# A proposed Mining System for Recognition of Deep Venous Thrombosis and Analyzing the Predisposing Factors

**Assi. Lect. Sarah Saadoon Jasim**
Computer Science/Technical college of

## *Abstract*:

DVT is a very common problem with a very serious complication like pulmonary embolism (PE) which carries a high mortality, and many other chronic and annoying complications ( like chronic DVT, post-phlebitic syndrome, and chronic venous insufficiency) ,and it has many risk factors that affect its course, severity ,and response to treatment. Most of those risk factors are modifiable, and a better understanding of the relationships between them can be beneficial for better assessment for liable patients, prevention of disease, and the effectiveness of our treatment modalities. Male to female ratio was nearly equal, so we didn't discuss the gender among other risk factors. Data taken from 200 patients with DVT were recognized and analyzed by proposed mining system which contains two stages the first one was NN(back-propagation) trained on patients records with DVT, the second stage was to apply the association rule program to extract the relations among the Predisposing Factors. Immobility was the most important risk factor. Alcoholism and Smoking add more risk to immobile of post operative patient. Age per se has no effect. 100% of patients with long bone fracture were immobile. Fever occurred in one third of post operative patients who develop DVT.

The proposed miming system allows better and faster recognition and analysis of more data, which saves time and effort, and discovers the relations among many factors to one or more than one factors. So, we get the above mentioned relations, which are important for the future management of DVT.

## 1. Introduction:

### 1.1 Knowledge Discovery And Data Mining:

With the enormous amount of data stored in files, databases, and other repositories, is increasingly important, if not necessary, to develop powerful means for analysis and perhaps interpretation of such data and for the extraction of interesting knowledge that could help in decision-making.

Data Mining, also popularly known as Knowledge Discovery in Databases (KDD), refers to the nontrivial extraction of implicit, previously unknown and potentially useful information from data in databases. While data mining and knowledge discovery in databases (or KDD) are frequently treated as synonyms, data mining is actually part of the knowledge discovery process.

The Knowledge Discovery in Databases process comprises of a few steps leading from raw data collections to some form of new knowledge. The iterative process consists of the following steps:

- Data cleaning: also known as data cleansing, it is a phase in which noise data and irrelevant data are removed from the collection.
- Data integration: at this stage, multiple data sources, often heterogeneous, may be combined in a common source.
- Data selection: at this step, the data relevant to the analysis is decided on and retrieved from the data collection.
- Data transformation: also known as data consolidation, it is a phase in which the selected data is transformed into forms appropriate for the mining procedure.
- Data mining: it is the crucial step in which clever techniques are applied to extract patterns potentially useful.
- Pattern evaluation: in this step, strictly interesting patterns representing knowledge are identified based on given measures.
- Knowledge representation: is the final phase in which the discovered knowledge is visually represented to the user. This essential step uses visualization techniques to help users understand and interpret the data mining results [1, 2].

## 1.2 Data Mining Methods:

Various data mining methods are [3]:
- Neural Networks
- Genetic Algorithms
- Rough Sets Techniques
- Support Vector Machines
- Cluster Analysis
- Induction
- OLAP
- Data Visualization
- Association rule mining.

## 1.3 Neural Networks:

Neural networks are an approach to computing that involves developing mathematical structures with the ability to learn. The methods are the result of academic investigations to model nervous system learning.

Neural networks have the remarkable ability to derive meaning from complicated or imprecise data and can be used to extract patterns and detect trends that are too complex to be noticed by either humans or other computer techniques.

### *1.3.1 The Perceptron and Multi Layer Perceptrons (MLPs):*

The perceptron is the simplest form of a neural network, since a single layer perceptron operates as a single neuron, multilayer perceptron are an expansion of the perceptron idea and can be used to solve much more difficult problems. They consist of an input layer, one or more hidden layers of an output layer. The hidden layers give the network its power and allow for it to extract extra features from the input. One of the most popular methods used in training a multilayer perceptron is the (error) back-propagation algorithm, which includes two passes through each layer of the network: a forward pass and a backward pass. The back-propagation training algorithm involves three stages [4, 5]:

1. The feed-forward of the input training pattern.
2. The back-propagation of the associated (error).
3. The adjustment of the weights.

