

Discovering User Interests by Document Classification

Loc Nguyen

University of Science, Ho Chi Minh city, Vietnam
Email: ng_phloc@yahoo.com

User interest is one of personal traits attracting researchers' attention in user modeling and user profiling. User interest competes with user knowledge to become the most important characteristics in user model. Adaptive systems need to know user interests so that provide adaptation to user. For example, adaptive learning systems tailor learning materials (lesson, example, exercise, test...) to user interests. I propose a new approach for discovering user interest based on document classification. The basic idea is to consider user interests as classes of documents. The process of classifying documents is also the process of discovering user interests. There are two new points of view:

- The series of user access in his/her history are modeled as documents. So user is referred indirectly to as “document”.
- User interests are classes such documents are belong to.

Our approach includes four following steps:

1. Documents in training corpus are represented according to **vector model**. Each element of vector is product of term frequency and inverse document frequency. However the inverse document frequency can be removed from each element for convenience.
2. **Classifying training corpus** by applying decision tree or support vector machine or neural network. Classification rules (weight vectors W^*) are drawn from decision tree (support vector machine). They are used as classifiers.
3. Mining user's access history to **find maximum frequent itemsets**. Each itemset is considered a interesting document and its member items are considered as terms. Such interesting documents are modeled as vectors
4. Applying classifiers (see step 3) into these interesting documents in order to choose which classes are most suitable to these interesting documents. **Such classes are user interests**.

This approach bases on document classification but it also relates to information retrieval in the manner of representing documents. Hence section 1 discusses about vector model for representing documents. Support vector machine, decision tree and neural network on document classification are mentioned in section 2, 3, 4. Main technique to discover user interest is described in section 5. Section 6 is the evaluation

1. Vector model for representing documents

Suppose our corpus \mathbf{D} is the composition of documents $D_i \in \mathbf{D} = \{D_1, D_2, \dots, D_m\}$. Every document D_i contains a set of key words so-called *terms*. The number of times a term occurs in a document is called *term frequency*. Given the document D_i and term t_j , the term frequency tf_{ij} measuring the importance of term t_j within document D_i is defined as below:

$$tf_{ij} = \frac{n_{ij}}{\sum_k n_{ik}}$$

Where n_{ij} is the number of occurrences of term t_j in document D_i , and the denominator is the sum of number of occurrences of all terms in document D_i .

Suppose we need to search documents which are most relevant to the query having term t_j . The simple way is to choose documents which have highest term frequency tf_{ij} . However in situation that t_j is not a good term to distinguish between relevant and non-relevant documents and other terms occurring rarely are better ones to distinguish between relevant and non-relevant documents. This will tend to incorrectly emphasize documents containing term t_j more, without giving enough weight to other meaningful terms. So the *inverse document frequency* is a measure of general importance of the term. It is used to decrease the weight of terms occurring frequently and increase the weight of terms occurring rarely. The inverse document frequency of term t_j is the ratio of the size of corpus to the number of documents that t_j occurs.

$$idf_j = \log \frac{|\text{corpus}|}{|\{D: t_j \in D\}|}$$

Where $|\text{corpus}|$ is the total number of documents in corpus and $|\{D: t_j \in D\}|$ is the number of documents containing term t_j . We use \log function to normalize idf_j so that it is less than or equal 1.

The weight of term t_j in document D_i is defined as product of tf_{ij} and idf_j

$$w_{ij} = tf_{ij} * idf_j$$

This weight measure the importance of a term in a document over the corpus. It increases proportionally to the number of times a term occurs in the document but is offset by the frequency of this term in the corpus. In general this weight balances the importance of two measures: term frequency and inverse document frequency.

Suppose there are n terms $\{t_1, t_2, \dots, t_n\}$, each document D_i is modeled as the vector which is composed of weights of such terms.

$$D_i = (w_{i1}, w_{i2}, w_{i3}, \dots, w_{in})$$

Hence the corpus becomes a matrix $m \times n$, which have m rows and n cols with respect to m document and n terms. D_i is called document vector.

The essence of document classification is to use supervised learning algorithms in order to classify corpus into groups of documents; each group is labeled. In this chapter we apply two methods namely support vector machine, decision tree and neural network for document classification.

2. Document classification based on Support Vector Machine

2.1 Support vector machine

Support vector machine (SVM) [Cristianini, Shawe-Taylor 2000] is a supervised learning algorithm for classification and regression. Given a set of n -dimensional vectors in vector space, SVM finds the separating hyper-plane that splits vector space into sub-set of vector; each separated sub-set (so-called data set) is assigned one class. There is the constraint for this separating hyper-plane: “it must maximize the margin between two sub-sets”.

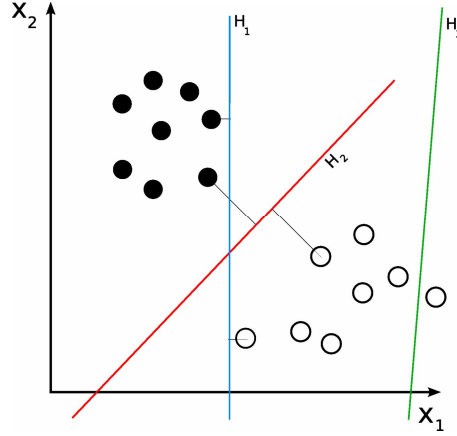


Fig. 1. Separating hyper-planes

Suppose we have some n -dimensional vectors; each of them belongs to one of two classes. We can find many $n-1$ dimensional hyper-planes that classify such vectors but there is only one hyper-plane that maximizes the margin between two classes. In other words, the nearest between a point in one side of this hyper-plane and other side of this hyper-plane is maximized. Such hyper-plane is called maximum-margin hyper-plane and it is considered as maximum-margin classifier.

