## **Coronary Artery Disease**



Angina Pectoris Myocardial infraction

## Sudden death of coronary artery disease

- The death rates that of from diseases in the a recent medical report are much higher than that of from accidents and disasters
  - Statistical reports show that 130(cancer), 77(cerebrovascular), and 37(cardiovascular) per 100,000 people in the Korean population have been dying as the result of the diseases every year
  - Increased carotid intima-media thickness (IMT) is associated with atherosclerosis risk factors and adverse cardiovascular outcomes

## Our objective

- The first aim is to implement Korean reference standard database
  - Normal and abnormal reference and standard data such as multiparametric features, clinical information
  - For different heart diseases, we develop diagnostic indexes and algorithms

#### The second aim is to develop biosignal data mining methodology for predicting and diagnosing the cardiovascular disease

- A multi-parametric measure of HRV, ST-segments and IMT measurements and
- A suitable prediction method to enhance the reliability of medical treatment
- We apply and analyze existing classification and statistical methods

# Development of intelligent automatic diagnosis system prototype (KOSEF project)



intelligent automatic diagnosis system architecture

- iDB 2008 -

A Data Mining Approach and Framework of Intelligent Diagnosis System for Coronary Artery Disease Prediction

Keun Ho Ryu<sup>1</sup>, Heon Gyu Lee<sup>1</sup>, Wuon-Shick Kim<sup>2</sup> <sup>1</sup>(khryu, hglee@dblab.chungbuk.ac.kr), <sup>2</sup>wskim@kriss.re.kr <sup>1</sup>Database/Bioinformatics Laboratory, Chungbuk National University, Korea <sup>2</sup>Korea Research Institute of Standards and Science, Korea

## Outline

#### Introduction

- <u>Application of data mining</u> <u>to medical data</u>
- Motivation and objective
- Background
  - ECG Analysis
  - ST-segments
  - Heart Rate Variability
  - Carotid arterial wall thickness
- Automated detection of Ischemic ECG beats
  - ST-segments feature extraction
  - CASE study1: Ischemia beat classification

- Heart Rate Variability
  - ECG processing, Linear/ Nonlinear features
  - Classification methods for diagnosing CAD
  - CASE study2 : AP and ACS diseases classification
- Dyslipidemia diagnosis methodology
  - Measuring Intima-Media thickness
  - Feature extraction from IMT image
  - CASE study3: dyslipidemia classification
- Discussion
  - On going work
    - Compressed Patricia FP-Tree for Frequent Itemsets Mining
- Conclusion

## Introduction ... : Why data mining?

Modern medicine generates, almost daily, huge amounts of heterogeneous data

- Medical data may contain signals like ECG, clinical information like temperature, cholesterol levels, etc., as well as the physician's interpretation
- Those who deal with such data understand that there is a widening gap between data collection and data comprehension
- As more and more medical procedures employ signal and imaging as preferred diagnostic tools

# There is a need to develop methods for efficient mining in databases of signals

## Introduction ... Application of data mining to medical data

- Cardiovascular disease
  - Myocardial ischemia : ST-segments from ECG signal
  - Coronary artery disease: Multi-parametric features of Heart rate variability from ECG signal



Myocardial ischemia

Angina Pectoris

Carotid wall thickness (intima-media thickness)



#### A Data Mining Framework of Intelligent Heart Disease Diagnosis System









#### Motivation : Why need multi-parametric features

- Coronary heart disease is in danger of becoming the most fatal disease among Koreans because of <u>smoking</u> and the adoption of <u>westernized eating habits</u>
- It is very important social issue to detect early heart problems and predict accurate heart disease
- The various diagnosis indexes have different physiological meanings
  - There is not a master feature that explains the whole characteristics of diagnosis indexes at one time
  - Researchers must select to use among various features
  - Need multi-parametric feature with accuracy
    - It will be helpful to predict the heart disease

## **Objective**

A helpful diagnostic supplementary tool, using various features, that takes into consideration all possible diagnostic features

