

# Exploratory Text Analysis

# How can you analyse text?

⊗ Matching ↗ Regular Expressions  
    ↗ LIKE .match()  
SQL

\* Extracting "1<sup>1</sup>2|1<sup>2</sup>2|1<sup>3</sup>00"  
    ↗ Regex .findall()

\* Replacement Regex .sub()  
    "£15.00" → "€15.00"

... These techniques all  
use Explicit patterns  
provided by the analyst

Eg.

Regular Expression      ↴ Slash

FIND " [0-9][0-9] / [0-9][0-9]"

two digits      two digits

# Aside: Advanced Statistical Techniques

⇒ Patterns not provided by analyst

“Natural language Processing”

④ Sentiment Analysis

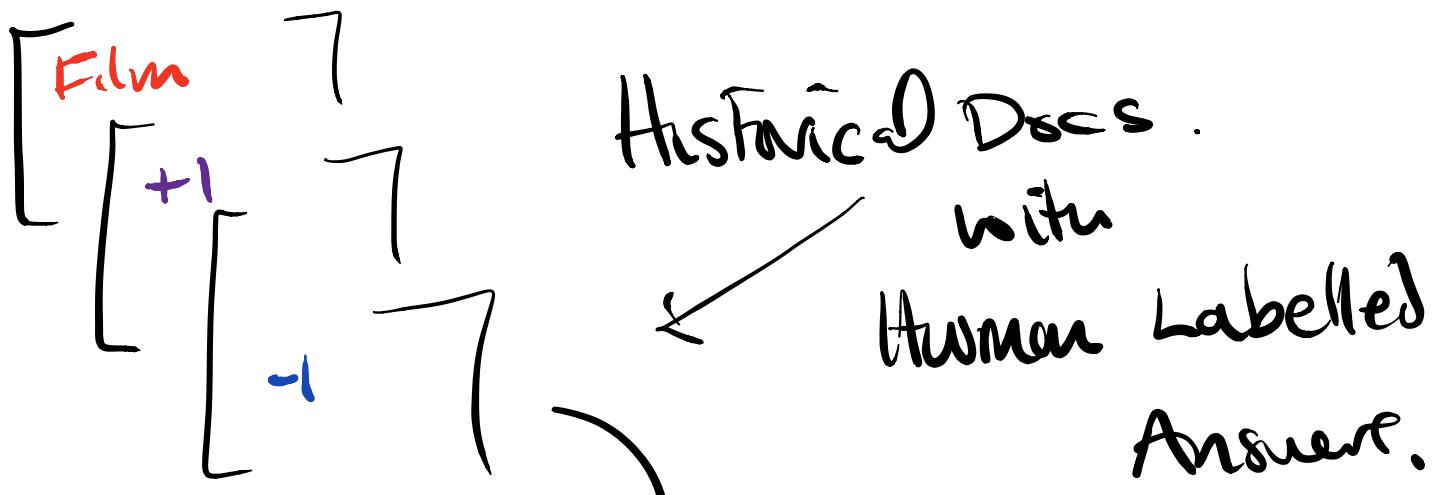
“I had a great time” → +1

④ Topic Analysis

“I went to see Marvel” → “Cinema”

NB. NLP “Text” <sup>Machine</sup> → Interpretation

# Aside : NLP



Eg. word FREQUENCIES

$f(x_{\text{FILM}}, x_{\text{SCIFI}}, x_{\dots}) = y$

# What are Regular Expressions?

a language for describing  
text patterns

Used to -

- Match
- Extract
- Replace.

# Examples : #1 Repeated Nums.

"My BP is 20 mmHg; I'm paying £10.00 for magnesium"

<sup>REPEATED</sup>  
[0-9]+ mm Hg

Number      "mmHg"



# Examples : #2

“My food came to: £100.50,  
I spent £10.00 on the  
taxi!”

; £ [0-9]<sup>+</sup>. [0-9]<sup>+</sup>

“: £” Num                  Num

REPEATED                  REPEATED

# Regular Expression Syntax

Pattern	means	Eg
a £	a	Abba
.	any single symbol	£10.00
[012]	any single sym. in 0, 1, 2	Hello
[0-9]	Any single sym. in 0, ..., 9	
[a-z]	" " " lowercase	
[A-Z]	" " " uppercase .	
a single symbol!		

