



Hình 2: Ví dụ về các bước tiền xử lý dữ liệu

2.2 Biểu diễn dữ liệu

Mô hình Bag of Words (BoW) là một mô hình được sử dụng phổ biến trong lĩnh vực phân loại văn bản. Mô hình này thường sử dụng để xử lý ngôn ngữ tự nhiên, được dùng để biểu diễn tài liệu, xem tài liệu là một tập hợp các từ (words) mà không quan tâm đến thứ tự cũng như cấu trúc cú pháp của chúng.

Một văn bản được biểu diễn dạng véc-tơ (có n thành phần là các từ tương ứng) mà giá trị thành phần thứ j là tần số xuất hiện từ thứ j trong văn bản. Nếu xét tập D gồm m văn bản và từ điển có n từ vựng, thì D có thể được biểu diễn thành bảng D kích thước $m \times n$, dòng thứ i của bảng là véc-tơ biểu diễn văn bản thứ i tương ứng.

Giả sử dữ liệu có 15.000 tweets với 20.000 đặc trưng (từ vựng), thông thường mỗi tweet sẽ được lưu trữ như sau:

Chi mục	1	2	3	...	20000
Tần số xuất hiện	0	1	0	...	0

Nếu mỗi tần số xuất hiện của 1 đặc trưng tốn khoảng 2 bytes để lưu trữ, vậy mỗi tweet tốn 40.000 bytes. 15.000 tweets tốn 600.000.000 bytes.

Đối với nghiên cứu này, chúng tôi đề xuất cách lưu trữ tiết kiệm bộ nhớ, tương tự như LibSVM [Chang & Lin, 2011], chỉ lưu những từ có tần số xuất hiện lớn hơn 0. Cách lưu trữ như sau:

<label> <index-1>:<value-1> <index-2>:<value-2> ...

Trong đó:

<label> là lớp ban đầu của tweet, 1 là tích cực, 0 là tiêu cực.