# Introduction to Vector Space Classification

## 1 Vector Space Model

The vector space model is a framework used to represent text documents as vectors in a multi-dimensional space. Each term in the document vocabulary corresponds to a dimension, and documents are points in this space. Typically, each document vector is weighted using a measure like *tf-idf* (term frequency-inverse document frequency), which assigns values based on term relevance and frequency across documents. This model allows for efficient document comparisons using vector math.

### Example

Consider a document containing terms like *machine*, *learning*, and *algorithm*. In a vector space, this document might have a vector representation such as:

$$\vec{d} = (0.4, 0.2, 0.1)$$

where each value is the *tf-idf* score of each term in the document. The higher the *tf-idf* score, the more significant the term is for that document.

## 2 Document Representation

Documents are represented as vectors in a space where each dimension corresponds to a unique term from the vocabulary. This model is efficient because it reduces text to numerical data, allowing for comparison of document similarities using measures like cosine similarity or Euclidean distance.

### Example

A document discussing *Artificial Intelligence* might have high scores in dimensions representing *AI*, *machine learning*, and *data*. Another document discussing the *History of Computers* may score higher on terms like *computing* and *evolution*.

# 3  Contiguity Hypothesis

The Contiguity Hypothesis is central to the vector space model. It proposes that documents within the same class occupy contiguous regions in the vector space, and these regions do not overlap with those of other classes. This hypothesis allows for clear boundaries between classes, aiding classification.

### Example

Documents on *sports* might cluster around terms like *games*, *teams*, and *scores*, while *technology* articles might cluster around *AI*, *devices*, and *computing*. If regions overlap, it may indicate ambiguity or shared topics between classes, which can affect classification accuracy.

# 4  Classification Methods in Vector Spaces

The document introduces two key classification methods within this vector space framework:

- **Rocchio Classification**: This method involves creating centroids (average vectors) for each class and assigning new documents to the nearest centroid's class. This technique is straightforward and efficient but struggles with irregularly shaped classes.

- **k-Nearest Neighbor (kNN) Classification**: kNN classifies a document based on the most common class among its nearest neighbors in the vector space. This approach can manage more complex class shapes but is computationally intensive.

### Example

In classifying emails, a centroid for *spam* might have high weights on terms like *free*, *offer*, and *win*, while a *non-spam* centroid might emphasize terms relevant to work or personal communication. If a new email's vector is closer to the *spam* centroid, it would be classified as spam.

# 5  Rocchio and k-Nearest Neighbor (kNN) Classification

In this section, we dive into specific classification techniques in vector spaces, focusing on **Rocchio and k-Nearest Neighbor (kNN) Classification**. We'll explore their underlying principles and provide examples to illustrate their use in text classification.

## 5.1 Rocchio Classification

**Rocchio Classification** is a centroid-based classification technique within the vector space model, assigning new documents to the class whose centroid they are closest to. A **centroid** is the central point of all documents in a given class, computed as the average vector of these documents. This technique is most effective when classes form approximately spherical clusters in the vector space.

**Process:**

1. **Calculate Class Centroids**: Each class is represented by a centroid vector, computed as the mean of the vectors of all documents in that class. Mathematically, for class $c$ with documents $D_c$, the centroid $\mu(c)$ is:

$$\mu(c) = \frac{1}{|D_c|} \sum_{d \in D_c} v(d)$$

   where $v(d)$ represents the vector of document $d$ in class $c$.

2. **Assign Document to Class**: A new document is classified by assigning it to the class with the nearest centroid, determined by distance measures such as *Euclidean distance*.

**Example**

Consider two classes, "Sports" and "Politics." If "Sports" documents emphasize terms like *team*, *score*, and *game*, the "Sports" centroid will reflect these terms with high weights. A new document mentioning "team victory" would likely be classified as "Sports" since its vector is closer to the "Sports" centroid than to the "Politics" centroid.

**Advantages**

- **Efficiency**: Rocchio is computationally simple and allows quick classification by pre-computing centroids.

- **Interpretability**: Class centroids provide an interpretable central theme for each class.

**Limitations**

- **Limited Flexibility**: Rocchio assumes that classes have spherical, non-overlapping regions in vector space, making it less suitable for classes with complex or overlapping shapes.

## 5.2   k-Nearest Neighbor (kNN) Classification

Unlike Rocchio, **k-Nearest Neighbor (kNN)** classification is a **non-parametric** method that determines class membership based on a document's proximity to $k$ neighboring documents in the vector space. The class is determined by the majority class among the $k$ nearest neighbors.

**Process:**

1. **Identify k Nearest Neighbors**: For a given document $d$, find the $k$ closest documents in the training set based on *distance* (often *Euclidean distance* or *cosine similarity*).

2. **Vote on Class Membership**: Assign $d$ to the class that appears most frequently among its $k$ neighbors.

