

Information Theory

CS5154/6054

Yizong Cheng

10/13/2022

13 *Text classification and Naive Bayes*

13.5 Feature selection

13.5.1 Mutual information

13.5.2 χ^2 Feature selection

Odds Ratio and Retrieval Status Value

$$(11.18) \quad c_t = \log \frac{p_t(1 - u_t)}{u_t(1 - p_t)}$$

(11.19)

documents		relevant	nonrelevant	Total
Term present	$x_t = 1$	s	$\text{df}_t - s$	df_t
Term absent	$x_t = 0$	$S - s$	$(N - \text{df}_t) - (S - s)$	$N - \text{df}_t$
Total		S	$N - S$	N

Using this, $p_t = s/S$ and $u_t = (\text{df}_t - s)/(N - S)$ and

$$(11.20) \quad c_t = K(N, \text{df}_t, S, s) = \log \frac{s/(S - s)}{(\text{df}_t - s)/((N - \text{df}_t) - (S - s))}$$

When a Term Has Nothing to Do with Relevance

Table 2	R=1 (Relevant)	R=0 (on-relevant)	
Term t present	N11	N10	
Term t absent	N01	N00	

(11.19)

documents		relevant	nonrelevant	Total
Term present	$x_t = 1$	s	$df_t - s$	df_t
Term absent	$x_t = 0$	$S - s$	$(N - df_t) - (S - s)$	$N - df_t$
Total		S	$N - S$	N

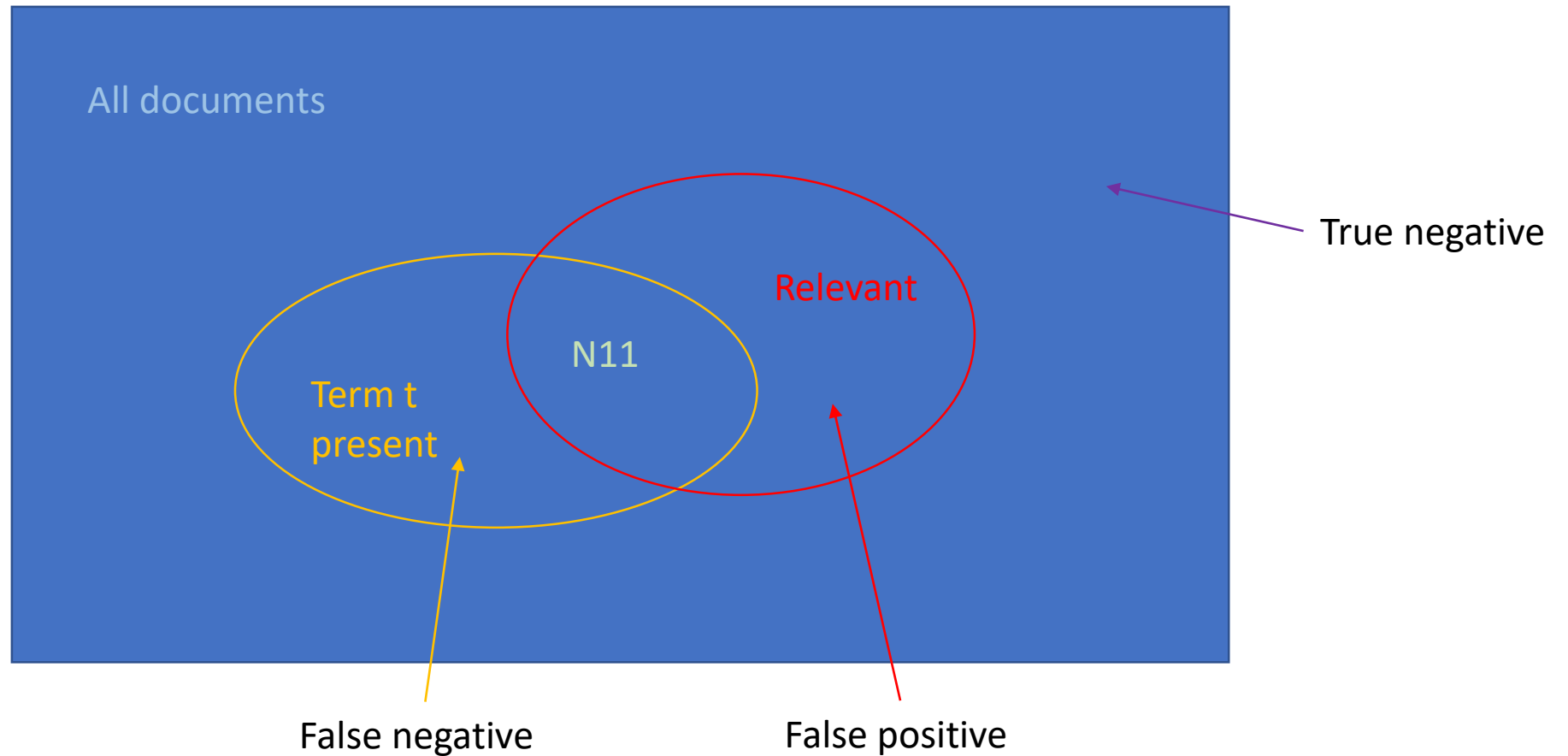
Using this, $p_t = s/S$ and $u_t = (df_t - s)/(N - S)$ and