### *1.4 Association Rules Mining:*

An Association Rule is a rule, which implies certain association relationships among a set of objects in a database. In this process people discover a set of association rules at multiple levels of abstraction from the relevant set(s) of data in a database. For example, one may discover a set of symptoms often occurring together with certain kinds of diseases and further study the reasons behind them. Since finding interesting association rules in databases may disclose some useful patterns for decision support, selective marketing, financial forecast, medical diagnosis, and many other applications, it has attracted a lot of attention in recent data mining research.

In this section we introduce the standard definition of A-priori Algorithm that is used to discover the association rules,The *support* of an itemset $I[$ sup($I$)], is defined as the number of transactions in the database containing $I$. The *minimum support* ( min_sup), is a user predefined threshold. An itemset is *frequent* if its support is not less than the *min_sup*. An itemset with k items is called a *k*-itemset.

Let $D$ be a set of transactions and $I = \{i_1, i_2, \ldots, i_m\}$, an itemset is a subset of $I$. Given $X$ and $Y$ are itemsets, an *association rule* is of the form $X \Rightarrow Y$, where $X \subset I$, $Y \subset I$, and $X \cap Y = \Phi$, where the $sup(X \cup Y) \geq min\_sup$, and the *confidence* of $X \Rightarrow Y$ is not less than a predefined threshold, *min_conf*, the *confidence* of $X \Rightarrow Y$ is $\dfrac{\sup(X \cup Y)}{\sup(X)}$.

After discovering all frequent item sets, the algorithm of generating association

rules uses the subsets of a frequent item set as antecedents to generate the rules [6, 7, 8].

The form $X \Rightarrow Y$ with confidence 60% where X = computer and Y= software for example means that (60% of the customers who purchased a computer also bought the software).

And the form of $X \Rightarrow Y$ *is* not equal to the form $Y \Rightarrow X$ because

Sup (X U Y)    NOT EQUAL    Sup (Y U X)         unless X=Y which is called perfect rule.

Sup(X)                              Sup(Y)

This must be taken into consideration when analyzing the results of association rules.

Also we used the Mean of the items confidence (MOC) which can be defined as follows:

MOC (A) = {conf (A→ B) + conf (A → C) + conf (A→ D)}/3.

The formula gives the mean of the occurrence of item A among the items B, C and D which means the influence of item A on the others [9].

## 1.5 *Deep Vein Thrombosis:*

Deep vein thrombosis, or DVT, is a blood clot that forms in a vein deep in the body. Blood clots occur when blood thickens and clumps together, [10,11].

Most deep vein blood clots occur in the lower leg or thigh. They also can occur in other parts of the body Figure (1).
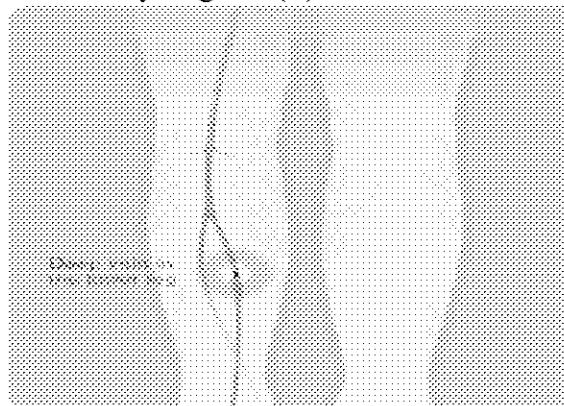


Figure (1) shows deep vein blood clots occur in the lower leg

A blood clot in a deep vein can break off and travel through the bloodstream. The loose clot is called an embolus. It can travel to an artery in the lungs and block blood flow. This condition is called pulmonary embolism, (PE).

PE is a very serious condition. It can damage the lungs and other organs in the body and cause death Figure (2), [10, 11, 12].