**Example**

Assume $k = 3$, with classes "Technology" and "Business." For a new document about *AI in the business sector*, its three nearest neighbors could be two "Business" documents and one "Technology" document. The document would be classified as "Business" by majority vote.

**Advantages**

- **Adaptability**: kNN can handle non-spherical and irregular class distributions better than centroid-based methods like Rocchio.

- **No Training Needed**: kNN uses the entire training set directly for classification, which can be advantageous if training time is limited.

**Limitations**

- **Efficiency**: kNN is computationally expensive at test time, especially with large datasets, as each document must be compared with all others.

- **Storage Requirements**: Since it doesn't reduce the training set, kNN requires significant memory.

## 5.3   Comparison and Use Cases

- **Rocchio** is best suited for applications where classes are distinct and have roughly spherical clusters in the vector space, such as distinguishing between clearly defined categories (e.g., separating science and literature articles).

- **kNN** is more versatile, performing well in cases with complex, overlapping classes, like sentiment analysis where documents may share vocabulary across categories.

# 6 Linear Classifiers and Their Applications in Text Classification

## 6.1 What are Linear Classifiers?

A linear classifier is a classification model that separates data points (or document vectors) by a linear boundary. In a two-dimensional space, this boundary is represented by a line, while in higher dimensions, it becomes a hyperplane. Linear classifiers assign a document to a particular class based on the side of the hyperplane on which it falls.

**Defining a Linear Classifier:**

The decision boundary of a linear classifier can be represented by the equation:

$$\vec{w} \cdot \vec{x} = b$$

where:

- $\vec{w}$ is the weight vector representing the orientation of the hyperplane,
- $\vec{x}$ is the vector representation of the document,
- $b$ is a bias term that shifts the hyperplane.

  A document $\vec{x}$ is classified into one class if $\vec{w} \cdot \vec{x} > b$ and another if $\vec{w} \cdot \vec{x} \leq b$.

**Example:**

Consider two classes, "Positive Sentiment" and "Negative Sentiment," represented by terms associated with each sentiment. A linear classifier might place terms like *good* and *great* on the positive side, and terms like *bad* and *poor* on the negative side. If a document vector $\vec{x}$ representing *good service* results in $\vec{w} \cdot \vec{x} > b$, it would be classified as "Positive Sentiment."

## 6.2 How Linear Classifiers Work in Vector Spaces

Linear classifiers operate in high-dimensional spaces, with each term in the vocabulary representing a dimension. The classifier tries to establish a hyperplane that best separates the documents of different classes by optimizing the weight vector $\vec{w}$ and the bias term $b$. This hyperplane aims to achieve maximal separation with minimal misclassification.

**Geometric Interpretation:**

In two-dimensional space, a linear classifier can be visualized as a line that divides data points belonging to different classes. In higher dimensions, this line generalizes to a hyperplane. The decision boundary created by the hyperplane effectively segments the vector space into regions associated with different classes.

**Figure Example (Hypothetical):**

Imagine a two-dimensional space where documents are classified as "sports" or "politics." A line could separate the two categories based on term weights. For example, if the x-axis represents terms related to sports and the y-axis represents terms related to politics, a diagonal line might separate sports articles on one side and politics articles on the other.

## 6.3 Applications of Linear Classifiers in Text Classification

Linear classifiers are highly effective in text classification tasks that involve high-dimensional, linearly separable data. Their simplicity and efficiency make them suitable for applications where quick and interpretable results are important.

**Use Cases:**

- **Spam Filtering**: A linear classifier can separate emails into "spam" and "not spam" based on word frequencies.

- **Sentiment Analysis**: In reviews or social media posts, linear classifiers can distinguish between "positive" and "negative" sentiments by analyzing the presence of specific words associated with each sentiment.

- **Topic Categorization**: Articles can be classified into topics (e.g., sports, health, technology) based on term associations, as terms in each topic typically form linearly separable groups.

## 6.4 Advantages and Limitations of Linear Classifiers

**Advantages:**

- **Efficiency**: Linear classifiers are computationally efficient, especially in high-dimensional text data, and can be trained quickly compared to complex nonlinear models.

- **Interpretability**: The weight vector $\vec{w}$ provides insight into which terms are most indicative of each class, aiding interpretability in applications like topic categorization.

- **Scalability**: Linear classifiers scale well with large datasets, making them ideal for text classification applications with thousands of terms and documents.

**Limitations:**

- **Limited Flexibility**: Linear classifiers assume classes are linearly separable, which may not hold for tasks with complex or overlapping categories

(e.g., multi-faceted topics where documents can belong to multiple categories).

- **Sensitivity to Class Overlap**: Linear classifiers may struggle with classes that have shared vocabulary or content overlap, leading to potential misclassification.