Linear and Nonlinear features of HRV, ST-segments, IMT

## Outline

#### Introduction

- Application of data mining to medical data
- Motivation and objective

#### Background

- ECG Analysis
- ST-segments
- Heart Rate Variability
- Carotid arterial wall thickness
- Automated detection of Ischemic ECG beats
  - ST-segments feature extraction
  - CASE study1: Ischemia beat classification

- Heart Rate Variability
  - ECG processing, Linear/ Nonlinear features
  - Classification methods for diagnosing CAD
  - CASE study2 : AP and ACS diseases classification
- Dyslipidemia diagnosis methodology
  - Measuring Intima-Media thickness
  - Feature extraction from IMT image
  - CASE study3: dyslipidemia classification

#### Discussion

- On going work
  - Development of Contrast data mining and frequent pattern mining methods

#### Conclusion

## Background : BioSignals (ECG:Electrocardiogram)

#### The most used signal in clinical practice is the Electrocardiogram (ECG)

- ECG is test that measures a heart's electrical activity
- ECG provides clinical information about the heart's status



## Background ... : A basic approach to ECG analysis

Simply a guide to understanding how clinicians identify abnormalities in the ECG

- 1. Identify the QRS complexes
- 2. Identify the P waves
- 3. Examine ST-T segments: are there abnormalities (such as elevation or depression)?
- 4. Examine the T wave
- 5. Examine the QT interval



ECG analysis is a routine part of any complete medical evaluation, due to the heart's essential role in human health and disease

#### Background ... : ST-segments and T amplitude

Several techniques the evaluate the ST-segment changes and the Twave alterations by different methodologies have been used to ischemic beat detection

- Parametric modeling, wavelet theory, set of rules, neural nets, decision analysis, genetic algorithms
- It is clinically important if elevated or depressed as it can be a sign of ischemia or infarction



http://dblab.chungbuk.ac.kr

#### Background ... : Heart Rate Variability

- HRV (Heart Rate Variability) has been used to assess autonomic control of the heart under physiological and pathological conditions
  - Control of Heart Rate is known to be affected by the sympathetic and parasympathetic system
- HRV is used as a clinical tool predicting heart function in both health and disease
- ARV is based on analysis of RR intervals in ECG signals
  - RRIs (RR intervals) are the series of time intervals between heartbeats



#### Background ... : Heart Rate Variability

HRV is one of the most promising indications of autonomic activity (sympathetic / parasympathetic activity)

- Reduced parasympathetic activity has been reported in patients with CAD (Coronary Artery Disease)
- This reduction in parasympathetic activity evaluated by HRV analysis

#### Two properties of HRV analysis

- Linear analysis : <u>time- and frequency-domain</u>
- Nonlinear analysis : <u>complexity and irregularity</u>

## Background ... : Heart Rate Variability

#### Linear property : Time- / Frequency- domain

- Linear time domain have been employed in quantification of the overall variability of HR
- Linear frequency domain variables provide markers of the heart regulation

#### Nonlinear property

- The heart rate has complexity which reflect the healthy condition in a living body
- The complexity of the human physiological system is reduced in bad health, but increased in good health
- This complexity can be analyzed by various nonlinear methods

#### Background ... : Effect of 3 recumbent positions on nervous system using HRV

There is a relationship between recumbent position and HRV

- In linear analysis, parasympathetic activity is increased and sympathetic activity is decreased in the right position
- In nonlinear analysis, complexity are increased in right position
- Among 3 recumbent positions (supine, right, left), the right position can enhance the parasympathetic activity
- Right recumbent position is recommended in patient to recover the depressed parasympathetic activity
- Because of the effect of the recumbent position on HRV, we considered that it is worthwhile to include Linear/Nonlinear properties of HRV for three positions

#### Background ... : Carotid arterial wall thickness

- The appearance of atherosclerosis has been highly associated with incidence of coronary heart disease
- the severity of carotid IMT is an independent predictor of coronary events such as myocardial infarction [American Heart Association]
  - Carotid IMT consists of intima thickness and media thickness