Let $\{X_1, X_2, \dots, X_n\}$ be the training set of vectors and let $y_i = \{1, -1\}$ be the class label of vector X_i . It is necessary to determine the maximum-margin hyper-plane that separates vectors belonging to $y_i=1$ from vectors belonging to $y_i = -1$. This hyper-plane is written as the set of point satisfying:

$$W^T \otimes X_i + b = 0 \quad (1)$$

Where \otimes denotes the scalar product and W is a weight vector perpendicular to hyper-plane and b is the bias. W is also called perpendicular vector or normal vector. It is used to specify hyper-plane.

The value $\frac{b}{|W|}$ is the offset of the hyper-plane from the origin along the weight vector W .

To calculate the margin, two parallel hyper-planes are constructed, one on each side of the maximum-margin hyper-plane. Such two parallel hyper-planes are represented by two following equations:

$$\begin{aligned} W^T \otimes X_i + b &= 1 \\ W^T \otimes X_i + b &= -1 \end{aligned}$$

To prevent vectors falling into the margin, all vectors belonging to two class $y_i=1$, $y_i=-1$ have two following constraints respectively:

$$\begin{aligned} W^T \otimes X_i + b &\geq 1 \quad (\text{for } X_i \text{ of class } y_i=1) \\ W^T \otimes X_i + b &\leq -1 \quad (\text{for } X_i \text{ of class } y_i=-1) \end{aligned}$$

These constraints can be re-written as:

$$Y_i (W^T \otimes X_i + b) \geq 1 \quad (2)$$

For any new vector X , the rule for classifying it is computed as below:

$$f(X_i) = \text{sign}(W^T \otimes X_i + b) \in \{-1, 1\} \quad (3)$$

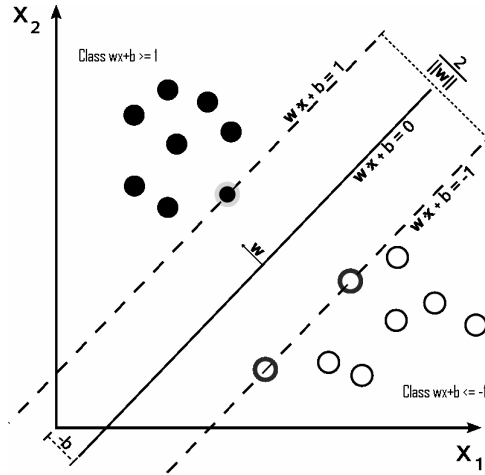


Fig. 2. Maximum-margin hyper-plane

Because maximum-margin hyper-plane is defined by weight vector W , it is easy to recognize that the essence of constructing maximum-margin hyper-plane is to solve the following constrained optimization problem:

$$\underset{W,b}{\text{Minimize}} \frac{1}{2} |W|^2 \text{ subject to } y_i (W^T \otimes X_i + b) \geq 1, \forall i$$

This constraints can be re-written as below:

$$\underset{W,b}{\text{Minimize}} \frac{1}{2} |W|^2 \text{ subject to } 1 - y_i (W^T \otimes X_i + b) \leq 0, \forall i \quad (4)$$

The reason of minimize $\frac{1}{2} |W|^2$ is that distance between two parallel hyper-planes is

$\frac{2}{|W|}$ and we need to maximize such distance in order to maximize the margin for maximum-margin hyper-plane. Then maximizing $\frac{2}{|W|}$ is to minimize $\frac{1}{2} |W|^2$.

Because $|W|$ is the norm of W being complex to compute, we substitute $\frac{1}{2} |W|^2$ with

$\frac{1}{2} |W|^2$. This is the *constrained optimization problem* or *quadratic programming* (QP) *optimization problem*. Note that $|W|^2$ can be computed by scalar product of it and itself:

$$|W|^2 = W^T \otimes W$$

The way to solve QP problem in formula III.3.2 is through its Lagrange dual. Suppose we want to minimize $f(x)$ subject to $g(x)=0$, it exists a solution x_0 to set of equations:

$$\begin{cases} \frac{\partial}{\partial x} (f(x) + \alpha g(x))_{x=x_0} = 0 \\ g(x) = 0 \end{cases} \text{ where } \alpha \text{ is Lagrange multiplier.}$$

For multi constraints $g_i(x) = 0$ ($i=\overline{1,n}$), there is a Lagrange multiplier for each constraint

$$\begin{cases} \frac{\partial}{\partial x} (f(x) + \sum_{i=1}^n \alpha_i g_i(x))_{x=x_0} = 0 \\ g_i(x) = 0 \end{cases}$$

In case of $g_i(x) \leq 0$, there is no change. If x_0 is a solution to the constrained optimized problem:

$$\underset{x}{\text{Minimize}} f(x) \text{ subject to } g_i(x) \leq 0 \forall i$$

Then there must exist $\alpha_i \geq 0 \forall i$ so that x_0 satisfies following equations:

$$\begin{cases} \frac{\partial}{\partial x} (f(x) + \sum_{i=1}^n \alpha_i g_i(x)) = 0 \\ g_i(x) \leq 0 \end{cases} \quad x=x_0 \quad (5)$$