Figure (2) shows how a blood clot in a deep vein of the leg can break off, travel to the lungs, and block blood flow.

Blood clots in the thighs are more likely to break off and cause PE than blood clots in the lower legs or other parts of the body. Blood clots also can form in veins closer to the skin's surface. However, these clots won't break off and cause PE, [10,11,12].

### 1.5.1 Causes of Deep Vein Thrombosis:

Blood clots can form in the body deep veins if:

- A vein's inner lining is damaged. Injuries caused by physical, chemical, or biological factors can damage the veins. Such factors include surgery, serious injuries, inflammation, and immune responses.
- Blood flow is sluggish or slow. Lack of motion can cause sluggish or slow blood flow. This may occur after surgery, if the patient were ill and in bed for a long time, or if he were traveling for a long time.
- The blood is thicker or more likely to clot than normal. Some inherited conditions (such as factor V Leiden) increase the risk of blood clotting. Hormone therapy or birth control pills also can increase the risk of clotting.
- The risk factors for deep vein thrombosis (DVT) include:
- A history of DVT.
- Conditions or factors that make the blood thicker or more likely to clot than normal. Some inherited blood disorders (such as factor V Leiden) will do this. Hormone therapy or birth control pills also increase the risk of clotting.
- Injury to a deep vein from surgery, a broken bone, or other trauma.
- Slow blood flow in a deep vein due to lack of movement. This may occur after surgery, if the patient were ill and in bed for a long time, or if he were traveling for a long time.
- Pregnancy and the first 6 weeks after giving birth.
- Recent or ongoing treatment for cancer.
- A central venous catheter. This is a tube placed in a vein to allow easy access to the bloodstream for medical treatment.

- Older age. Being older than 60 is a risk factor for DVT, although DVT can occur at any age.
- Overweight or obesity.
- Smoking.
- The risk for DVT increases if the patient has more than one of the risk factors listed above, [10, 11, 12].

### *1.5.2 Signs and Symptoms of Deep Vein Thrombosis:*

Only about half of the people who have DVT have signs and symptoms. These signs and symptoms occur in the leg affected by the deep vein clot. They include:

- Swelling of the leg or along a vein in the leg
- Pain or tenderness in the leg, when standing or walking
- Increased warmth in the area of the leg that's swollen or painful
- Red or discolored skin on the leg
- Some people aren't aware of a deep vein clot until they have signs and symptoms of PE. Signs and symptoms of PE include:
- Unexplained shortness of breath
- Pain with deep breathing
- Coughing up blood
- Rapid breathing and a fast heart rate also may be signs of PE, [10, 11, 12].

The most common test for diagnosing deep vein blood clots is ultrasound. This test uses sound waves to create pictures of blood flowing through the arteries and veins in the affected leg Figure (3).



Figure (3) shows blood clots in the leg and arms (ultrasound image)

## *2. Methods:*

### *2.1 Study Design:*

A prospective observational study was conducted from August the 1$^{st}$ 2011, to Feb the 1$^{st}$ 2013.

### *2.2 Study sample:*

The sample was collected from outpatient clinics in two general hospitals, and a private vascular clinic at Baghdad, 57 patients (Binary records) who were diagnosed to have DVT in previous research [9] were included in this study .Data that was collected include the name, age, gender, history of surgical intervention, history of long bone fracture, history of immobility history of

febrile illness before the DVT, smoking habits, alcohol intake, family history of DVT, dopplex scan results.

Seven predisposing factors (age<=40, long bone fractures, immobility, febrile illness, history of surgery before DVT, smoking, alcoholism) were chosen to be applied by the proposed system. The total number of patients studied was 200. We exclude pregnancy and child birth factor to avoid bias between the two sexes. We took the age of<=40 years because we have lower mean of age than Western countries, and we have more young people affected because of less firm preventive measures taken preoperatively. 57 records were analyzed by the association rule and mean of confidence by previous research [9] before programming the proposed system and then added 143 new not recognized records to test the proposed system and to make a comparison between the old and the new results.

