# Sentiment Analysis and Basic Feature Extraction

the movie was <u>great</u>! would <u>watch again</u>!

the film was <u>awful</u>; I'll never <u>watch again</u>!

① text $\overline{x} \Rightarrow f(\overline{x})$ feature extraction

② $\{f(\overline{x}^{(i)}), y^{(i)}\}_{i=1}^{D}$ dataset of $D$ labeled exs,

$\qquad\qquad\qquad \Rightarrow$ train classifier

# Feature extraction

the movie was great

Bag-of-words: Assume 10,000 words in vocabulary

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$$
the   a   of   at  ... movie ... was ... great ...

4 1s
9996 0s

Counts (how many "the" are present)

presence/absence (0/1)

n-gram: sequence of n consecutive words

Bag-of-ngrams

2-grams: <u>the movie</u>, <u>movie was</u>, <u>was great</u>

tf-idf
↳ tf × idf

tf: count of the term

idf: inverse document frequency $\log \dfrac{N \quad \text{total # docs}}{\{D: w \in D\} \quad \text{# docs with } w}$

# Preprocessing

① Tokenization

was great!
was greaT

→ was ␣ great ␣ !

[ ... great .... great !...]

wasn't → was ␣ n't

② [Sometimes] Stopword removal (the, of, a, ...)

③ [Sometimes] Casing (lowercasing, truecasing)

④ Handling unknown words    Durrett ⟹ UNK

⑤ Indexing: map each {word, n-gram} into $\mathbb{N}$

use a map