

# Machine Learning and Real World Data

## Sentiment Classification

Is a standard task in NLP.

### Definitions

- A **type** is a unique word.
- A **token** is an instance of a type.
- A **tokeniser** split text into tokens.

## Naive Bayes Classification

**Machine learning** is a program that adapts its behaviour after being exposed to new data: without explicit programming, and implicitly from the data alone.

### Definitions

- **Features** are observable properties of the data.
- **Classes** are labels associated with the data (e.g. positive/negative)
- **Classification** is a function that maps features to classes.

The classifier is given a set of features  $O$  and classes  $c_1, \dots, c_n \in C$ , and gives  $P(c_i|O)$  for each  $c_i$ .

- $O = \{w_1, w_2, \dots, w_n\}$  are the words in the review.
- $C = \{\text{pos}, \text{neg}\}$

Choose the  $P(c_i|O)$  with the highest probability

$$\begin{aligned} c_{NB} &= \operatorname{argmax}_{c \in C} P(c|O) \\ &= \operatorname{argmax}_{c \in C} \frac{P(O|c)P(c)}{P(O)} \\ &= \operatorname{argmax}_{c \in C} P(O|c)P(c) \end{aligned}$$

As  $P(O)$  is a constant and does not affect  $\operatorname{argmax}$ .

$$\begin{aligned} P(O|c) &= P(w_1|c) \times P(w_2|c) \times \cdots \times P(w_n|c) \\ c_{NB} &= \operatorname{argmax}_{c \in C} P(c) \prod_{i=1}^n P(w_i|c) \\ &= \operatorname{argmax}_{c \in C} \left( \log P(c) + \sum_{i=1}^n \log P(w_i|c) \right) \end{aligned}$$

Summing results in less floating point precision errors than multiplying.

In training, collect information needed to calculate  $P(c)$  and  $P(w_i|c)$

$$\begin{aligned} P(c) &= \frac{N_c}{N_{\text{total}}} \\ P(w_i|c) &= \frac{\text{count}(w_i \text{ in } c)}{\sum_{w \in V} \text{count}(w \text{ in } c)} \end{aligned}$$

**Laplace smoothing** prevents  $P(w_i|c)$  from being zero by adding 1 to the count of each type.

$$\begin{aligned}
P(w_i|c) &= \frac{\text{count}(w_i \text{ in } c) + 1}{\sum_{w \in V} (\text{count}(w \text{ in } c) + 1)} \\
&= \frac{\text{count}(w_i \text{ in } c) + 1}{\sum_{w \in V} \text{count}(w \text{ in } c) + |V|}
\end{aligned}$$


---

Name	Description
Zipf's law	$f_w \approx \frac{k}{(r_w + \beta)^\alpha}$ <ul style="list-style-type: none"> <li><math>f_w</math> is the frequency of the word</li> <li><math>r_w</math> is the frequency rank of the word</li> <li><math>k, \alpha</math> and <math>\beta</math> are language dependent constants</li> </ul> <p>Note: <math>\beta</math> is the rank shift.</p>
Heap's law	$u_n = kn^\beta$ <ul style="list-style-type: none"> <li><math>u_n</math> is the number of types (vocabulary size)</li> <li><math>n</math> is the number of tokens</li> <li><math>\beta</math> and <math>k</math> are language dependent constants</li> </ul>

In Naive Bayes Classification, only seen types receive a probability estimate. Adding 1 redistributes some probability mass to unseen types.

---

### Significance Testing

- The **null hypothesis**: the two result sets comes from the same distribution.
- Rejecting the null hypothesis means the observed results is unlikely to have happened by chance.

Choose a **significance level**  $\alpha$ , reject the null hypothesis if the probability of observing the event under the null hypothesis is less than  $\alpha$ .

In a **binomial distribution**  $B(N, q)$

$$\begin{aligned}
P(X = k) &= \binom{N}{k} q^k (1 - q)^{N-k} \\
P(X \leq k) &= \sum_{i=0}^k \binom{N}{i} q^i (1 - q)^{N-i}
\end{aligned}$$

A **two-tailed test** tests if the two systems performs equally well: the  $\alpha$  in each tail is halved.

	Actual = same	Actual = different
Predicted = same	Correct	Type II error: <ul style="list-style-type: none"> <li>Use a more powerful test (e.g. permutation test rather than sign test)</li> <li>Use more data</li> </ul>
Predicted = different	Type I error	Correct

Significance testing cannot show two distributions are the same.

### Note

For testing sentiment classifiers, ignoring ties will lead to the null hypothesis being incorrectly rejected. Add 0.5 to the count of positive and negative results in case of ties.

## Overtraining

Overtraining is where more training makes the classifier perform worse on unseen data.

Am I overtraining?

- If you are using large amounts of new test data, not overtraining.
- If incrementally improving the classifier on the same small test data, overtraining.

Overtraining is caused by finding characteristic features of each class that are hard to generalise.

## N-fold cross-validation

1. Split data into  $N$  equal folds.
2. For each fold  $X$ , train on other folds, test on fold  $X$  only.
3. Average all the accuracy for the final accuracy.

It is good if each splits performs equally well, calculate the variance:

$$\text{var} = \frac{1}{n} \sum_i^n (x_i - \mu)^2$$

Consider the  $N$  experiments as one overall experiment.

Cross validation	Description
Stratified cross-validation	Each split mirrors the distribution of classes in the overall data.
Jack-knitting	Each individual data point is a split.
Dependency-sensitive cross-validation	Fold in a way that known characteristics of a data are isolated (e.g. one split per genre)

Cross-validation does not solve the problem of overtraining.

Instead, a **validation corpus** (a separate set of data not used for training or testing) can be used to

- Tweak parameters before training
- Check if training is making the system perform worse on the validation corpus (is it overtraining?)