## 2.3 *The Proposed system:*

The proposed system contains two stages as shown in the following diagram:



Figure (4) shows the Proposed system diagram

*2.3.1.* *Neural networks stage:* This neural network(back-propagation) structure is:

- Input layer consists of 9 neurons.
- Hidden layer consists of 7 neurons.
- Output layer consists of 3 neurons.

The number of neurons of any layer can affect the number of connections to
the related layer. As the number of connection increases, this can be time and
cost consuming, because the process of each connection needs one
multiplication process (wij * input i) the weight Matrix (input layer) = 9*7 = 63
connection (full connection). If the number of hidden layer increase (for
example = 8) then the weights matrix w= 9*8 = 72 (full connection). For
example, if this process takes five millisecond then for 72 connections needs
(5*72) mill second (this for time) and the same criteria can be applied to space.
This base can be applied for all layers then as result figure (5).



Figure (5) MLP proposed for the proposed system

### 2.3.2. Association Rules Mining stage:

The proposed system applied the Interactive KDD system for fast mining
association rules and mean of confidence [7, 8]. All the above techniques have
been collected in the proposed system.

### 2.3.3. The implementation:

Some chosen patterns by the doctors have been applied to train the neural
network as shown in Figure (6).



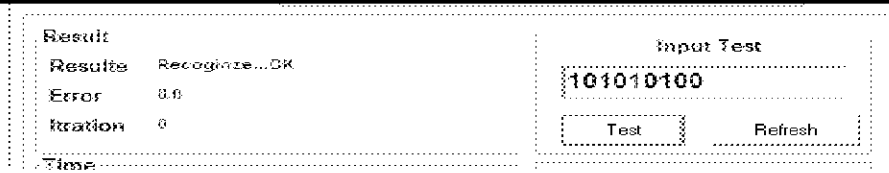Figure (6) Training the neural network

Figure (7) shows the recognition of patterns

All the passed patterns from the first stage (NN) are gathered in the input file of the second stage (AR) to be analyzed as shown in Figure (8).
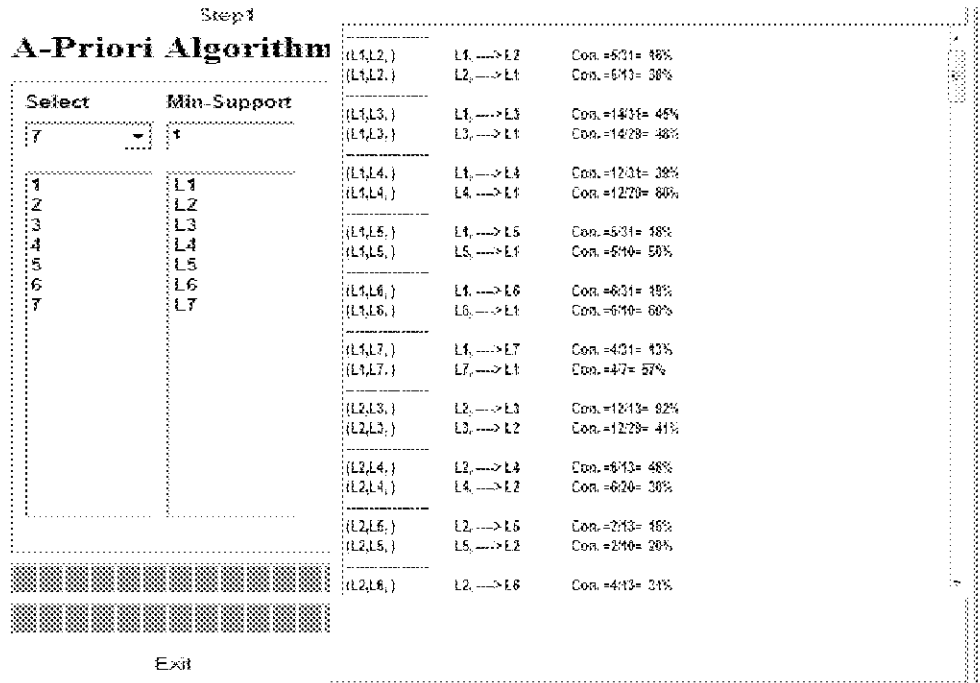


Figure (8) some of the extracted relations to be analyzed

## 3. Results:

**3.1** By analyzing our output relations regarding risk factors of vein thrombosis by the proposed system, we got so many relations among those factors, from which we can extract the significant relations that could have an effect on the pathogenesis, severity, management and the course of the disease.

