

# The characteristics of hepatitis B virus sequence with virology and clinical pathology

Yi-Chun Lin<sup>a</sup> Koun-Tem Sun<sup>b</sup> Wen-Chun Lin<sup>c</sup> Ting-Tsung Chang<sup>d</sup> Yueh-Min Huang<sup>a</sup>

<sup>a</sup> Department of Engineering Science <sup>c</sup>Institute of Molecular Medicine <sup>d</sup>Institute of Basic Medical Sciences,

<sup>acd</sup>National Cheng Kung University ; <sup>b</sup>Department of Information and Learning Technology, Nation University of Tainan, Tainan, Taiwan

ktsun@mail.nutn.edu.tw

## Abstract

The hepatitis B virus infection is the invisible killer of pathological changes of liver, and endangered people's health for a long time. Recently, related studies have found that difficult genotypes and HBsAg seroconversion effect the pathological change of liver. In this paper, we explored the variation of hepatitis B virus (HBV) on different genotypes and the HBsAg seroconversion in patients. Then, we develop a technique for association rules to solve this problem.

## 1. Samples and data preprocessing

To investigate the HBV DNA sequences, we collected analyzing samples form routine diagnostic blood specimens amount to twenty-two patients and obtained the amino acid sequences amount to three hundred and fifty-eight from NCBI database.

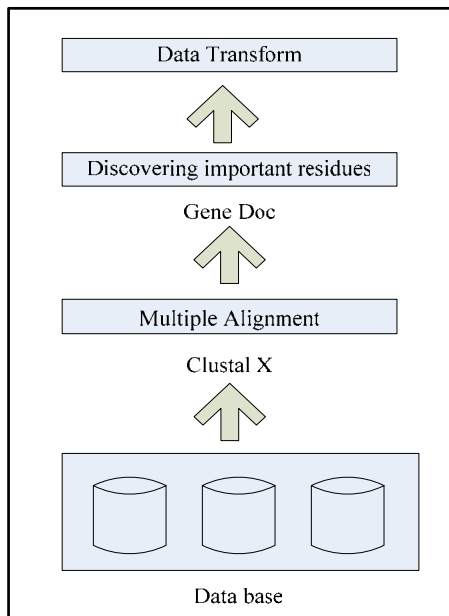


Figure 1. The flow chart of data preprocessing

## 2. Mining specific positions of HBV sequence

Of particular interest in study is employed association analysis to find the specific positions of HBV sequence and which have relation. Then, we further prune the discovered associations to remove those insignificant

associations and find a set of useful rules.

## 3.1 Association rules of distinguish genotypes

We discovered that the rules related to genotyping appear in p gene, S gene, PreS1 gene, respectively.

Table 1. The specific positions in p gene and s gene

Number	Gene	Association rules
rule1	P	{position_104=N,position_234=G,position_584=H,position_199=P}
rule3	S	{position_8=L,position_47=V,position_49=L,position_56=Q,position_57=I,position_59=S,position_64=C,position_77=-,position_85=C}
rule7	PreS1	{position_10=K,position_35=K,position_39=E,position_45=L,position_48=H,position_51=N,position_54=D,position_57=K}

Experimentally, the genotype A is classified according to that when position\_104 shows “K” of P gene, position\_47 shows “L” and position\_59 shows N of S gene at the same time. Therefore, we can corollary the following rules.

- {position\_104=K, position\_49=L, position\_59=N} => Genotype A.
- {position\_199=P, position\_47=V, position\_59=S} => Genotype B.
- {position\_104 =R, position\_199=Q } => Genotype C.
- {position\_199 =Q, position\_47=V, position\_49=L} => Genotype D.
- {position\_104 =K, position\_199=H, position\_584=N} => Genotype E.
- {position\_104 =L, position\_234=N} => Genotype F.
- {position\_104=K,position\_199=Q,position\_234=R, position\_47=V,position\_49=P } => Genotype G.
- {position\_199=A, position\_234=N, position\_584=A} => Genotype H.

