# Understanding Integrated Gradients

Mannheim is a beautiful city

[CLS] | mannheim | is | a | beautiful | city | [SEP]

Vector representation

| 101 | 25116 | 2003 | 1037 | 3376 | 102 |

Step 2 ✗     Goal

Step 1 ✗

Start

## BASELINE

| 0 | 0 | 0 | 0 | 0 | 0 |

Choose a baseline that does not affect the model classification, for instance, a sequence of zeros

The algorithm calculates the model prediction for each input sequence at each step and compares it to the baseline. The difference in the model's prediction at each step is then multiplied by the corresponding step size. We can then derive the integrated gradients for each input feature by summing up the results of these calculations to identify how important each feature is for the prediction.

Cosima_meyer