### 3.2 Association rules results:

Where L1 = AGE <=40 years.

L2 = long bone fractures.

L3 = immobility.

L4 = surgery before DVT.

L5 = Febrile illness.

L6 = Smoking.

L7 = Alcoholism.

From those so many relations that have been extracted and those who extracted by the previous research [9], doctors chosen the most significant relations (rules) with high confidence to be compared Figure (9).

1. All patients (100%) with long bone fractures were immobile, where the proposed system shows the same (100%) (L2 ⇒ L3 conf 100%).

2. All alcoholic patients (100%) were immobile, where the proposed system shows the same (100%) (L7 ⇒ L3 conf 100%).

3. All smoker patients <=40 years, with surgery before DVT were immobile, where the proposed system shows the same (100%) (L1 L4 L6 ⇒L3 conf 100%).

4. 70% of patients with surgery before DVT were immobile, while (50%) of immobile patients had surgery before DVT, where the proposed system shows (69%) and (55%) respectively (L4 ⇒ L3 conf 69%, L3 ⇒ L4 conf 55%).

5. 60 % of smoker patients were <=40 years, where the proposed system shows (59%) (L6 ⇒ L1 conf 59%).

6. 75% of alcoholics with surgery before DVT, had a long bone fracture, where the proposed system shows (67%) (L4 L7 ⇒ L2 conf 67%).

7. 75% of patients<=40 years, with long bone fractures, were smokers, where the proposed system shows (71%) (L1 L2 ⇒ L6 conf 71%).



Figure (9) Some of the chosen relations

## 3.3 Mean of confidence and the effects of factors:

Mean of confidence of the previous research [9] (occurrence of the other factors with each single factor), Figure (10).

Age 22%, long bone fracture 42%, immobility 32%, surgery before DVT 32%, febrile illness 28%, smoking 43% and alcoholism 58%.
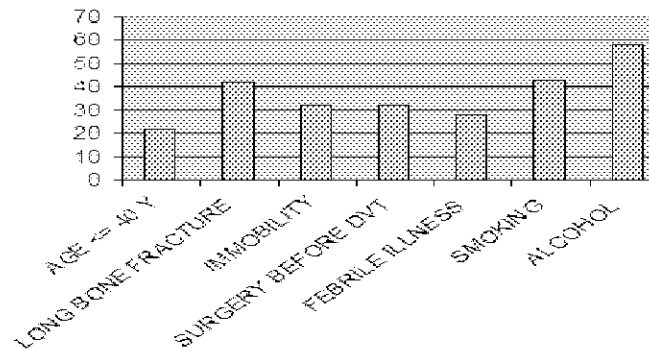
Figure (10) shows the Mean of confidence of previous research [9]

Mean of confidence of the proposed system (occurrence of the other factors with each single factor), Figure (11).

Age 19%, long bone fracture 39%, immobility 28%, surgery before DVT 27%, febrile illness 25%, smoking 39% and alcoholism 50%.
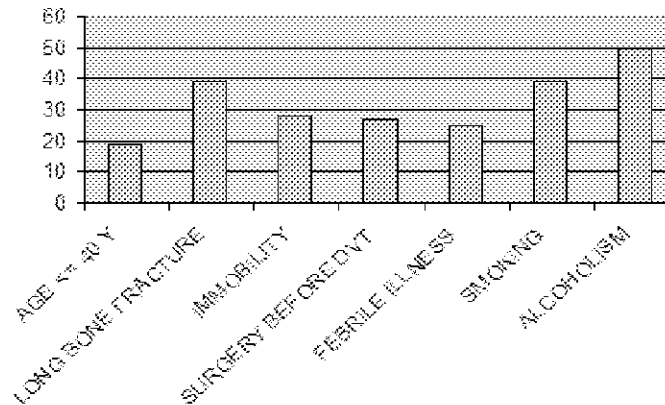