The function $f(x) + \sum_i \alpha_i g_i(x)$ is called Lagrangian denoted L . Let $f(X_i)$ be $\frac{1}{2} \|W\|^2$ and let $g_i(X)$ be $1 - y_i(W^T \otimes X_i + b)$, the constraint in formula III.3.4 becomes:

$$\underset{\lambda}{\text{Maximize}} (\underset{W, b}{\text{Minimize}} L(W, b, \lambda)) \quad (6)$$

where

$$L(W, b, \lambda) = \frac{1}{2} \|W\|^2 + \sum_{i=1}^n \lambda_i (1 - y_i (W^T \otimes X_i + b)) \quad (7)$$

The Lagrangian L is minimized with respect to the primal variables W and b and maximized with respect to the dual variables λ . Setting the gradient of L w.r.t W and b to zero, we have

$$\begin{cases} \frac{\partial L(W, b, \lambda)}{\partial W} = 0 \\ \frac{\partial L(W, b, \lambda)}{\partial \lambda} = 0 \end{cases} \Leftrightarrow \begin{cases} W - \sum_{i=1}^n \lambda_i y_i X_i = 0 \\ \sum_{i=1}^n \lambda_i y_i = 0 \end{cases} \quad (8)$$

Substituting (III.3.8) into (III.3.7) we have:

$$L(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^1 \lambda_i \lambda_j y_i y_j X_i^T X_j \quad (9)$$

Where $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_n)^T$

If (9) is re-written in matrix notation then the constrained optimization problem in (4), (6) leads to dual problem:

$$\text{Maximize } L(\lambda) = \lambda^T I - \frac{1}{2} \lambda^T S \lambda \text{ subject to } \lambda \geq 0 \text{ and } \lambda^T y = 0 \quad (10)$$

Where S is a symmetric $n \times n$ matrix with elements $s_{ij} = y_i y_j X_i^T X_j$.

Suppose λ^* is a solution of (III.3.1) with condition (III.3.8). In other words, $L(\lambda^*)$ is maximized; so the weight vector representing maximum-margin hyper-plane is recovered by λ^* and X_i :

$$W^* = \sum_{i=1}^n \lambda_i y_i X_i \quad (11)$$

So the bias b is computed as below:

$$b^* = y_i - W^{*T} \otimes X_i \quad (12)$$

The rule for classification in (3) becomes:

$$f(X_i) = R = \text{sign}(W^{*T} \otimes X_i + b^*) \quad (13)$$

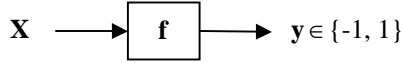


Fig. 3. Classification function

It means that whenever we need to determine to which class a new vector X_i belongs, it is only to substitute X_i into $W^{*T} \otimes X_i + b^*$ and check the value of this expression. If the value is less than or equal -1 (≤ -1) then X_i belongs to class $y_i = -1$. Otherwise, if the value is greater than or equal 1 (≥ 1) then X_i belongs to class $y_i = 1$. Hence the function $(W^{*T} \otimes X_i + b^*)$ is called classification function or classification rule.

The Lagrange multipliers are non-zero when $W^T \otimes X_i + b$ is equal 1 or -1 , vectors X_i in this case are considered support vectors they are closest to the maximum-margin hyper-plane. These vectors lie on parallel hyper-planes. So this approach is called support vector machine.

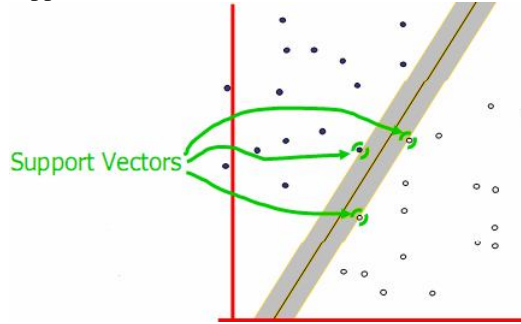


Fig. 4. Support vectors

However the way to find λ^* in (10) is quadratic programming (QP) problem. There are many approaches to solve this problem; for instance, the sequential minimal optimization is described as below:

- A QP with two variables is trivial to solve
- At each iteration, a pair of (λ_i, λ_j) is picked up and QP is solved with these variables; repeat until convergence

2.2. Applying Support Vector Machine into Document Classification

Give corpus $\mathbf{D} = \{D_1, D_2, D_3, \dots, D_m\}$ in which every document D_i is modeled by a *tf-idf* weight vector. Suppose there are n terms $\{t_1, t_2, \dots, t_n\}$, we have:

$$D_i = (d_{i1}, d_{i2}, \dots, d_{in})$$

where d_{ij} is product of term frequency and inverse document frequency $d_{ij} = tf_{ij} * idf_j$

If the index of document is ignored, document D is represented as below:

$$D = (d_1, d_2, \dots, d_n)$$

Given k classes $\{l_1, l_2, \dots, l_k\}$, there is demand of classifying documents into such classes. The technique of classification based on SVM is *two-class* classification in which the classes is $+1, -1$ for $y_i = +1, -1$ respectively. So we need to determine unique hyper-plane referred to as *two-class* classifier. It is possible to extend *two-class* classification to *k-class* classification by constructing k *two-class* classifier. In means that we must specify k couple of optimal weight vector W_i^* and bias b_i^* . Each couple (W_i^*, b_i^*) being a *two-class* classifier is the representation of class l_i . The process of finding (W_i^*, b_i^*) in training corpus \mathbf{D} is described in the section of support vector machine.