## 6.5  Linear vs. Nonlinear Classification

Nonlinear classifiers, such as kNN and support vector machines with nonlinear kernels, can classify data that isn't linearly separable by creating more complex decision boundaries. However, in high-dimensional text classification, linear classifiers often perform well because of the natural separation that occurs in high-dimensional spaces.

### Example:

For the task of language identification in text, linear classifiers can often separate languages based on unique word patterns in each language, as each language's vocabulary and syntax typically form distinct clusters.

# 7  Nonlinear Classifiers and the Bias-Variance Tradeoff

## 7.1  Nonlinear Classifiers in Text Classification

Unlike linear classifiers, **nonlinear classifiers** can separate classes that don't have a simple linear boundary between them. These classifiers are particularly useful when the classes are more complex, multimodal, or overlap in various regions of the vector space.

### k-Nearest Neighbor (kNN) as a Nonlinear Classifier

kNN, which we discussed earlier, is a prime example of a nonlinear classifier. It does not assume any specific shape for class boundaries. Instead, it adapts to the local structure of the data, drawing complex, non-linear boundaries around classes based on the nearest neighbors.

**Example**: In a dataset with multiple overlapping classes, such as topics with nuanced distinctions (e.g., "sports" and "entertainment" articles that discuss celebrities from both fields), kNN can accurately assign a class label based on the local clustering of similar documents.

### Nonlinear Classifier Applications

- **Sentiment Analysis with Overlapping Sentiment Classes**: Sentiment can vary across a spectrum (positive, neutral, negative), where

certain reviews might mix sentiments (e.g., "The product is good, but delivery was bad"). Nonlinear classifiers can better model this nuance.

- **Multifaceted Document Classification**: For documents that cover multiple subjects or themes (e.g., a "sports technology" article), nonlinear classifiers can capture this multifaceted nature, whereas linear models might misclassify based on dominant terms alone.

## 7.2 The Bias-Variance Tradeoff

The **bias-variance tradeoff** is a fundamental concept in machine learning that reflects the balance between a model's ability to generalize to new data and its accuracy on training data. This tradeoff impacts the choice of classifier, especially in text classification where high-dimensional data can introduce significant noise.

### Understanding Bias and Variance

- **Bias** refers to the error introduced by assuming a simplified model structure. High-bias models make strong assumptions (like linear separability) and often fail to capture complex patterns, leading to **underfitting**.

- **Variance** measures the model's sensitivity to fluctuations in the training set. High-variance models are highly flexible, capturing complex patterns, but they may also fit noise, leading to **overfitting**.

**Example**:

- A linear classifier like Rocchio has **high bias** (it assumes linear separability) but **low variance** (it is relatively stable across different training sets).

- kNN, a nonlinear model, has **low bias** (it doesn't assume a specific shape for class boundaries) but **high variance** (classification decisions can vary significantly based on small changes in the data).

### The Tradeoff in Text Classification

Balancing bias and variance is crucial in text classification due to the high dimensionality of text data:

- **High-dimensional vector spaces** increase the chance of linear separability (benefiting high-bias models).

- **Noisy or sparse data** increases the risk of overfitting for high-variance models, making it essential to carefully choose the classifier based on the dataset's complexity and size.

**Choosing a Classifier Based on the Bias-Variance Tradeoff**:

1. **If the dataset is small or the classes are clearly separable** (e.g., formal documents categorized by specific subjects), a linear classifier (high bias, low variance) may be optimal.

2. **If the dataset is large and contains overlapping or multimodal classes** (e.g., social media sentiment with mixed emotions), a nonlinear classifier like kNN or a neural network (low bias, high variance) might perform better.

**Example: Sentiment Classification Tradeoff**

For a large dataset of product reviews with nuanced sentiments, a high-variance model (such as a neural network or kNN) could capture subtle differences, like distinguishing between "very satisfied" and "slightly satisfied." However, a simpler, high-bias model like a linear classifier might be sufficient and even preferable if the goal is to identify broad sentiment categories (e.g., positive vs. negative) with a limited dataset.

## 7.3 Managing the Bias-Variance Tradeoff

Several strategies help balance the tradeoff:

- **Regularization**: Adding a penalty for model complexity can control variance in nonlinear models, reducing the risk of overfitting.

- **Feature Selection**: Reducing the dimensionality by selecting only the most informative terms helps reduce both bias and variance, as it simplifies the model while retaining key features.

- **Cross-validation**: This method helps gauge a model's performance on unseen data, ensuring that it doesn't overfit or underfit during training.

**Example of Regularization in Text Classification**:

In logistic regression (a linear classifier), regularization helps by assigning lower weights to terms that aren't significant predictors, allowing the classifier to generalize better and control variance.

