✓ penguin
- albatross
- tern
- full name
- nothing before or after
- ✗ all occurrences

Global Modifier

Regular Expression

/penguin/ig

Subject

King penguin

Emperor penguin

Wandering albatross

Arctic Tern

Narwhal

Rockhopper Penguin

Weddell seal

Match “penguin” as many times as possible with global modifier

Pattern Goals

- ✓ penguin
- albatross
- tern
- full name
- nothing before or after
- ✓ all occurrences

Matching the Full Name of Just Penguin Animals

Regular Expression

```
/\w+\spenguin/ig
```

Subject

King penguin
Emperor penguin
Wandering albatross
Arctic Tern
Narwhal
Rockhopper Penguin
Weddell seal

Pattern Goals

- ✓ penguin
- albatross
- tern
- ✓ full name
- nothing before or after
- ✓ all occurrences

Our Anchors Aren't Working!

Regular Expression

```
/^\w+spenguin$/ig
```

Now nothing is matching!

Pattern Goals

- ✖ penguin
- ✖ albatross
- ✖ tern
- ✖ full name
- ✖ nothing before or after
- ✖ all occurrences

Subject

```
King penguin
Emperor penguin
Wandering albatross
Arctic Tern
Narwhal
Rockhopper Penguin
Weddell seal
```

Anchors in a Normal Subject String

Subject

King penguin
Emperor penguin
Wandering albatross
Arctic Tern
Narwhal
Rockhopper Penguin
Weddell seal



anchors to beginning
of string



anchors to end of string



We were hoping it would do this!

Anchors in a Multiline String

Regular Expression

```
/^\w+spenguin$/mig
```

Subject

King penguin
Emperor penguin
Wandering albatross
Arctic Tern
Narwhal
Rockhopper Penguin
Weddell seal

↖
anchors to beginning
of line and not entire
subject

Pattern Goals

- ✓ penguin
- ✓ albatross
- ✓ tern
- ✓ full name
- ✓ nothing before or after
- ✓ all occurrences

↘
anchors to end of
line and not entire
subject

m = multiline modifier

Capturing the Full Animal Name

Regular Expression

```
/^\w+\s(penguin|albatross)$/mig
```



Adds “albatross” to search pattern

Subject

King penguin
Emperor penguin
Wandering albatross
Arctic Tern
Narwhal
Rockhopper Penguin
Weddell seal

Pattern Goals

- ✓ penguin
- ✓ albatross
- tern
- ✓ full name
- ✓ nothing before or after
- ✓ all occurrences

Capturing the Full Animal Name

Regular Expression

```
/^\w+\s(penguin|albatross|tern)$/mig
```

Subject

King penguin
Emperor penguin
Wandering albatross
Arctic Tern
Narwhal
Rockhopper Penguin
Weddell seal

Adds “tern” to search pattern



Pattern Goals

- ✓ penguin
- ✓ albatross
- ✓ tern
- ✓ full name
- ✓ nothing before or after
- ✓ all occurrences

Contents

01 The String Story

02 Crew Emails

03 Confirmative

04 Multi-line Strings

→ 05 Capture Groups

Every Crew Has a Home



1 Reindeer Lane, North Pole, AK 99705
120 East 4th Street, Juneau, AK 99705



BREAKING THE ICE
WITH
REGULAR EXPRESSIONS

Building the Number and Street Name

Regular Expression

```
\d+\s[\w\s]+\w{4,6},\s
```

Some numbers and 2 words followed by a comma and a space

Possible Subjects

1 Reindeer Lane,



Number and street name fully matched

Building the City Pattern

Regular Expression

```
[\\w\\s]+, \\s
```

city name of any number of characters followed by a comma and a space

Possible Subjects

North Pole,



city is completely matched

Building the State Pattern

Regular Expression

```
\w{2}\s
```

2-letter state followed by a space

Possible Subjects

```
AK
```



State is completely matched

Building the Zip Pattern

Regular Expression

\d{5}

5-digit ZIP

Possible Subjects

99705



ZIP code is completely matched

Our Fully Matching Pattern

Regular Expression

```
/^\d+\s[\w\s]+[^\s]{4,6},\s[\w\s]+,\s[^\s]{2}\s\d{5}$/ig
```