Table 1. k couple (W_i^*, b_i^*) corresponds with k class $\{l_1, l_2, \dots, l_k\}$

Class	Weight vector	Bias	Classification rule
l_1	W_1^*	b_1^*	$R_1 = W_1^{*T} \otimes X + b_1^*$
l_2	W_2^*	b_2^*	$R_2 = W_2^{*T} \otimes X + b_2^*$
...
l_k	W_k^*	b_k^*	$R_k = W_k^{*T} \otimes X + b_k^*$

For example, classifying document $D = (d_1, d_2, \dots, d_n)$ is described as below:

1. For each classification rule $R_i = W_i^{*T} \otimes X + b_i^*$, substituting each D into such rule. It means that vector X in such rule is replaced by document D .

$$Expression_i = W_i^{*T} \otimes D + b_i^*$$

2. Suppose there is a sub-set of rules $\{R_{h+1}, R_{h+2}, \dots, R_{h+r}\}$ that the value of expression $Expression_i = W_i^{*T} \otimes D + b_i^*$ is greater than or equal 1. We can conclude that document D is belongs to r classes $\{l_{h+1}, l_{h+2}, \dots, l_{h+r}\}$ where $\{l_{h+1}, l_{h+2}, \dots, l_{h+r}\} \subseteq \{l_1, l_2, \dots, l_k\}$

Table 2. Classifying document D

Classification Expression	Value
$W_1^{*T} \otimes D + b_1^*$	$\geq 1 : D \in l_1$ $\leq -1 : D \notin l_1$
$W_2^{*T} \otimes D + b_2^*$	$\geq 1 : D \in l_2$ $\leq -1 : D \notin l_2$
...	...
$W_k^{*T} \otimes D + b_k^*$	$\geq 1 : D \in l_k$ $\leq -1 : D \notin l_k$

3. Document classification based on Decision Tree

Given a set of classes $\mathbf{C} = \{\text{computer science, math}\}$, a set of terms $\mathbf{T} = \{\text{computer, programming language, algorithm, derivative}\}$ and the corpus $\mathbf{D} = \{\text{doc1.txt, doc2.txt, doc3.txt, doc4.txt, doc5.txt}\}$. The training data is showed in following table in which cell (i, j) indicates the number of times that term j (column j) occurs in document i (row i).

Table 3. Term frequencies of documents

	<i>computer</i>	<i>programming language</i>	<i>algorithm</i>	<i>derivative</i>	<i>class</i>
<i>doc1.txt</i>	5	3	1	1	computer
<i>doc2.txt</i>	5	5	40	5	math
<i>doc3.txt</i>	20	5	20	55	math
<i>doc4.txt</i>	20	55	5	20	computer
<i>doc5.txt</i>	15	15	4	0.3	math
<i>doc6.txt</i>	35	10	45	10	computer

Table 4. Normalized term frequencies

	<i>computer</i>	<i>programming language</i>	<i>algorithm</i>	<i>derivative</i>	<i>class</i>
<i>doc1.txt</i>	0.5	0.3	0.1	0.1	computer
<i>doc2.txt</i>	0.05	0.05	0.4	0.5	math
<i>doc3.txt</i>	0.2	0.05	0.2	0.55	math
<i>doc4.txt</i>	0.2	0.55	0.05	0.2	computer
<i>doc5.txt</i>	0.15	0.15	0.4	0.3	math
<i>doc6.txt</i>	0.35	0.1	0.45	0.1	computer

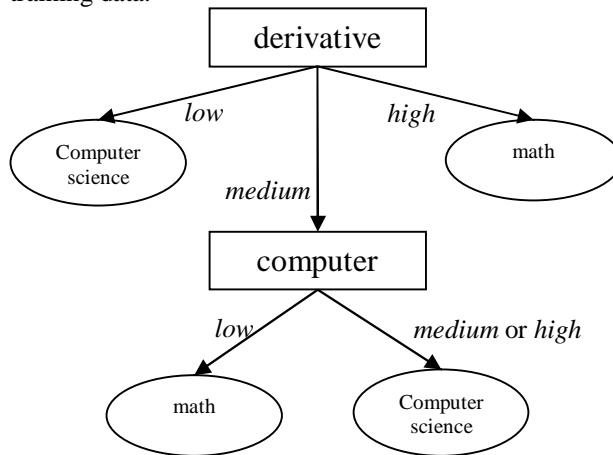
Because the expense of real number computation is so high, all term frequencies are changed from real number into nominal value:

1. $0 \leq \text{frequency} < 0.2$: *low*
2. $0.2 \leq \text{frequency} < 0.5$: *medium*

3. $0.5 \leq \text{frequency}$: *high***Table 5.** Nominal term frequencies

	<i>computer</i>	<i>programming language</i>	<i>algorithm</i>	<i>derivative</i>	class
<i>doc1.txt</i>	high	medium	low	low	computer
<i>doc2.txt</i>	low	low	medium	high	math
<i>doc3.txt</i>	medium	low	medium	high	math
<i>doc4.txt</i>	medium	high	low	medium	computer
<i>doc5.txt</i>	low	low	medium	medium	math
<i>doc6.txt</i>	medium	low	medium	low	computer

The basic idea of generating decision tree [Mitchell 1997] is to split the tree into two sub-trees at the most informative node. Such node is chosen by computing its entropy or information gain. Following figure shows the decision tree generated from our training data.

