# Introduction to Word2Vec:

**Uzair Ahmad**

In this tutorial, we'll use the popular `Gensim` library in Python to perform various tasks using word embeddings, specifically the pre-trained Word2Vec model.

## Setup

First, let's get our environment set up. Install the required packages:

```
pip install gensim
```

## Loading Pre-trained Word2Vec Model

For simplicity, we'll use the pre-trained Google News vectors:

```python
from gensim.models import KeyedVectors

# Load the Google News vectors
model = KeyedVectors.load_word2vec_format('https://s3.amazonaws.com/dl4j-distribution/GoogleNews-vectors-negative300.bin.gz', binary=True, limit=100000)
```

Or You can manually download the Google News vectors from the [Gensim GitHub repository](). Once you download it, you can load it using the path to the downloaded file.

```python
import gensim.downloader as api

model = api.load("word2vec-google-news-300")

```

## 1. Finding Similarity

You can find the similarity between two words:

```python
similarity = model.similarity('king', 'queen')
print(f"Similarity between 'king' and 'queen': {similarity:.4f}")
```

```
man:woman::king:queen
```

## 2. Finding Analogies

The classic example is: man👧:king:?

```
1   result = model.most_similar(positive=['woman', 'king'], negative=['man'])
2   print(f"man:woman::king:{result[0][0]}")
```

```
1   man:woman::king:queen
```

## 3. Finding Synonyms (Most Similar Words)

Let's say we want the 5 most similar words to 'computer':

```
1   synonyms = model.most_similar('computer', topn=5)
2   for word, score in synonyms:
3       print(f"{word}: {score:.4f}")
```

```
1   computers: 0.7979
2   laptop: 0.6640
3   laptop_computer: 0.6549
4   Computer: 0.6473
5   com_puter: 0.6082
```

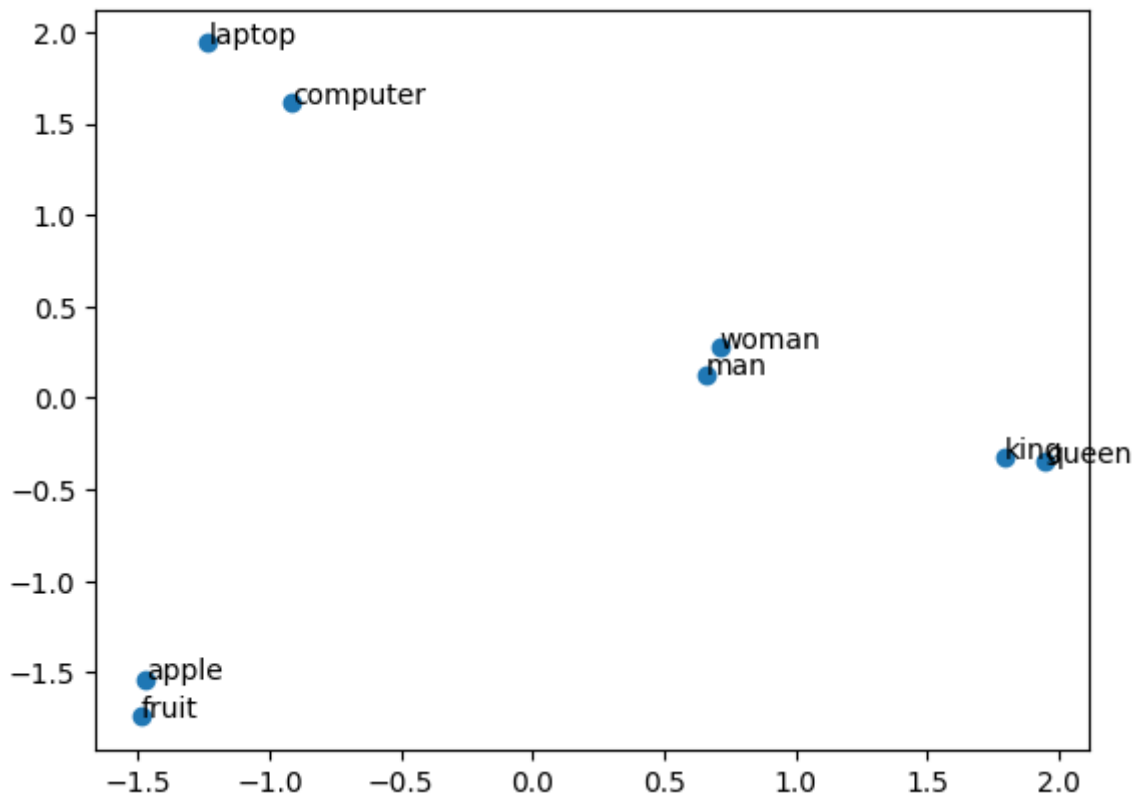## 4. Odd-One-Out

Find the word that doesn't belong:

```
1   odd_word = model.doesnt_match(['apple', 'banana', 'cherry', 'computer'])
2   print(f"The odd one out is: {odd_word}")
```

```
1   The odd one out is: computer
```

## 5. Visualizing Word Embeddings

Using PCA (Principal Component Analysis), you can reduce the dimensionality of word vectors and plot them:

```
1   import matplotlib.pyplot as plt
2   from sklearn.decomposition import PCA
3
4   words = ['king', 'queen', 'man', 'woman', 'computer', 'laptop', 'fruit',
        'apple']
5   vectors = [model[word] for word in words]
6
7   # Reduce dimensions
8   pca = PCA(n_components=2)
9   result = pca.fit_transform(vectors)
10
11  # Plotting
12  plt.scatter(result[:, 0], result[:, 1])
13  for i, word in enumerate(words):
14      plt.annotate(word, xy=(result[i, 0], result[i, 1]))
15  plt.show()
```

This tutorial provided a brief overview of how you can work with word embeddings using the Gensim library. There's a lot more you can do with these embeddings, from using them in deep learning models to improving information retrieval systems. Happy coding!