# Text Classification

Ch.5 Text Analytics with Python. Dipanjan Sarkar. 2019. Apress
Ch.2 NLP with Python. Steven Bird et al. 2009. O'Reilly Media
Ch.1  Mastering NLP with Python. Deepti Chopra et al. 2016. Packt

# Text Classification Definition

$D$  – a set of records $\{X_1,...,X_N\}$

$C$ –  a set of labels $\{c_1...c_n\}$

$T$ – a Text classification system

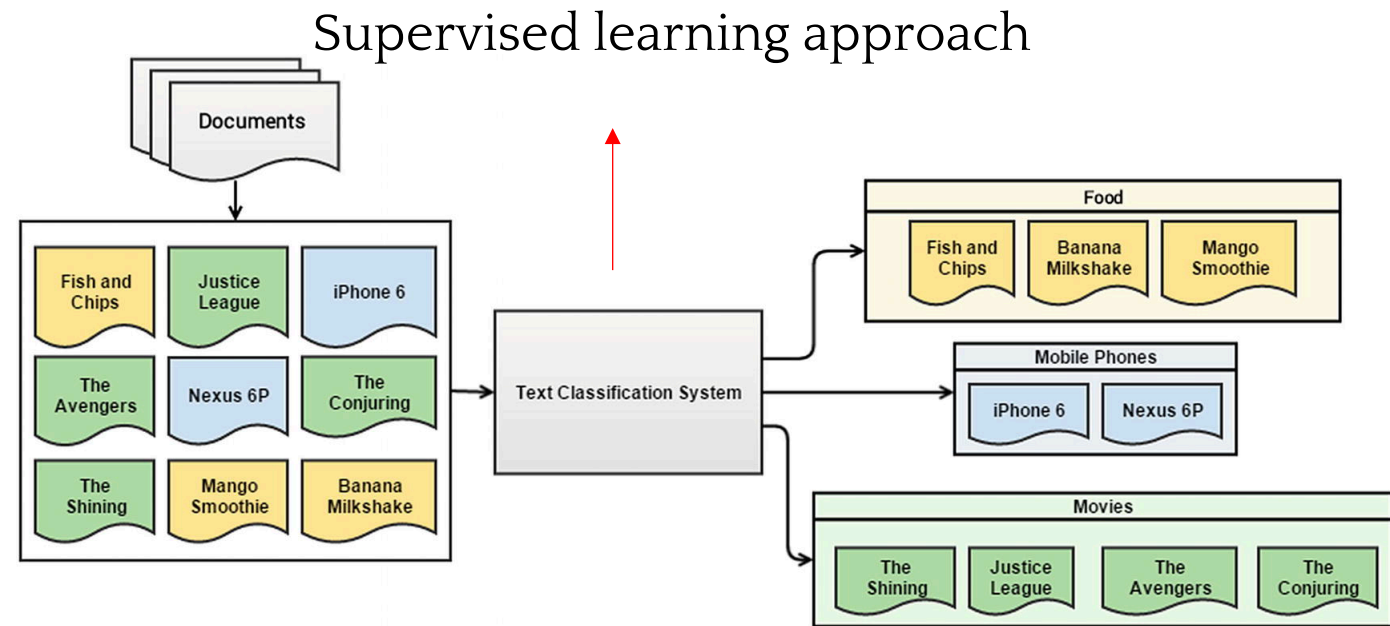$$T : D \rightarrow C_x$$

The training model
predicts a class label

Supervised learning approach



(Dipanjan Sarkar. 2019. Ch.5)

(Charu Aggarwal, 2014, Ch.11)

# Text Classification

**Text Classification Variants**

**Content-based classification**

Topic Weights (% of content to determine the document class)

**Request-based classification**

User behavior (requests)

**Text Classification Approaches**

**Supervised machine learning**

Requires training on prelabeled data samples (training data). Model is used to predict labels in future test data.

**Unsupervised machine learning**

Does not require training on prelabeled data samples. The focus is more on pattern mining and finding latent substructures in the data.

**Classification**

(Categorical variables)

**Regression**

(Continuous variables)