(11.20)
$$c_t = K(N, df_t, S, s) = \log \frac{s/(S - s)}{(df_t - s)/((N - df_t) - (S - s))}$$

$$C_t = \log (N11/N01)/(N10/N00)$$

$$C_t = 0 \text{ iff } N11 \cdot N00 = N01 \cdot N10$$

Two Ways to Partition (Cluster) Documents



Expected Counts for Independent Events

Actual	R=1 (Relevant)	R=0 (on-relevant)	
Term t present	N11	N10	N1.
Term t absent	N01	N00	N0.
	N.1	N.0	N

S of N documents are relevant and df_t of N documents contains term t. The expected number of document both relevant and containing t is $S df_t / N = (N11 + N01)(N11 + N10)/N$. If this expected number is also the actual N11, then events t present and relevance are independent.

Quiz problem: under what condition do we have $(N11 + N01)(N11 + N10)/N = (N11 + N01)(N11 + N10)/(N11 + N01 + N10 + N00) = N11$?

Expected Counts from Observed Counts

Example 13.3: Consider the class *poultry* and the term export in Reuters-RCV1. The counts of the number of documents with the four possible combinations of indicator values are as follows:

	$e_c = e_{poultry} = 1$	$e_c = e_{poultry} = 0$
$e_t = e_{export} = 1$	$N_{11} = 49$	$N_{10} = 27,652$
$e_t = e_{export} = 0$	$N_{01} = 141$	$N_{00} = 774,106$

Example 13.4: We first compute E_{11} for the data in Example 13.3:

$$\begin{aligned}
 E_{11} &= N \times P(t) \times P(c) = N \times \frac{N_{11} + N_{10}}{N} \times \frac{N_{11} + N_{01}}{N} \\
 &= N \times \frac{49 + 141}{N} \times \frac{49 + 27652}{N} \approx 6.6
 \end{aligned}$$

where N is the total number of documents as before.

We compute the other $E_{e_t e_c}$ in the same way:

	$e_{poultry} = 1$	$e_{poultry} = 0$
$e_{export} = 1$	$N_{11} = 49 \quad E_{11} \approx 6.6$	$N_{10} = 27,652 \quad E_{10} \approx 27,694.4$
$e_{export} = 0$	$N_{01} = 141 \quad E_{01} \approx 183.4$	$N_{00} = 774,106 \quad E_{00} \approx 774,063.6$

X^2 : Difference between Expected and Observed

$$(13.18) \quad X^2(\mathbb{D}, t, c) = \sum_{e_t \in \{0,1\}} \sum_{e_c \in \{0,1\}} \frac{(N_{e_t e_c} - E_{e_t e_c})^2}{E_{e_t e_c}}$$