[^012]

ba\$

^Ab

(May | June)

Alternative

any single symbol

NOT: 0, 1, 2

occurring at the End

Abba Abba

occurring at the start

Abba Abba

EITHER :

May OR  
June

I'm considering  
May, June  
or July

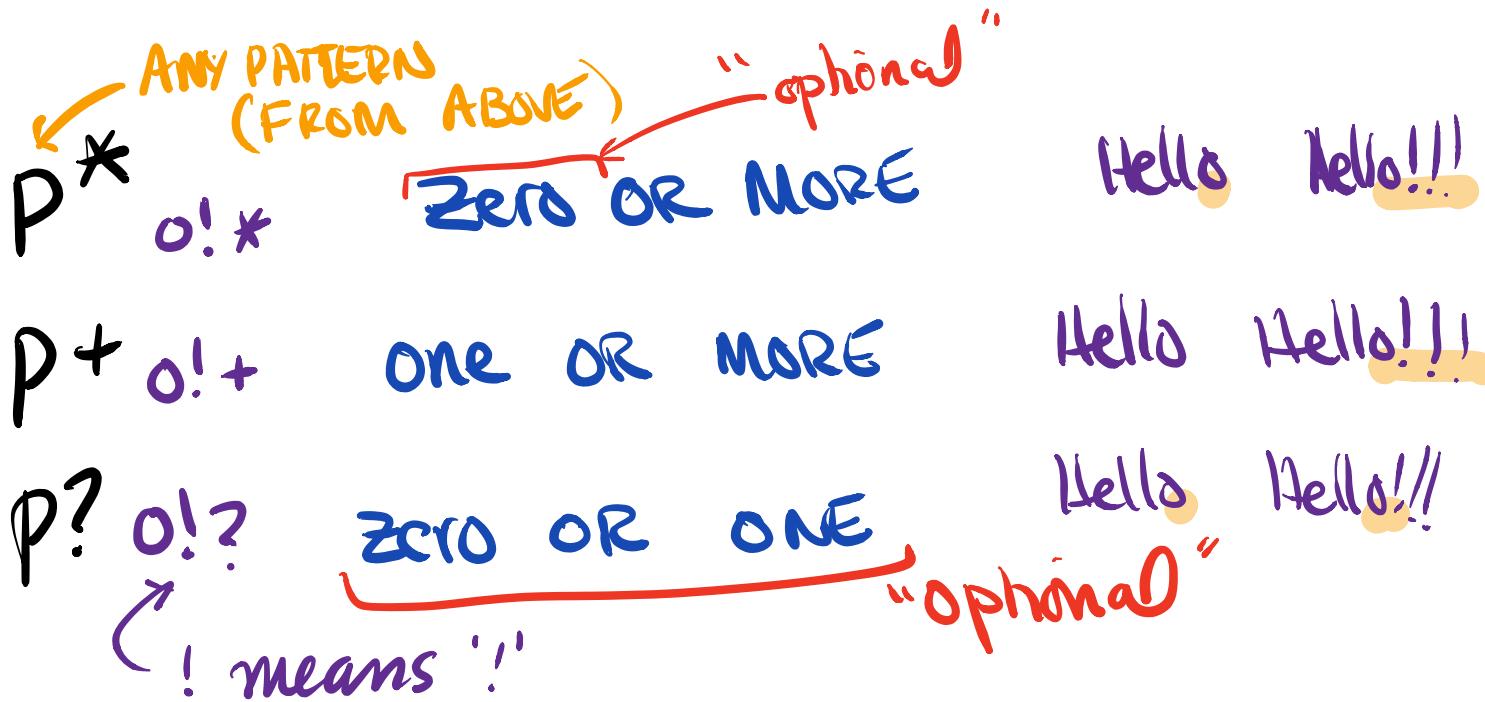
# Quantifier → Quantity

"How Many?"

Note:

prior patterns

DO NOT REPEAT!



# Example : #3 Form Fields

"Name: Michael; Age: 31; Location:   "

: [a-zA-Z0-9]\*<sup>\*</sup> ← Can be Missing

Examples which Match.

{ :;  
:a;  
:a'a'a=  
:01a;

# Aside: Special & literal

Note that '.' means 'ANY symbol'  
and '\$' means 'END'

...  
So How can we say "literally" \$  
or "literally" . ?

TWO WAYS

1. OR [.]

























