# Building the model II

## 4.0 Calculating word probabilities

Calculate word probabilities

Example: "I am happy because I am learning"

$$P(w) = \frac{C(w)}{V}$$

$$\boxed{P(\text{am}) = \frac{C(\text{am})}{V} = \frac{2}{7}}$$

$P(w)$   Probability of a word

$C(w)$   Number of times the word appears

$V$   Total size of the corpus

| Word | Count |
|------|-------|
| I | 2 |
| am | 2 |
| happy | 1 |
| because | 1 |
| learning | 1 |

Total :   7

Note that you are storing the count of words and then you can use that to generate the probabilities. For this week, you will be counting the probabilities of words occurring.  If you want to build a slightly more sophisticated auto-correct you can keep track of two words occurring next to each other instead. You can then use the previous word to decide. For example which combo is more likely, *there friend* or *their friend?* For this week however you will be implementing the probabilities by just using the word frequencies. Here is a summary of everything you have seen before in the previous two videos.

1. Identify a misspelled word
2. Find strings n edit distance away
   - Insert
   - Delete
   - Switch
   - Replace
3. Filter candidates
4. Calculate word probabilities

$$P(w) = \frac{C(w)}{V}$$

deah → dear ☑
yeah
dear
dean
... *etc*