#### Background : Classification for Diagnosing Disease (SVM)

#### SVM (Support Vector Machine)

- A new classification method for both linear and nonlinear data
- It uses a nonlinear mapping to transform the original training data into a higher dimension



## Background ... : Classification for Diagnosing Disease (SVM)

- With the new dimension, it searches for the linear optimal separating *hyperplane* (i.e., "decision boundary")
- SVM finds this hyperplane using *support vectors* ("essential" training tuples) and *margins* (defined by the support vectors)



## Background : Classification for Diagnosing Disease (MDA)

- Multiple Discriminant Analysis is an analysis of dependence method that is a special case of canonical correlation
- With more than two groups, there will potentially be more than one discriminant function that can be used to explain the differences among groups
  - For example, if we want to discriminate among three groups, two canonical discriminant functions will be derived
  - The first discriminant function separates group 1 from groups 2 and 3, and the second discriminant function separates group 2 from group 3.
- We can obtain classification function for prediction through the discriminant analysis
- Classification function is generated for each group
  - If new case we have to classify comes into existence, this subject will belong to a group that has the highest value of classification function

#### **Background** ....

## : Classification for Diagnosing Disease (CMAR)

#### Associative classification (CMAR: Classification based on

|  | multiple assoc  | 🔮 CARs  |   |  |                         |      |  |  |
|--|---|---|---|--|-------------------------|------|--|--|
|  | Association rule                                      | File Open   | Support   | confidence   | List Rules              | tion |  |  |
|  | <ul> <li>Search for str<br/>(conjunctions)</li> </ul> | (2) {1.092 < APEN(R) <= 1.13 0.358 < MT <= 0.446} -> {CLASS(I)=CAD} 100.0 (3) {0.9 < H_EXT(S) <= 1.009 0.623 < IMT <= 0.712} -> {CLASS(I)=CAD} 100.0 (4) {917.773 < M_NN(L) <= 976.974} -> {CLASS(I)=CAD} 100.0 (17) (0.9 < H_EXT(S) <= 1.009 0.358 < MT <= 0.446} -> (CLASS(I)=CAD) 20.0 |   |  |                         |      |  |  |
|  | <ul> <li>Classification</li> </ul>                    | (19) {2.202 < sampen1(S) <= 2.251} -> {CLASS(I)=Normal} 88.88<br>(21) (720 427 < M NN(I) <= 709 272) -> (CLASS(I)=Normal) 87.5  |   |  |                         | 21   |  |  |
|  | <u>p₂ ^ p₁ →</u>                                      | (22) {2.157 < sampen1(3)<br>(22) {2.157 < sampen1(3)<br>(23) {18.62 < SDNN(L) <<br>(24) {0.322 < NLE(R) <=  | 2) {2.157 < sampen1(S) <= 2.202 0.358 < MT <= 0.446} -> {CLASS(I)=CAD} 87.5<br>3) {18.62 < SDNN(L) <= 22.976} -> {CLASS(I)=CAD} 87.5<br>4) {0.322 < NLE(R) <= 0.374} -> {CLASS(I)=CAD} 87.5 |  |                         |      |  |  |
| It may overco         (27) {0.663 < Dfe2(S) <= 0.836} -> {CLASS(I)=Normal} 84.61         (27) {0.663 < Dfe2(S) <= 0.836} -> {CLASS(I)=Normal} 84.61         (27) {0.663 < Dfe2(S) <= 0.836} -> {CLASS(I)=Normal} 84.61         (27) {0.663 < Dfe2(S) <= 0.836} -> {CLASS(I)=Normal} 84.61         (27) {0.663 < Dfe2(S) <= 0.836} -> {CLASS(I)=Normal} 84.61         (27) {0.663 < Dfe2(S) <= 0.836} -> {CLASS(I)=Normal} 84.61         (27) {0.663 < Dfe2(S) <= 0.836} -> {CLASS(I)=Normal} 84.61         (27) {0.663 < Dfe2(S) <= 0.836}          (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0.836}         (27) {0.663 < Dfe2(S) <= 0 |   |   |   |  |                         | h-   |  |  |
|  | tree induction  | (29) {0.756 < NLP(S)} -><br>(31) {799.372 < M_NN(L<br>(32) {1.018 < Dfe1(R) <=  |   | ne   |                         |      |  |  |
| •  | CMAR classific  | (33) <u>{0.535 &lt; NLF(S) &lt;=</u><br>(34) {0.535 < NLF(S) <=<br>(35) {0.535 < NLF(S) <=<br>(37) {0.275 < F1_AS <= 1  | 0.634 1.022 < LF_HF(S) <=<br>0.634} -> {1.022 < LF_HF<br>0.634} -> {CLASS(I)=CAD<br>1.083} -> {CLASS(I)=Norm  | <u>= 1.698} -&gt; {CLASS(I)=CAD</u><br>(S) <= 1.698 CLASS(I)=CAD<br>} 80.0<br>nal} 77.77 | <u>} 80.0</u><br>} 80.0 |      |  |  |
|  |   | (38) {0.329 < IT <= 0.339<br>(40) {1.026 < APEN(S) <=<br>(44) {1.383 < Dfe1(L) <=<br>(45) {1.092 < APEN(R) <<br>(55) {1.092 < APEN(R) <   | I} -> {CLASS(I)=CAD} 77.1<br>= 1.06 0.9 < H_EXT(S) <= 1<br>1.906} -> {CLASS(I)=Norr<br>= 1.13} -> {CLASS(I)=CAD<br>= 1.13 0.358 < MT <= 0.448   | 77<br>1.009} -> {CLASS(I)=CAD}<br>mal} 76.92<br>)} 76.92<br>}} -> {CLASS(I)=CAD} 100.    | 0                       |      |  |  |