STATISTICAL
SIGNIFICANCE

X^2 is a measure of how much expected counts E and observed counts N deviate from each other. A high value of X^2 indicates that the hypothesis of independence, which implies that expected and observed counts are similar, is incorrect. In our example, $X^2 \approx 284 > 10.83$. Based on Table 13.6, we can reject the hypothesis that *poultry* and *export* are independent with only a 0.001 chance of being wrong.⁸ Equivalently, we say that the outcome $X^2 \approx 284 > 10.83$ is *statistically significant* at the 0.001 level. If the two events are dependent, then the occurrence of the term makes the occurrence of the class more likely (or less likely), so it should be helpful as a feature. This is the rationale of χ^2 feature selection.

An arithmetically simpler way of computing X^2 is the following:

$$(13.19) \quad X^2(\mathbb{D}, t, c) = \frac{(N_{11} + N_{10} + N_{01} + N_{00}) \times (N_{11}N_{00} - N_{10}N_{01})^2}{(N_{11} + N_{01}) \times (N_{11} + N_{10}) \times (N_{10} + N_{00}) \times (N_{01} + N_{00})}$$

This is equivalent to Equation (13.18) (Exercise 13.14).

Statistical Significance (p-Value)

► **Table 13.6** Critical values of the χ^2 distribution with one degree of freedom. For example, if the two events are independent, then $P(X^2 > 6.63) < 0.01$. So for $X^2 > 6.63$ the assumption of independence can be rejected with 99% confidence.

p	χ^2 critical value
0.1	2.71
0.05	3.84
0.01	6.63
0.005	7.88
0.001	10.83

FileEditViewHistoryBookmarksToolsHelp

sklearn.feature_selection.chi2 — ×

+

←→↺🏠

🔒https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.chi2.html

📄150%☆

🔔S🔍

☰

^

sklearn.feature_selection.chi2

sklearn.feature_selection.chi2(X, y) [source]

Compute chi-squared stats between each non-negative feature and class.

This score can be used to select the `n_features` features with the highest values for the test chi-squared statistic from `X`, which must contain only non-negative features such as booleans or frequencies (e.g., term counts in document classification), relative to the classes.

Recall that the chi-square test measures dependence between stochastic variables, so using this function “weeds out” the features that are the most likely to be independent of class and therefore irrelevant for classification.

Read more in the [User Guide](#).

Parameters:

X : {array-like, sparse matrix} of shape (n_samples, n_features)
Sample vectors.

y : array-like of shape (n_samples,)
Target vector (class labels).

Returns:

chi2 : ndarray of shape (n_features,)
Chi2 statistics for each feature.

p_values : ndarray of shape (n_features,)
P-values for each feature.

Toggle Menu

From Counts to Probability Distribution

Actual	R=1 (Relevant)	R=0 (on-relevant)	
Term t present	N_{11}/N	N_{10}/N	$N_{1.}/N$
Term t absent	N_{01}/N	N_{00}/N	$N_{0.}/N$
	$N_{.1}/N$	$N_{.0}/N$	$N/N = 1$

$$N = N_{11} + N_{01} + N_{10} + N_{00}$$

Independent	R=1 (Relevant)	R=0 (on-relevant)	
Term t present	$N_{1.} \cdot N_{.1} / N^2$	$N_{1.} \cdot N_{.0} / N^2$	$N_{1.}/N$
Term t absent	$N_{0.} \cdot N_{.1} / N^2$	$N_{0.} \cdot N_{.0} / N^2$	$N_{0.}/N$
	$N_{.1}/N$	$N_{.0}/N$	$N/N = 1$

Kullback-Leibler Divergence

- A measure of difference between two probability distributions (over the same set of outcomes).
- $KL(p \parallel q) = \sum_k p_k \log (p_k/q_k)$
- Theorem: $KL(p \parallel q) \geq 0$ with equality iff $p = q$.