**Fig. 5.** Decision tree

We can extract classification rules from this decision tree:

Table 6. Classification rules deriving from decision tree induction

<i>Rule 1</i>	If frequency of term “derivative” is <i>low</i> then document belongs to class <i>computer science</i>
<i>Rule 2</i>	If frequency of term “derivative” is <i>medium</i> and frequency of term “computer” is <i>medium</i> or <i>high</i> then document belongs to class <i>computer science</i>
<i>Rule 3</i>	If frequency of term “derivative” is <i>medium</i> and frequency of term “computer” is <i>low</i> then document belongs to class <i>math</i> .
<i>Rule 4</i>	If frequency of term “derivative” is <i>high</i> then document belongs to class <i>math</i>

Suppose the numbers of times that terms *computer*, *programming language*, *algorithm* and *derivative* occur in document D are 5, 1, 1, and 3 respectively. We need to determine which class document D belongs to. D is normalized as term frequency vector.

$$D = (0.5, 0.1, 0.1, 0.3)$$

Changing real number into nominal value, we have:

$$D = (high, low, low, medium)$$

According to rule 2 in above table, D is *computer science* document because in document vector D , frequency of term “derivative” is *medium* and frequency of term “computer” is *high*.

4. Document classification based on Neural Network

4.1. Artificial Neural Network

Artificial Neural Network (ANN) is the mathematical model based on biological neural network. It consists of a set of processing units which communicate together by sending signals to each other over a large number of weighted connections. Such processing units are also called neurons or cells or variables. Each unit is responsible for receiving input from neighbors or external sources and using this input to compute an output signal which is propagated to other units. However each unit also adjusts the weights of connections. There are three types of units:

- Input units receive data from outside the network. These units structure the input layer.
- Hidden units own input and output signals that remain within the neural network. These units structure the hidden layer. There can be one or more hidden layers.
- Output units send data out of the network. These units structure the output layer.

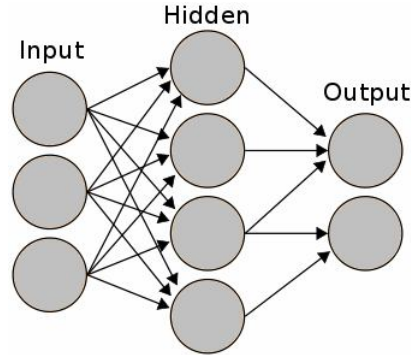


Fig. 6. Neural network

Each connection between unit i and unit j is defined by the weight w_{ij} determining the effect which the signal of unit i on unit j . Suppose input unit, hidden unit and output

unit is denoted as x , h , y respectively. In the topology, unit y is the composition of other units h which in turn are the compositions of others units x . The composition of a unit is represented as a weighted sum which will be evaluated to determine the output of this unit. If such unit is the output unit, its output is the output of neural network. The process of computing the output of a unit includes two following steps:

- An adder so-called summing function sums up all the inputs multiplied by their respective weights. It is essential to compute the weighted sum. This activity is referred to as linear combination.
- An activation function controls the amplitude of the output of the neuron. This activity aims to determine the output.

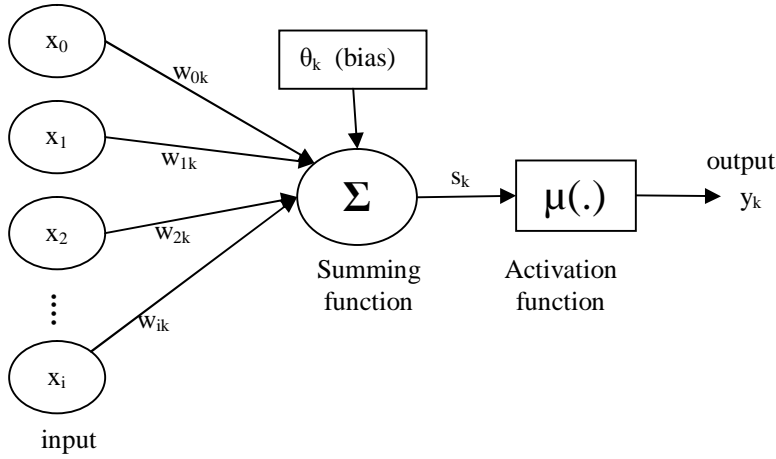


Fig. 7. The process of computing the output of a unit

Note that values of units are arbitrary but they should range from 0 to 1 (sometimes -1 to 1 range). In general, every unit has following aspects:

- A set of inputs connects to it. Each connection is defined by a weight
- Its weighted sum is computed by summing up all the inputs modified by their respective weights
- A bias value is added to weighted sum. This weighted sum is also called activation value.
- Its output is the outcome of activation function on weighted sum. Activation function is crucial factor in neural network.

Activation function is the squashing function which “squashes” a large weighted sum into possible smaller values ranging from 0 to 1 (sometimes -1 to 1 range). There are three types of activation function:

- *Threshold function* takes on value 0 if weighted sum is less than 0 and otherwise.

$$\text{So } \mu(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases}$$

- *Piecewise-linear function* takes on values according to amplification factor in a

$$\mu(x) = \begin{cases} 0 & \text{if } x \leq -\frac{1}{2} \\ x & \text{if } -\frac{1}{2} < x < \frac{1}{2} \\ 1 & \text{if } x \geq \frac{1}{2} \end{cases}$$

- *Sigmoid function* takes on values in range $[0, 1]$ or $[-1, 1]$. The formula of sigmoid function which is $\mu(x) = \frac{1}{1 + e^{-x}}$

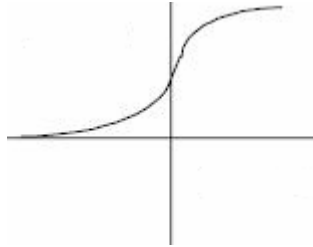


Fig. 8. Sigmoid function

There are two topologies of neural network:

- *Feed-forward neural network*. It is directed acyclic graph in which flow of signal from input units to output units is one-way flow so-called feed-forward. There are no feedback connections
- *Recurrent neural network*. The graph contains cycles; so there are feedback connections in network.