#### Background ....

## : Classification for Diagnosing Disease (Bayes.)

- Let *D* be a training set of tuples and their associated class labels, and each tuple is represented by an *n*-*D* attribute vector  $X = (x_1, x_2, ..., x_n)$
- Suppose there are m classes  $C_{1}, C_{2}, ..., C_{m}$ .
- Classification is to derive the maximum posteriori, i.e., the maximal  $P(C_i/X)$
- This can be derived from Bayes' theorem

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i) P(C_i)}{P(\mathbf{X})}$$

A simplified assumption: attributes are conditionally independent -> Naïve Bayesian

$$P(\mathbf{X} | C_i) = \prod_{k=1}^{n} P(x_k | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$

#### Background ... Classification for Diagnosing Disease (Bayes.)

- Bayesian belief network allows a subset of the variables conditionally independent
  - A graphical model of causal relationships
  - Represents dependency among the variables
  - Gives a specification of joint probability distribution
    - Nodes: random variables
    - Links: dependency
  - Network search algorithm
    - Tree Augmented Naïve Bayes (TAN)





#### 30-Sep-08

#### http://dblab.chungbuk.ac.kr

## Outline

#### Introduction

- Application of data mining to medical data
- Motivation and objective
- Background
  - ECG Analysis
  - ST-segments
  - Heart Rate Variability
  - Carotid arterial wall thickness

#### Automated detection of Ischemic ECG beats

- ST-segments feature extraction
- <u>CASE study1: Ischemia beat</u> <u>classification</u>

- Heart Rate Variability
  - ECG processing, Linear/ Nonlinear features
  - Classification methods for diagnosing CAD
  - CASE study2 : AP and ACS diseases classification
- Dyslipidemia diagnosis methodology
  - Measuring Intima-Media thickness
  - Feature extraction from IMT image
  - CASE study3: dyslipidemia classification
- Discussion
  - On going work
    - Development of Contrast data mining and frequent pattern mining methods
    - Conclusion

## Automated detection of Ischemic beats

Flowchart of the overall ischemic beat classification methodology





- Location of the J point was detected
  - If (heart rate  $\leq$  120 beats/min) then J80
  - If (heart rate > 120 beats/min) then J60
- ST-segment deviation: refer to the amplitude deviation of the STsegment from the isoelectric line (the line defining the level of zero amplitude)
- ST-segment slope: the slope of the line connecting the J and J80 (or J60) points
- ST-segment area: the area between the ECG trace, the isoelectric line and the points J and J80
- T-wave amplitude: the amplitude deviation of the T-wave peak from the isoelectric line