## Uncertainty and Agreement

$$\overline{P}_a \text{ observed agreement} = \frac{1}{N} \sum_{i=0}^{N-1} \frac{\#\text{observed pairs of agreements on item } i}{\#\text{possible pairs}}$$

where  $N$  is the number of items to be classified.

$$\overline{P}_e \text{ chance agreement} = \sum_{c \in C} P(c)^2$$

## Fleiss' Kamma

$$\kappa = \frac{\overline{P}_a - \overline{P}_e}{1 - \overline{P}_e}$$

- $\kappa = 1$  then complete agreement
- $\kappa = 0$  then no agreement beyond what is expected by chance
- $\kappa = 0.8$  means very good agreement

#### Note

$\kappa$  can be negative.

## Social Networks

#### Definitions

- **Distance** is the length of the shortest path between two nodes.
- The **diameter** of a graph is the maximum distance between any pair of nodes.
- The **degree** of a node is the number of neighbours it has.

Natural networks often have the **small world property**.

#### Definition

A **small world network** is one with

- Not many connections
- There are closely clustered regions
- Clusters are connected by only a few links

The measurable characteristics are

- **High clustering coefficient:** a node's neighbours are likely to be neighbours of each other.

#### Definitions

- **Triadic closure** is if  $A \leftrightarrow B$  and  $A \leftrightarrow C$ , then very likely  $B \leftrightarrow C$ .
- **Global clustering coefficient** is

$$\frac{\text{\# of closed triads}}{\text{\# of possible closed triads}}$$

- **Small average path length**, typically grows logarithmically with network size.
- **Sparse connectivity:** has few edges compared to the fully connected graph.

#### Definitions

- **Giant component** is a connected component containing most of the nodes in a graph.
- **Weak ties** are distant links, opposite of a **strong link**.

#### Note

Links that keep two giant components together are often weak ties.

- **Bridge** is an edge that connects two components which would otherwise be unconnected.
- **Local bridge** is an edge when removed, the two nodes will have no neighbours in common. (shortest path between the two nodes increases)

## Betweenness Centrality

### Definition

A **gatekeeper** are crucial edges in linking densely connected regions in the graph.

Define

- $\sigma(s, t)$  : number of shortest path between nodes  $s$  and  $t$

$$\sigma(s, t) = \sum_{u \in \text{pred}(t)} \sigma(s, u)$$

where  $\text{pred}(t) = \{u \mid (u, t) \in E \wedge d(s, t) = d(s, u) + 1\}$

- $\sigma(s, t|v)$  number of those paths passing through  $v$

### Note

- If  $s = t$ , then  $\sigma(s, t) = 1$
- If  $v \in \{s, t\}$ , then  $\sigma(s, t|v) = 0$

### Definitions

Betweenness centrality of  $v$

$$C_B(v) = \sum_{s, t \in V} \frac{\sigma(s, t|v)}{\sigma(s, t)}$$

Another way to calculate  $C_B(v)$  :

Define **pairwise dependencies**  $\delta(s, t|v)$  and **one-sided dependencies**  $\delta(s|v)$

$$\begin{aligned} \delta(s, t|v) &= \frac{\sigma(s, t|v)}{\sigma(s, t)} \\ \delta(s|v) &= \sum_{t \in V} \delta(s, t|v) \\ &= \sum_{\substack{(v, w) \in E \\ w: d(s, w) = d(s, v) + 1}} \frac{\sigma(s, v)}{\sigma(s, w)} \cdot (1 + \delta(s|w)) \\ C_B(v) &= \sum_{s \in V} \delta(s|v) \end{aligned}$$

### Algorithm 1: Brandes Algorithm

```

1: for  $s \in V$  do
2:    $\triangleright$  Initialisation
3:   for  $w \in V$  do
4:      $\text{pred}[w] \leftarrow$  empty list
5:      $\text{dist}[w] \leftarrow \infty$ 
6:      $\sigma[w] \leftarrow 0$ 
7:      $\delta[w] \leftarrow 0$ 
8:   end
9:    $Q \leftarrow$  empty list
10:   $Q \leftarrow \text{enqueue } s$ 
11:   $S \leftarrow$  empty stack
12:

```

---

```

13:   $\triangleright$  Calculate  $\sigma$  for all  $v \in V$ 
14:  while  $Q$  not empty do
15:     $v \leftarrow$  dequeue  $Q$ 
16:     $S \leftarrow$  push  $v$ 
17:
18:    for  $w : (v, w) \in E$  do
19:       $\triangleright$  Finding  $w$  for the first time
20:      if  $\text{dist}[w] = \infty$  then
21:         $\text{dist}[w] \leftarrow \text{dist}[v] + 1$ 
22:         $Q \leftarrow$  enqueue  $w$ 
23:      end
24:
25:       $\triangleright$  Update number of shortest paths through  $w$ 
26:      if  $\text{dist}[w] = \text{dist}[v] + 1$  then
27:         $\sigma[w] \leftarrow \sigma[w] + \sigma[v]$ 
28:         $\text{pred}[w] \leftarrow$  append  $v$ 
29:      end
30:    end
31:
32:     $\triangleright$  calculate  $C_B[w]$ 
33:    while  $S$  not empty do
34:      pop  $w \leftarrow S$ 
35:      for  $v \in \text{pred}[w]$  do
36:         $\delta[v] \leftarrow \delta[v] + \frac{\sigma[v]}{\sigma[w]} \cdot (1 + \delta[w])$ 
37:      end
38:      if  $w \neq s$  then
39:         $C_B[w] \leftarrow C_B[w] + \delta[w]$ 
40:      end
41:    end
42:  end
43: end

```

---

For bidirectional graphs, divide the value of each  $C_B[w]$  by 2.

---