Figure (11) shows the Mean of confidence of the proposed system

## 4. Discussion:

DVT is a very common problem with a very serious complication(PE)which carries a high mortality, and many other chronic and annoying complications ( like chronic DVT, post-phlebitic syndrome, and chronic venous insufficiency) ,and it has many risk factors that affect its course, severity ,and response to treatment. Most of those risk factors are modifiable, and a better understanding of the relationships between them can be beneficial for better assessment, control of disease progress, and the effectiveness of our treatment modalities. Male to female ratio was nearly equal , so we didn't discuss the gender among other risk factors.

The proposed system allows better and faster analysis of more data with an interactive powerful system, which saves time and effort, and diagnoses DVT and discovers the relations among many factors to one or more than one factors. So, we get the above mentioned relations, which are important for the future management of DVT.

## 5. Conclusions and Recommendations:

1) The system gives results which are so close to the results of the previous research [9] after adding 143 recognized records which mean that the system goes with the standards of the DVT researches.

2) Alcoholism factor of the previous research [9] and the proposed system was the same; this shows the high influence of alcohol among DVT factors.

3) Immobility is a very important risk factor for DVT after surgery, so, early mobilization of patients after surgery will reduce the incidence of DVT.

4) Smoking, add more risk for DVT, when combined with immobility after surgery. So, cessation of smoking preoperatively will reduce the risk of DVT.

5) Occurrence of age with other factors gives no important relations because DVT could occur at any age.

6) The MOC (mean of confidence) gives almost the same effectiveness of the alcoholism, smoking and long bone fracture among the other factors, Figure (10, 11).

7) Occurrence of long bone fracture with immobility was 100%, and this is quite correct Because all patients with long bone fracture could not move, and they either have a cast external or internal fixator, which prevent them from movement for few weeks at least, While its occurrence with other factors wasn't so significant. This indicates that MOC is a sensitive indicator and goes with clinical applications.

8) About one third of patients with surgery before DVT had history of fever. This means that prevention and control of post operative fever will reduce the incidence of DVT after surgery.

9) Genetic factor must be taken in consideration with the other factors because the system rejects the full zeroes pattern (patient record), which that means there was patient with no factors mentioned, has DVT.

10) We can enlarge the neural network and expand pattern by adding new factors, such as contraceptive pills, genetic factor, gender, and pregnancy.

11) Applying negative association rule mining for low confidence patterns.

## *References*

[1] M. S. Chen, J. Han, and P. S. Yu. Data mining: An overview from a database perspective. IEEE Trans. Knowledge and Data Engineering, 8:866-883, 1996.

[2] U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy. Advances in Knowledge Discovery and Data Mining. AAAI/MIT Press, 1996.

[3] H Lu, R Setiono, H Liu. Effective Data Mining Using Neural Network. IEEE
Transactions on Knowledge and Data Engineering, 1996, 8(6):957-961.

[4] L. Fasett, "Fundamentals of neural networks" Prentice Hall International
Inc., 1994.

[5] C. gershenson "Artificial Neural networks for beginners",
Gershenson@sussex.ac.uk.
1993.

[6] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules", In
Proc. of Int. Conf. on Very Large Data Bases (VLDB'94), Santiago, Chile,
September 1994, pp. 487-499.

[7] G. Johnny, Developing A-priori Algorithm for Fast Mining Association Rules
, Al-Taqani, Refereed Scientific Journal, Foundation of Technical Education,
Vol. 22 , No.5, 2009, Baghdad-Iraq.

[8] G. Johnny, Interactive KDD system for fast mining association rules, Al-
Taqani, Refereed Scientific Journal, Foundation of Technical Education, Vol.
22 , No.5, 2009, Baghdad-Iraq.