# ST-seg

#### ST-segments feature extraction ...



http://dblab.chungbuk.ac.kr

#### CASE study – I : Ischemia beat classification

- Data generation
  - Download ST-segments data from *PhysioNet database* (<u>http://physionet.org</u>)
  - Total # of data : Ischemia(204), Normal(198)
  - Selected features for our experiment
    - ST-deviation, Slope, Area, T amplitude, Class (Ischemia/Normal)

#### Experimental results



| Classifier | TP    | FP    | Precision | Recall | Class    |
|------------|-------|-------|-----------|--------|----------|
| CMAR       | 0.924 | 0.122 | 0.885     | 0.924  | Ischemia |
|            | 0.878 | 0.076 | 0.919     | 0.878  | Normal   |
| SVM        | 0.909 | 0.096 | 0.909     | 0.909  | Ischemia |
|            | 0.904 | 0.091 | 0.904     | 0.904  | Normal   |
| TAN        | 0.818 | 0.234 | 0.786     | 0.818  | Ischemia |
|            | 0.766 | 0.182 | 0.8       | 0.766  | Normal   |
| C4.5       | 0.768 | 0.16  | 0.835     | 0.768  | Ischemia |
|            | 0.84  | 0.232 | 0.775     | 0.84   | Normal   |
| MDA        | 0.838 | 0.277 | 0.761     | 0.838  | Ischemia |
|            | 0.723 | 0.162 | 0.81      | 0.723  | Normal   |

## Outline

- Introduction
  - Application of data mining to medical data
  - Motivation and objective
- Background
  - ECG Analysis
  - ST-segments
  - Heart Rate Variability
  - Carotid arterial wall thickness
- Automated detection of Ischemic ECG beats
  - ST-segments feature extraction
  - CASE study1: Ischemia beat classification

#### Heart Rate Variability

- <u>ECG processing, Linear/ Nonlinear</u> <u>features</u>
- Classification methods for diagnosing CAD
- CASE study2 : AP and ACS diseases classification
- Dyslipidemia diagnosis methodology
  - Measuring Intima-Media thickness
  - Feature extraction from IMT image
  - CASE study3: dyslipidemia classification

#### Discussion

- On going work
  - Development of Contrast data mining and frequent pattern mining methods

#### Conclusion

#### HRV analysis : ECG processing

Measuring ECG for each 3 position

 The ECG signals are recorded during 5min. for each 3 recumbent position (Supine, Right, Left)



#### HRV analysis : ECG processing ...

#### Time series of RR intervals

- RR intervals are the series of time intervals between heartbeat
- The ECG are retrieved to measure the RR intervals by using S/W for the detection of R waves : Extract the R-peaks from the ECG recordings
- RR intervals time series were resampled at rate of 4Hz to obtain power spectral density (frequency domain analysis)



#### HRV analysis : Linear Feature Extraction

- Time domain of Linear Features
  - Simple time domain features include :
    - **RRm**: the mean **RR** intervals
    - SDRR: the standard deviation of all RR intervals
    - **SDSD**: the standard deviation of differences between RR intervals
- Frequency domain of Linear Features
  - After Calculating the heart rate, performed the power spectrum analysis
  - Define areas of spectral peaks as follows:
    - 1 Total Power (TP): 0~0.4Hz
    - very Low Frequency (VLF) power: 0~0.04Hz
    - 3 Low Frequency (LF) power: 0.04~0.15Hz
    - High Frequency (HF) power: 0.15~0.4Hz.
    - Ratio LF/HF : LF to HF ratio

Normalized LF (nLF): 
$$nLF = \frac{(TP - VLF)}{LF} \times 100$$
 Normalized HF (nHF):  $nHF = \frac{(TP - VLF)}{HF} \times 100$ 

