[a-z] :   a...z

[0-9] :   0-9

[^Ss] :   neither S nor s   I ...
                                 ~ ete.

State s? :   the previous expression is
                              optional

colou?r :   color  or  colour

oo*h! :   0 or more  previous char

0+h! :   1 or more previous char

FP: incorrectly catches the words

FN: incorrectly misses the words

$$/ \$ [0-9] + /$$

$$/ \$ [0-9] + \backslash . [0-9][0-9] /$$

$$/ \$ [0-9] + (\backslash . [0-9][0-9])? /$$

$$/ (\char94 | \backslash w) \$ [0-9] + (\backslash . [0-9]$$
$$[0-9])?/$$

```
/ [0-9]+GB /

/[0-9]+\. [0-9]+ GB/

/ [0-9]+(\. [0-9]+ +) ? G B/

/ ( ^ | \W ) [0-9]+ (\. [0-9]+ +)?GB
                                    ␣ᵢ*/
```

$/ \backslash b$ [5-9] [0-9] {2,} (\. [0-9]+)?

←* GB \b

500
600
1000 ✗

3 digits

( [5-9] [0-9] [0-9] |

4 or more

[1-9] [0-9] + {3,} )

( (?: Some | a few) ) ( people ( cats )

like some \|/

?: as long as the ?: are there, it won't be registered.

$$\left( \; [\overline{A - Z}_{a-z}] + \right/$$

$$\left/ \; [a-z] + b \; \right/$$

$$\left/ \; b^*(a\ b)^* \; b^* \; \right/$$

( [A-Za-z ]+) \b ( \1 )

13 Blue

/ [0-9]+ [a-z A-Z ]+ /

/\b [0-9] + [a-z A-Z]+ \b/

(^ | \b) [0-9] + [a-z A-Z] +(\b|$)

/

grotto , raven

\b [a-zA-z]grotto [^a-z A-Z] \b

[A-Z a-z]

Word types = # of vocab.

Word instances = Total # of running
                                    words.

$$|V| = k N^{\beta}$$

# of types          # of running words
                                    instances

# Word Normalisation

U.S.A / USA

am, is, be, are

## Lemmatization

Lemma = shared rule of words

## Morphemes

= smallest meaningful units that
make up words

# Morphemes

- stems — core meaning
- affixes — adher to stems

# Stemming (chopping off affixes)

<span style="color:red">simplifies ver. of lemmetization</span>

# Porter Stemmer

ATIONAL → ATE ,

ING → ϵ

SSE S

Sentence seg:

!? is very obvious, but:

"." (period) can be used in Dr. Inc.
(etc...)

→ So, tokenise them first;
┌ either part of the word
└ sentence boundary,

then sentence seg.

# Minimum Edit Distance :

to <mark>measure how similar two strings are.</mark>

→ used for spell Correction.

$$P\left( W_3 \mid W_1, W_2 \right)$$

$$\frac{C\left( W_1, W_2, W_3 \right)}{C\left( W_1, W_2 \right)}$$

| | # | E | X | E | C | U | T | I | O | N |
|---|---|---|---|---|---|---|---|---|---|---|
| N | 9 | | | | | | | | | |
| O | 8 | 7 | 8 | 9 | 10 | 11 | 10 | 9 | 8 | 9 |
| I | 7 | 6 | 7 | 8 | 9 | 10 | 9 | 8 | 9 | 10 |
| T | 6 | 5 | 6 | 7 | 8 | 9 | 8 | 9 | 10 | 11 |
| N | 5 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 10 |
| E | 4 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 9 |
| T | 3 | 4 | 5 | 6 | 7 | 8 | 7 | 8 | 9 | 8 |
| N | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 7 | 8 | 7 |
| I | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 6 | 7 | 8 |
| # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| # | E | X | E | C | U | T | I | O | N | |

