L₁: So it's gonna be forever
L₂: Or it's gonna go down in flames
L₃: You can tell me when it's over

**- Taylor Swift, Blank Space**

## Text Similarity

**tf-idf** is a classical measure of text similarity, allowing us to establish a base similarity metric between similar regions of text.
   **tf**, **t**erm **f**requency asks "how often does the word appear?"
   **idf**, **i**nverse **d**ocument **f**requency asks "how rare is the world?"

## Calculating tf-idf

**First,** we calculate the tf-idf for every word in every line:

(a): How many regions does our document have?    $\#_{regions}$:

(b): How many regions does [**it's**] appear?    $\#(it's)_{regions}$:

(c): How many times does [**it's**] appear in L₁?    **tf** = L₁(it's):

(d): Calculate the idf of [**it's**] on L₁:    **idf** = log( (a) / (b) ):

(e): Calculate tf-idf of [**it's**] on L₁:    **tf-idf** = (c) * (d):

Second, we calculate the **cosine similarity** of each tf-idf:

$$similarity = \cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|} = \frac{\sum\limits_{i=1}^{n} A_i \times B_i}{\sqrt{\sum\limits_{i=1}^{n} (A_i)^2} \times \sqrt{\sum\limits_{i=1}^{n} (B_i)^2}}$$

How often does [**it's**] appear in L1?

| | #regions | #regions(word) | tf L₁ | tf L₂ | tf L₃ | idf L₁ | idf L₂ | idf L₃ | tf-idf L₁ | tf-idf L₂ | tf-idf L₃ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **so** | | | | | | | | | | | |
| **it's** | | | | | | | | | | | |
| **gonna** | | | | | | | | | | | |
| **be** | | | | | | | | | | | |
| **forever** | | | | | | | | | | | |
| **or** | | | | | | | | | | | |
| **go** | | | | | | | | | | | |
| **down** | | | | | | | | | | | |
| **in** | | | | | | | | | | | |
| **flames** | | | | | | | | | | | |
| **you** | | | | | | | | | | | |
| **can** | | | | | | | | | | | |
| **tell** | | | | | | | | | | | |
| **me** | | | | | | | | | | | |
| **when** | | | | | | | | | | | |
| **over** | | | | | | | | | | | |

How do we set up a cosine similarity?