## HRV analysis : Nonlinear Feature Extraction

#### Poincare Plots

- The Poincare Plots is constructed by plotting each RR interval against the previous one
- SD1: the standard deviation of the distances of points from y=x axis
- SD2: the standard deviation of the distances of points from y= -x + RR<sub>avg</sub>
- SD2/SD1: the ratio SD2/SD1 is the measure of heart activity

#### Approximate Entropy

- The ApEn describes the complexity and irregularity of the signal
- The value of ApEn is low in regular time series (bad health) / high in complex irregular ones (good health)

### HRV analysis : Nonlinear Feature Extraction ...

#### Hurst Exponent, H

- H is the measure of the smoothness of a time series based on asymptotical behavior of the rescaled range of the process
- H = 1 : regular behavior (bad health)
- 0.5 < H < 1.0 : normal (good health : 0.7)</p>
- H = 0.5 : random behavior

#### CASE study – II : AP and ACS classification

Real ECG datasets : Patients in five university hospitals

- Patients with stenosis by using angiography
  - Luminal narrowing >= 50% : CAD Group
  - The others: Normal Group

CAD group was divided into two group by cardiologists, AP (Angina Pectoris) and ACS (Acute Coronary Syndrome)







Collect 193 ECG data from our ECG DB

99 patients with CAD group and 94 normal group were studied

지금카드시스테

#### CASE study – II ... : Classifiers Evaluation Results

#### Analysis : Use of multiparametric features of HRV

| Classifier | TP rate | FP rate | Precision | Recall | F-measure | Class   |
|------------|---------|---------|-----------|--------|-----------|---------|
|            | 0.925   | 0.103   | 0.933     | 0.925  | 0.929     | AP      |
| SVM        | 0.820   | 0.072   | 0.781     | 0.820  | 0.800     | Control |
|            | 0.815   | 0.025   | 0.855     | 0.815  | 0.835     | ACS     |
| Davasian   | 0.957   | 0.200   | 0.881     | 0.957  | 0.917     | AP      |
| (TAN)      | 0.730   | 0.044   | 0.839     | 0.730  | 0.781     | Control |
| (1AN)      | 0.692   | 0.031   | 0.804     | 0.692  | 0.744     | ACS     |
|            | 0.792   | 0.255   | 0.828     | 0.792  | 0.810     | AP      |
| MDA        | 0.520   | 0.066   | 0.712     | 0.520  | 0.601     | Control |
|            | 0.692   | 0.163   | 0.437     | 0.692  | 0.536     | ACS     |
|            | 0.918   | 0.182   | 0.886     | 0.918  | 0.902     | AP      |
| C4.5       | 0.780   | 0.059   | 0.804     | 0.780  | 0.792     | Control |
|            | 0.631   | 0.051   | 0.695     | 0.631  | 0.661     | ACS     |
|            | 0.945   | 0.164   | 0.899     | 0.945  | 0.922     | AP      |
| CMAR       | 0.740   | 0.044   | 0.841     | 0.740  | 0.787     | Control |
|            | 0.692   | 0.054   | 0.703     | 0.692  | 0.698     | ACS     |



|        | Classifier | Class I abel | Predicted Class |         |        |
|--------|------------|--------------|-----------------|---------|--------|
|        | Clussifier |              | AP              | Control | ACS    |
|        |            | AP           | 92.55%          | 5.88%   | 1.57%  |
|        | SVM        | Control      | 13%             | 82%     | 5%     |
| Actual |            | ACS          | 6.15%           | 12.31%  | 81.54% |
| Class  | CMAR       | AP           | 94.51%          | 1.57%   | 3.92%  |
|        |            | Control      | 17%             | 74%     | 9%     |
|        |            | ACS          | 15.38%          | 15.38%  | 69.24% |

#### **Confusion matrix**

43 /57

#### http://dblab.chungbuk.ac.kr

## Outline

#### Introduction

- Application of data mining to medical data
- Motivation and objective
- Background
  - ECG Analysis
  - ST-segments
  - Heart Rate Variability
  - Carotid arterial wall thickness
- Automated detection of Ischemic ECG beats
  - ST-segments feature extraction
  - CASE study1: Ischemia beat classification