[9] G. Johnny and Dr M. ALassel "Association Rules Mining Analysis of the
Deep Vein Thrombosis Predisposing Factors" 2nd conference, Technical
College of
Management – Baghdad 28-29 /11/2012.

[10] Schwartz's Principles of surgery f.Charles Brunicardi ,8th
Edition,chap.24,2004.

[11] Sabiston DC:Disorder of the arterial system.In Textbook of surgery edited
by David C. Sabiston .14 ed ,vol.1 philadelphia W,.B.Saunders
Company,1991 1:1618-1722.

[12] Galanaud JP, Bosson JL, Quéré I (2011). "Risk factors and early outcomes
of patients with symptomatic distal vs. proximal deep-vein thrombosis". Curr
Opin Pulm Med 17 (5): 387–
91. doi:10.1097/MCP.0b013e328349a9e3.PMID 21832920.

# نظام تنقيبي مقترح لتمييز مرض تخثر الاوردة العميقة وتحليل العوامل المسببة

## الخلاصة:

أن تخثر الأوردة العميقة هي مشكلة صحية شائعة ولها مضاعفات خطيرة مثل الأنصمام الرئوي (خثرة الوريد الرئوي)، والذي يحمل معدلات وفيات مرتفعة للغاية، ولها أيضاً مضاعفات أخرى عديدة مزعجة (مثل تخثرالأوردة المزمن، متلازمة ما بعد التخثر، وعدم كفاءة الأوردة المزمن)، وهناك عدة عوامل خطورة لهذا المرض وهذه تؤثر على مسيرة المرض وشدته واستجابته للعلاج. وحيث أن معظم هذه العوامل قابلة للتعديل، فان فهماً أفضل لعوامل الخطورة هذه والعلاقات فيما بينها سيؤدي حتماً الى تقييم أفضل لمعرفة الأشخاص المعرضين للأصابة، والوقاية من المرض، وكفاءة أساليب العلاج. كانت نسبة الرجال للنساء متساوية تقريباً، لذا لم نتناول جنس المرضى ضمن عوامل الخطورة.

تم تمييز وتحليل المعطيات المأخوذة من 200 مريضاً مصاباً بتخثر الأوردة العميقة بنظام تنقيبي مقترح مكون من مرحلتين الاولى شبكة عصبيه مدربه على قيود مرضى مصابين بالمرض أما المرحله الثانيه فهي تطبيق نظام القواعد العلائقيه لايجاد العلاقات بين العوامل المسببه للمرض. كان عامل عدم الحركة الأكثر تأثيراً، وتبين ان الكحول والتدخين يضيف عامل نسبة خطورة أكثرلمرضى ما بعد العمليات الجراحية من عديمي الحركة.لم يتبين أي تأثير للعمر على حدوث التخثر. فيما كان 100% من المرضى المصابين بكسور في العظام الطويلة هم عديمي الحركة. وجدت الحمى في ثلث المرضى ما بعد العمليات الجراحية ممن أصيبوا بالتخثر في الأوردة العميقة.

يتيح النظام التنقيبي المقترح تمييزا وتحليلاً أسرع وأفضل لبيانات أكثر،و يوفر الوقت والجهد ويكتشف علاقات عديدة بين عدةعوامل وعامل واحد أو أكثر. أستخدمنا هذا النظام لتمييز وتحليل البيانات وحصلنا على العلاقات التي ذكرناها أعلاه، والتي لها أهمية في معالجة مرض تخثر الأوردة العميقة مستقبلاً.

في هذا البحث تم اعتماد المنهجيه العالميه لمثل هكذا بحوث كما هو متبع في بحوث المكتبه الافتراضيه وتم عرض النتائج النظريه على أطباء ذوي اختصاص مما أغنى الاستنتاجات بملاحظاتهم الطبيه.

*الكلمات المفتاحية:* – القواعد العلائقية، تخثر الاوردة العميقة، عدم الحركة، الحمى، كسور العظــام، الجراحــة قبــل التخثر، التدخين، الكحول، الشبكات العصبيه، الخوارزميات الجينية.