XECU
~~INXTENTION~~TION    4+4=8

→ EXECUTION

$$\sqrt[N]{\dfrac{1}{P(a_1, a_2)}}$$

$$P(am \mid I)$$

$$= \dfrac{P(I\ am)}{P(I)} = \dfrac{2}{3}$$

$$P(C) = \frac{C\left(\overset{W_1,\,W_2}{\phantom{x}}\right) + 1}{C\left(W_2\right) + V}$$

just plain boring (—)

entirely predictable and lack energy (—)

no surprises and very few laughs

(—)

very powerful (+)

the most fun film of the summer (+)

---

predicatable ~~with~~ no fun

$C(-) = 14$   $P(-) =$

$C(+) = 9$   $P(+) =$

$$P\left(\text{predictable} \mid (-)\right)$$

$$= \frac{C\left((-), \text{predictable}\right) + 1}{C\left((-)\right) + 20} = \frac{1+1}{14+20}$$

$$P\left(\text{predictable} \mid (+)\right) \qquad = \frac{1}{17}$$

$$= \frac{C\left(+, \text{pre}\right) + 1}{C\left(+\right) + 20} = \frac{0+1}{9+20}$$

$$= \frac{1}{29}$$

$$P\left(\text{no} \mid -\right) = \frac{1+1}{14+20} = \frac{1}{17}$$

$$P\left(\text{no} \mid +\right) = \frac{0+1}{9+20} = \frac{1}{29}$$

$$P\left(\text{fun} \mid -\right) = \frac{0+1}{14+20} = \frac{1}{34}$$

$$P\left(\text{fun} \mid +\right) = \frac{1+1}{9+20} = \frac{2}{29}$$

(-)
$$\left(\frac{1}{17} \times \frac{1}{17} \times \frac{1}{34}\right) \times \frac{3}{5} = 6.1 \times 10^{-5}$$

(+)
$$\left(\frac{1}{29} \times \frac{1}{29} \times \frac{2}{29}\right) \times \frac{2}{5} = 3.2 \times 10^{-5}$$

$P(\text{neg}) = \dfrac{3}{8}$    $c(-) = 9$

$P(\text{mid}) = \dfrac{2}{8} = \dfrac{1}{4}$    $c(N) = 9$

$P(\text{pos}) = \dfrac{3}{8}$    $c(+) = 7$

$V = 24$

Not compelting enough

not interesting enough (-)

failed to impress (-)

dull and uninspiring (-)

it was okay (N)

neither good nor bad, just average (N)

quite compelling (+)

exceptionally good (+)

a thrilling experience (+)

Not

$$P(\text{Not} \mid -) = \frac{1+1}{9+24} = \frac{2}{33}$$

$$P(\text{Not} \mid N) = \frac{0+1}{9+24} = \frac{1}{33}$$

$$P(\text{Not} \mid +) \quad \frac{0+1}{7+24} = \frac{1}{31}$$

---

Compelling

$$P(\text{Com} \mid -) = \frac{1}{9+24} = \frac{1}{33}$$

$$P(\text{Com} \mid N) = \frac{1}{9+24} = \frac{1}{33}$$

$$P(\text{Com} \mid +) = \frac{1+1}{7+24} = \frac{2}{31}$$

$$P\left(\text{enough} \mid (-)\right) = \frac{1+1}{9+24} = \frac{2}{33}$$

$$P\left(\text{enough} \mid (N)\right) = \frac{1}{33}$$

$$P\left(\text{enough} \mid (+)\right) = \frac{1}{31}$$

$$P\left(\text{test} \mid -\right) = \frac{4}{33^3} \times \frac{3}{8}$$
$$= 4.17 \times 10^{-5}$$

$$\mid N) = \frac{1}{33^3} \times \frac{2}{8}$$
$$= 6.95 \times 10^{-6}$$

$$\mid +) = \frac{1}{31^3} \times \frac{3}{8}$$
$$= 1.25 \times 10^{-5}$$

$$( 3, 2, 1, 3, 0, 4.19 )$$

$$( -5, 2.5, -1.2, 0.5, 2.0, 0.7 )$$

$$( -15, 5, -1.2, 1.5, 0, 2.933 )$$

$$-10 + \frac{\cancel{(-1.2 + 1.5)}}{0.3} + 2.933$$

$$3.333 - 10$$

$$-6.7$$

$$-0.67$$