Mutual Information

- Given two random variables (X = term t present and Y = relevant),
- $I(X; Y) = \text{KL}(p(X, Y) \parallel p(X)p(Y))$
- In our case, both random variables are binary and the estimates of the probabilities are based on counts N_{11} , N_{01} , N_{10} , and N_{00} .
- lir uses U as the random variable for term t present and C as the random variable for relevance.

$$(13.17) \quad I(U; C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_{1.}N_{.1}} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_{0.}N_{.1}} \\ + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_{1.}N_{.0}} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_{0.}N_{.0}}$$

Example 13.3: Consider the class *poultry* and the term *export* in Reuters-RCV1. The counts of the number of documents with the four possible combinations of indicator values are as follows:

	$e_c = e_{poultry} = 1$	$e_c = e_{poultry} = 0$
$e_t = e_{export} = 1$	$N_{11} = 49$	$N_{10} = 27,652$
$e_t = e_{export} = 0$	$N_{01} = 141$	$N_{00} = 774,106$

After plugging these values into Equation (13.17) we get:

$$\begin{aligned}
 I(U;C) &= \frac{49}{801,948} \log_2 \frac{801,948 \cdot 49}{(49+27,652)(49+141)} \\
 &\quad + \frac{141}{801,948} \log_2 \frac{801,948 \cdot 141}{(141+774,106)(49+141)} \\
 &\quad + \frac{27,652}{801,948} \log_2 \frac{801,948 \cdot 27,652}{(49+27,652)(27,652+774,106)} \\
 &\quad + \frac{774,106}{801,948} \log_2 \frac{801,948 \cdot 774,106}{(141+774,106)(27,652+774,106)} \\
 &\approx 0.0001105
 \end{aligned}$$

<i>UK</i>		<i>China</i>		<i>poultry</i>	
london	0.1925	china	0.0997	poultry	0.0013
uk	0.0755	chinese	0.0523	meat	0.0008
british	0.0596	beijing	0.0444	chicken	0.0006
stg	0.0555	yuan	0.0344	agriculture	0.0005
britain	0.0469	shanghai	0.0292	avian	0.0004
plc	0.0357	hong	0.0198	broiler	0.0003
england	0.0238	kong	0.0195	veterinary	0.0003
pence	0.0212	xinhua	0.0155	birds	0.0003
pounds	0.0149	province	0.0117	inspection	0.0003
english	0.0126	taiwan	0.0108	pathogenic	0.0003

<i>coffee</i>		<i>elections</i>		<i>sports</i>	
coffee	0.0111	election	0.0519	soccer	0.0681
bags	0.0042	elections	0.0342	cup	0.0515
growers	0.0025	polls	0.0339	match	0.0441
kg	0.0019	voters	0.0315	matches	0.0408
colombia	0.0018	party	0.0303	played	0.0388
brazil	0.0016	vote	0.0299	league	0.0386
export	0.0014	poll	0.0225	beat	0.0301
exporters	0.0013	candidate	0.0202	game	0.0299
exports	0.0013	campaign	0.0202	games	0.0284
crop	0.0012	democratic	0.0198	team	0.0264

► **Figure 13.7** Features with high mutual information scores for six Reuters-RCV1 classes.

sklearn.metrics.mutual_info_score

```
sklearn.metrics.mutual_info_score(labels_true, labels_pred, *,  
contingency=None)
```

[\[source\]](#)

Mutual Information between two clusterings.

The Mutual Information is a measure of the similarity between two labels of the same data. Where $|U_i|$ is the number of the samples in cluster U_i and $|V_j|$ is the number of the samples in cluster V_j , the Mutual Information between clusterings U and V is given as:

$$MI(U, V) = \sum_{i=1}^{|U|} \sum_{j=1}^{|V|} \frac{|U_i \cap V_j|}{N} \log \frac{N|U_i \cap V_j|}{|U_i||V_j|}$$

Any Two Sets Have a Mutual Information

- The set of retrieved and the set of relevant documents.
- The set of documents containing term t_1 and that containing term t_2 .
 - A distance measure between terms?
 - $I(X; Y)$ may not be the same as $I(Y; X)$.
- The set of terms in document d_1 and that in document d_2 .
 - Divergence between documents?
- Objective for optimization:
 - Modifying a set of parameters so that the mutual information generated with these parameters against a target distribution is minimized.