It is necessary to evolve neural network by modifying the weights of connections so that they become more accurate. In other words, such weights should not be fixed by experts. The neural network should be trained by feeding it teaching patterns and letting it change its weights. This is learning process. There are three types of learning methods:

- *Supervised learning*. The network is trained by providing it with input and matching output patterns. These patterns are known as classes.
- *Unsupervised learning*. The output is trained to respond to clusters of pattern within the input. There is no a priori set of categories into which the patterns are to be classified.
- *Reinforcement learning*. The learning machine does some action on the environment and gets a feedback response from the environment. The learning system grades its action rewarding or punishable based on the environmental response and accordingly adjusts weights. Weight adjustment is continued until there is no change in weights.

Reinforcement learning is the intermediate form between supervised learning and unsupervised learning. We apply neural network into classifying corpus and such supervised learning algorithm used in this chapter is back-propagation algorithm.

4.2. Back-propagation algorithm for classification

The back-propagation algorithm [4] is a famous supervised learning algorithm for classification, which is used in feed-forward neural network. It processes iteratively data tuples in training corpus and compares the network's prediction for each tuple to the actual class of the tuple. For each time it feeds a training tuple, the weights are modified in order to minimize the mean squared error between network's prediction and actual class. The modifications are made in backward direction, from output layer through hidden layer down to hidden layer. Back-propagation algorithm includes following steps:

- *Initializing the weights:* the weights are initialized as random real number which should be in space $[0, 1]$. Each bias associated to each unit is also initialized
- *Propagating the input value forward.* Training tuple is fed to input layer. Given the unit j , if unit j is input unit, its input value denoted I_j and its output value denoted O_j are the same.

$$O_j = I_j$$

Otherwise if unit j is hidden unit or output unit, its input value I_j is the weighted sum of all output values of units from previous layer. The bias is also added to this weighted sum

$$I_j = \sum_i w_{ij} O_i + \theta_j$$

Where w_{ij} is the weight of connection from unit i in the previous layer to unit j , O_i is the output value of unit i from the previous layer and θ_j is the bias of unit j .

The output value of hidden unit or output unit O_j is computed by applying activation function to its input value (weighted sum). Suppose activation function is sigmoid function. We have:

$$O_j = \frac{1}{1 + e^{-I_j}}$$

- *Propagating the error backward:* The error is propagated backward by updating the weights and biases to reflect the error of network's prediction. Given unit j , if unit j is output unit then its error is computed as below:

$$Err_j = O_j(1 - O_j)(T_j - O_j)$$

If unit j is hidden unit, the weighted sum of the errors of the units connected to it in the next higher layer is considered when its error is computed. So the error of hidden unit is computed as below:

$$Err_j = O_j (1 - O_j) \sum_k Err_k w_{jk}$$

Where w_{jk} is the weight of the connection from unit j to a unit k in the next higher layer and Err_k is the error of unit k .

- *Updating the weights and biases* is based on the errors. The weights are updated so as to minimize the errors. Given Δw_{ij} is the change in weight w_{ij} , the weight w_{ij} is updated as below:

$$\Delta w_{ij} = (l) Err_j O_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij}$$

Where l is learning rate ranging from 0 to 1. Learning rate helps to avoid getting stuck at a local minimum in decision space and helps to approach to a global minimum.

The bias θ_j of unit j is updated as below:

$$\Delta \theta_j = (l) Err_j$$

$$\theta_j = \theta_j + \Delta \theta_j$$

- *Terminating condition:* There are following terminating conditions:
 - All Δw_{ij} in some iteration are smaller than given threshold
 - Iterating through all possible training tuples.

4.3. Applying neural network into document classification

Given a set of classes $\mathbf{C} = \{\text{computer science, math}\}$, a set of terms $\mathbf{T} = \{\text{computer, programming language, algorithm, derivative}\}$. Suppose all input variables (units) are binary or Boolean, every document is represented as a set of input variables. Each term is mapped to an input variable in which value 1 indicates the existence of this term in document and otherwise value 0 indicates the lack of this term in document. So the input layer consists of four input units: “computer”, “programming language”, “algorithm” and “derivative”.

The hidden layer is constituted of two hidden units: “computer science”, “math”. These units (variables) are also binary or Boolean. The output layer has only one unit named “document class” which is binary or Boolean (0 – documents belong to *computer science* class and 1 – documents belong to *math* class). The evaluation function used in network is sigmoid function. Our topology is feed-forward neural network (showed in figure 4) in which the weights can be initialized arbitrarily.

Note that we denote Boolean value as 0 and 1 (instead of *true* and *false*) for convenience when representing neural network which only accepts numeric value for units.