# 8 Classification in Multi-Class and Multi-Label Settings

## 8.1 Multi-Class Classification

In **multi-class classification**, a document is assigned to **exactly one of several classes** that are mutually exclusive. Multi-class problems are common in text classification, where documents need to be sorted into distinct categories, such as "Sports," "Politics," or "Entertainment."

**Approaches to Multi-Class Classification**

1. **One-vs-All (OvA)**: This method trains multiple binary classifiers, each distinguishing one class from all others. For a test document, each classifier predicts whether it belongs to its designated class, and the document is assigned to the class with the highest confidence score.

   **Example**: For a dataset with "Sports," "Politics," and "Health," OvA trains one binary classifier for "Sports vs. Not Sports," another for "Politics vs. Not Politics," and so on.

2. **One-vs-One (OvO)**: This method trains a binary classifier for every pair of classes. For $J$ classes, OvO requires $\frac{J(J-1)}{2}$ classifiers, each trained to distinguish between two classes. At test time, each classifier "votes" on the class, and the document is assigned to the class with the majority vote.

   **Example**: With classes "Sports," "Politics," and "Health," OvO requires classifiers for "Sports vs. Politics," "Sports vs. Health," and "Politics vs. Health."

**Applications**: Multi-class classification is suitable for:

- **Topic categorization**: Where documents belong to a single topic.

- **Language detection**: Where each document is in one language from a defined set.

## 8.2 Multi-Label Classification

**Multi-label classification** differs from multi-class classification in that a document can belong to **multiple classes simultaneously**. This approach is essential for documents that span multiple themes or categories, as they cannot be neatly assigned to a single class.

**Approaches to Multi-Label Classification**

1. **Independent Binary Classifiers**: A binary classifier is trained for each class, deciding whether or not a document belongs to that class independently of the other classes. This approach allows each classifier to make its decision without interference from other labels.

   **Example**: In a research article dataset, a document could be classified under both "Machine Learning" and "Healthcare" if it discusses AI applications in medicine.

2. **Chained Classifiers**: In this approach, each classifier's output is used as an additional feature for subsequent classifiers. This method helps capture label dependencies by leveraging information from one classifier to aid others.

**Example**: For a multi-label dataset of news articles, if an article is classified under "Finance," it may have a higher probability of also being classified under "Economics" due to the natural dependency between these labels.

3. **Ensemble Methods**: Multi-label classification can benefit from ensemble techniques like **random forests** or **gradient boosting** that combine predictions from multiple classifiers to improve accuracy.

**Applications**: Multi-label classification is essential for:

- **Tagging systems**: For instance, an article on social media might be tagged with "Technology," "Innovation," and "Startups."

- **Sentiment analysis**: Documents may carry multiple sentiments (e.g., positive and negative sentiments toward different aspects of a product).

## 8.3 Adapting Linear and Nonlinear Classifiers for Multi-Class and Multi-Label Settings

**Linear Classifiers**

- **One-vs-All Linear Classifiers**: Linear models like logistic regression and support vector machines are easily adapted to multi-class classification through OvA or OvO methods. They perform well in high-dimensional spaces where classes are separable by linear boundaries.

- **Multi-Label with Independent Binary Linear Classifiers**: In multi-label settings, linear classifiers are trained independently for each label, creating simple and interpretable models that perform well with sparse data.

**Nonlinear Classifiers**

- **kNN and Decision Trees**: Nonlinear classifiers can model complex boundaries for both multi-class and multi-label problems, handling overlapping classes and dependencies between labels more effectively than linear models.

- **Neural Networks**: For multi-label classification, neural networks with sigmoid activation in the output layer can model dependencies between labels and provide probabilistic outputs for each class.

**Example**: In a multi-label classification task for product reviews, a neural network might classify a review as containing both "positive sentiment" (towards product quality) and "negative sentiment" (towards delivery), offering a nuanced understanding of the document.

## 8.4 Evaluation Metrics for Multi-Class and Multi-Label Classification

Evaluation metrics vary depending on whether the classification is multi-class or multi-label.

**Multi-Class Metrics**

- **Accuracy**: Measures the percentage of correct predictions out of all predictions.

- **Confusion Matrix**: A table that shows correct classifications and misclassifications, helping identify where models perform well or need improvement.

**Multi-Label Metrics**

- **Hamming Loss**: Measures the fraction of labels incorrectly predicted.

- **Subset Accuracy**: Measures how often the entire set of labels for a document is predicted correctly.

- **Precision, Recall, and F1 Score**: Calculated for each label and then averaged, either as a macro-average (equal weight for each label) or micro-average (weighted by label frequency).

**Example of Hamming Loss**: If a document with true labels $[1, 0, 1, 1]$ is predicted as $[1, 1, 0, 1]$, the Hamming loss is $\frac{1}{4} = 0.25$, as two out of four labels are incorrect.