Anchor to beginning of subject



Anchor to end of subject

Possible Subjects

1 Reindeer Lane, North Pole, AK 99705



Matching Groups

/learnbydoing)/

Possible Subjects

learnbydoing

Match Groups

1.

learnbydoing

Matching Groups

```
/learn(bydoing)/
```

Possible Subjects

```
learnbydoing
```

Match Groups

1.

```
learnbydoing
```

2.

```
bydoing
```

Each group is returned

Matching Groups

/learn((by)doing)/

Possible Subjects

learnbydoing

Match Groups

1.

learnbydoing

2.

bydoing

3.

by



Each group returns captured matches



Each group is returned

Matching Groups — Street

Regular Expression



only return group of first section

```
/^\d+\s[\w\s]+\w{4,6}),\s[\w\s]+,\s\w{2}\s\d{5}$/i
```

Possible Subjects

1 Reindeer Lane, North Pole, AK 99705



Match Groups

1.

1 Reindeer Lane



Matching Groups — City

Regular Expression

```
/^\d+\s[\w\s]+\w{4,6}),\s([\w\s]+),\s\w{2}\s\d{5}$/i
```

Possible Subjects

1 Reindeer Lane, North Pole, AK 99705



Match Groups

1.

1 Reindeer Lane

2.

North Pole

Matching Groups — State

Regular Expression

```
/^\d+\s[\w\s]+\w{4,6}),\s([\w\s]+),\s(\w{2})\s\d{5}$/i
```

Third match group
for 2-letter state

Possible Subjects

1 Reindeer Lane, North Pole, AK 99705 

Match Groups

1.

1 Reindeer Lane

2.

North Pole

3.

AK



Matching Groups — Zip

Regular Expression

```
/^\d+\s[\w\s]+\w{4,6}),\s([\w\s]+),\s(\w{2})\s(\d{5})$/i
```

Final match group for zip code

Possible Subjects

1 Reindeer Lane, North Pole, AK 99705 

Match Groups

1.

1 Reindeer Lane

2.

North Pole

3.

AK

4.

99705



Problem: Need to Restrict Potential Streets

Regular Expression



Should only match “street” or “lane”

```
/^\d+\s[\w\s]+\w{4,6}),\s([\w\s]+),\s(\w{2})\s(\d{5})$/i
```

Possible Subjects

1 Reindeer Lane, North Pole, AK 99705 

15 Discovery Road, Juneau, AK 99705 



We only want to match “street” or “lane”

Solution: Restricting the Street

Regular Expression

```
/^\d+\s[\w\s]+(street|lane),\s([\w\s]+),\s(\w{2})\s(\d{5})$/i
```



Now we're only looking for "street" or "lane"

Possible Subjects

1 Reindeer Lane, North Pole, AK 99705



15 Discovery Road, Juneau, AK 99705



Now we do not match on "road"

Problem: Group Returns Unwanted Result

Regular Expression

Group used for evaluation purposes
returns unwanted value

```
/^(\d+\s[\w\s]+(street|lane)),\s([\w\s]+),\s(\w{2})\s(\d{5})$/i
```

Possible Subjects

1 Reindeer Lane, North Pole, AK 99705



Match Groups

1.

1 Reindeer Lane

2.

Lane



3.

North Pole

4.

AK

5.

99705

we don't want to capture just the type of street!

Making the Street Type a Non-capturing Group

Regular Expression

```
/^\d+\s[\w\s]+(?:street|lane),\s([\w\s]+),\s(\w{2})\s(\d{5})$/i
```



Possible Subjects

1 Reindeer Lane, North Pole, AK 99705



Match Groups

1.

1 Reindeer Lane

2.

North Pole

3.

AK

4.

99705



Now we have the correct number of groups