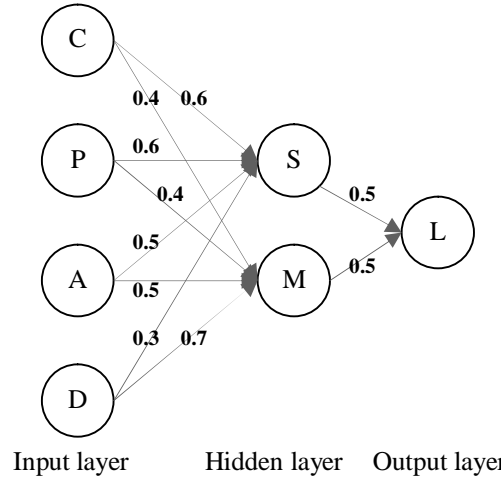


Fig. 9. The neural network for document classification

Note that C , P , A and D denote “computer”, “programming language”, “algorithm” and “derivative” respectively. S and M denote “computer science” and “math” respectively. L denotes “doc class”.

Given corpus $\mathbf{D} = \{doc1.txt, doc2.txt, doc3.txt, doc4.txt, doc5.txt\}$. The training data is showed in following table in which cell (i, j) indicates the number of times that term j (column j) occurs in document i (row i).

Table 7. Term frequencies of documents

	<i>computer</i>	<i>programming language</i>	<i>algorithm</i>	<i>derivative</i>	class
<i>doc1.txt</i>	5	3	1	1	computer
<i>doc2.txt</i>	5	5	40	5	math
<i>doc3.txt</i>	20	5	20	55	math
<i>doc4.txt</i>	20	55	5	20	computer
<i>doc5.txt</i>	15	15	4	0.3	math
<i>doc6.txt</i>	35	10	45	10	computer

Table 8. Normalized term frequencies

	<i>computer</i>	<i>programming language</i>	<i>algorithm</i>	<i>derivative</i>	class
<i>doc1.txt</i>	0.5	0.3	0.1	0.1	computer
<i>doc2.txt</i>	0.05	0.05	0.4	0.5	math
<i>doc3.txt</i>	0.2	0.05	0.2	0.55	math
<i>doc4.txt</i>	0.2	0.55	0.05	0.2	computer
<i>doc5.txt</i>	0.15	0.15	0.4	0.3	math
<i>doc6.txt</i>	0.35	0.1	0.45	0.1	computer

Given threshold $\alpha = 0.4$, if the frequency of a term j in document i is equal or greater than α , we consider that term j exists in document i . Otherwise there is no existence of term j in document i . So each document is represented as a Boolean vector. Each element in such vector has two values: 0 and 1 (0 – the respective term occurs in document and 1 – otherwise). So each Boolean vector is the manifest of the occurrences of terms in a document. Corpus **D** becomes a set of Boolean vectors.

Table 9. Boolean document vectors

	<i>computer</i>	<i>programming language</i>	<i>algorithm</i>	<i>derivative</i>	class
<i>doc1.txt</i>	1	0	0	0	computer
<i>doc2.txt</i>	0	0	1	1	math
<i>doc3.txt</i>	0	0	0	1	math
<i>doc4.txt</i>	0.2	1	0	0	computer
<i>doc5.txt</i>	0.15	0	1	0	math
<i>doc6.txt</i>	0.35	0	1	0	computer

Such vectors are fed to our neural network in figure 3 for supervised learning. Back-propagation algorithm is used to train network; thus Boolean document vectors are considered training tuples. We suppose that the weights in neural network are changed as in figure III.3.9 after training process.

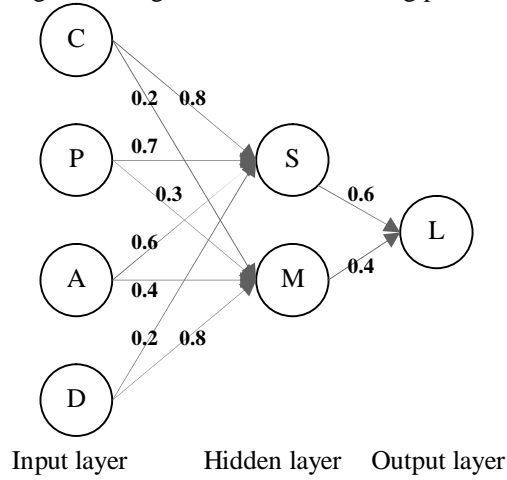


Fig. 10. Trained neural network

5. Discovering user interests based on document classification

Suppose in some library or website, user U do his search for his interesting books, documents... There is demand of discovering his interests so that such library or website can provide adaptive documents to him whenever he visits in the next time. This is adaptation process in which system tailors documents to each individual. Given there is a set of key words or terms $\{computer, programming language, algorithm, derivative\}$ that user U often looking for, the his searching history is showed in following table:

Table 10. User's searching history

<i>Date</i>	<i>Keywords (terms) searched</i>
Aug 28 10:20:01	computer, programming language, algorithm, derivative
Aug 28 13:00:00	computer, programming language, derivative, algorithm
Aug 29 8:15:01	computer
Aug 30 8:15:06	computer

This history is considered as training dataset for mining maximum frequent itemsets. The keywords are now considered items. A itemset is constituted of some items. The support of itemset x is defined as the fraction of total transaction which containing x . Given support threshold min_sup , the itemset x is called *frequent itemset* if its support satisfies the support threshold ($\geq min_sup$). Moreover x is *maximum frequent itemset* if x is frequent itemset and all super-itemsets of x are not frequent. Note that y is super-itemset of x if $x \subset y$. The itemset that has k items is called *k-itemset*. Tabel 9 shows the supports of 1-itemsets

Table 11. 1-itemsets

<i>1-itemset</i>	<i>support</i>
computer	4
programming language	2
algorithm	2
derivative	2

Applying algorithm Apriori, it is easy to find maximum frequent itemsets given $min_sup = 2$. The maximum frequent itemset that user searches are showed in below table:

Table 12. The maximum frequent itemset that user searches

N_o	<i>itemset</i>
1	computer, programming language, algorithm, derivative

We propose the new point of view: “**The maximum frequent itemsets are considered as documents and the classes of such documents are considered as user interests**”. Such documents may be called interesting documents. Which classes

such interesting documents belong to are user interests. It means that discovering user's interests involves in classifying interesting documents. Suppose we have a set of classes $\mathbf{C} = \{\text{computer science, math}\}$, a set of terms $\mathbf{T} = \{\text{computer, programming language, algorithm, derivative}\}$ and the set of classification rules in table 6. Each maximum frequent itemset that user searches is modeled as a document vector (so-called interesting document vector or user interest vector) whose elements are the support of its member items. Note that the supports of such items are showed in table 8.

Table 13. Interesting document vector

N_o	<i>vector</i>
1	(computer=4, programming language=2, algorithm=2, derivative=2)

Table 14. Interesting document vector is normalized

N_o	<i>vector</i>
1	(computer=0.4, programming language=0.2, algorithm=0.2, derivative=0.2)

Table 15. Nominal interesting document vector

N_o	<i>vector</i>
1	(computer= <i>medium</i> , programming language= <i>medium</i> , algorithm= <i>medium</i> , derivative= <i>medium</i>)

It is possible to use SVM or decision tree or neural network to classify documents. Hence we use decision tree as sample classifier for convenience because we intend to re-use classification rules in section III. Otherwise we must determine the weight vector W^* if applying SVM approach. However SVM approach is more powerful than decision tree with regard to document classification in case of huge training data.

Applying classification rule 2, the interesting document belongs to class *compute science* because the frequency of “derivative” and “computer” are *medium* and *medium*, respectively. So we can state that user U has only one interest: *computer science*.

Note that in case of using neural network for document classification, interesting document vector is specified as Boolean document vector (or Boolean user vector). For example, given threshold $\alpha = 0.4$, if the frequency of a term j in document i is equal or greater than α , we consider that term j exists in document i . Otherwise there is no existence of term j in document i . We have a Boolean vector. So user U is modeled as a Boolean document vector. Such vector is also called Boolean user vector: $U = (1, 0, 0, 0)$. According to table III.3.19, we have

Table 16. Boolean document vector (or Boolean user vector)

<i>Term</i>	<i>Existence</i>
computer	1
programming language	0
algorithm	0
derivative	0

The Boolean user vector is considered as a document and the classes of such document are considered as user interests". Now document (user Boolean vector) U becomes a data tuple which is fed to trained neural network in figure III.3.10. It is easy to know the class of document U by checking the value of output unit in trained neural network. Suppose such output value is 0, we can infer that document U belongs to class "computer science". So the interest of user U is "computer science".

6. Evaluation

Our approach includes following steps:

- Documents are represented as vectors
- Classifying documents by using decision tree or support vector machine or neural network
- Mining user's access history to find maximum frequent itemsets. Each itemset is considered an interesting document
- Applying classifiers into interesting documents in order to find their suitable classes. These classes are user interests.

Two new points of view are inferred from these steps:

- The series of user access in his/her history are modeled as documents. So user is referred indirectly to as document.
- User interests are classes that such documents are belong to

The technique of constructing vector model for representing document is not important to this approach. There are some algorithms of text segmentation for specifying all terms in documents. From this, it is easy to build up document vectors by computing term frequency and inverse document frequency. However the concerned techniques of document classification such as SVM, decision tree and neural network influence extremely on this approach. SVM and neural network is more effective than decision tree in case of huge training data set but it is not convenient for applying classifiers (weight vector W^*) into determining the classes of documents. Otherwise it is easy to use classification rules taken out from decision tree for this task.

References

1. N. Cristianini and J. Shawe-Taylor. An Introduction to Support Vector Machines. Cambridge University Press, 2000.
2. T. Mitchell. Machine Learning. McGraw-Hill International, 1997
3. C.Cortes and V.Vapnik. Support vector networks. Machine Learning, 20:273–297, 1995

4. Mohammad Alrifai, Peter Dolog, Wolfgang Nejdl. Learner Profile Management for Collaborating Adaptive eLearning Applications. APS'2006: Joint International Workshop on Adaptivity, Personalization and the Semantic Web at the 17th ACM Hypertext'06 conference, Odense, Denmark, August 2006
5. Christos Papatheodorou. Machine Learning in User Modeling. Lecture Notes in Artificial Intelligence (LNAI), No 2049: Springer-Verlag, 2001, pp. 286-294
6. Rojas, R. Neural Networks: A Systematic Introduction. Springer, Berlin 1996
7. Lippmann, R. P. An introduction to computing with neural nets. IEEE Transactions on Acoustics, Speech, and Signal Processing 1987
8. Christos Papatheodorou. Machine Learning in User Modeling. Lecture Notes in Artificial Intelligence (LNAI), No 2049: Springer-Verlag, 2001, pp. 286-294
9. Jiawei Han and Micheline Kamber. Data Mining: Concepts and Techniques. Second Edition. © 2006 by Elsevier Inc.