- Heart Rate Variability
  - ECG processing, Linear/ Nonlinear features
  - Classification methods for diagnosing CAD
  - CASE study2 : AP and ACS diseases classification
- Dyslipidemia diagnosis method
  - Measuring Intima-Media thickness
  - Feature extraction from IMT image
  - <u>CASE study3: dyslipidemia</u> classification
- Discussion
  - On going work
    - Development of Contrast data mining and frequent pattern mining methods
- Conclusion



## Measuring Intima-Media thickness: KRISS

- IMT can be measured by analyzing the ultrasound images of Common Carotid Artery
  - measure the IMT at the far wall of a designated 1cm length of the right common carotid
  - calculate the average values over 200 points



#### Measuring Intima-Media thickness ... : KRISS

CCA (Common Carotid Arterial) scanning with a highresolution ultrasound system is done

- Exam a cross-section of carotid artery
- Detect four points (a, b, c, d)



#### Cross-sectional view of carotid artery wall

### Feature extraction from IMT image: KRISS

- In order to classify patients with dyslipidemia from normal group, total 15 features are extracted from ultrasound image
- 1. Fv01 : Distance between *a* and *b*
- 2. Fv02 : Area of Fv01[sum(*a*:*b*)]
- 3. Fv03 : Value of *c*
- 4. Fv04 : Distance between a and c
- 5. Fv05 : Area of Fv04 [sum(*a*:*c*)]
- 6. Fv06 : Value of *d*
- 7. Fv07 : Distance between a and d
- 8. Fv08 : Area of Fv07[sum(*a*:*d*)]
- 9. Fv09 : Gradient of straight line from value point of *a* to value point of *c*
- 10. Fv10 : Gradient of straight line from value point of *c* to value point of *d*
- 11. Fv11 : Gradient of straight line from value point of *a* to value point of *b*
- 12. Fv12 : Angle acd
- 13. Fv13 : Angle *cdb*
- 14. Fv14 : Value of *c*-value of *a*
- 15. Fv15 : Value of *c*-value of *d*



#### CASE study – III : dyslipidemia classification using IMT features

#### Collect 600 real IMT dataset

- 328 patients with dyslipidemia and 272 patients with normal group were studied
- Performance evaluation: *Accuracy and confusion matrix*
- Result: suitable method for predicting patients with dyslipidemia : *Neural Net*



|        | Classifier | Class Label  | Predicted Class |        |  |
|--------|------------|--------------|-----------------|--------|--|
|        | Classifier |              | Dyslipidemia    | Normal |  |
|        | Neurol Net | Dyslipidemia | 299             | 29     |  |
|        |            | Normal       | 15              | 257    |  |
|        | SYM        | Dyslipidemia | 271             | 57     |  |
| Actual | 5 4 141    | Normal       | 48              | 224    |  |
| Class  | KNN        | Dyslipidemia | 261             | 67     |  |
|        | KININ      | Normal       | 64              | 208    |  |
|        | C4.5       | Dyslipidemia | 253             | 75     |  |
|        |            | Normal       | 104             | 168    |  |

In this work, in order to predict and diagnose cardiovascular disease such as ischemia, coronary artery disease (AP, ACS)

- Used ST-segments features of ECG
- Suggested a possibility that multi-parametric feature of HRV may be helpful to make a diagnosis the heart disease
- Extract useful features from carotid artery wall image

Also, we applied data mining techniques (classification methods and association rules) to enhance the reliability of the diagnosis and detect early symptoms of heart problems Development of Contrast data mining and frequent pattern mining methods

- Contrast data mining : Emerging Pattern Mining-Based Methodology for Automated Diagnosis of Cardiovascular Diseases
- Frequent pattern mining : Compressed Patricia FP-Tree for Frequent Itemsets Mining



- Propose a special type of Emerging Patterns (IEPs) and show that they are high quality patterns for building accurate classifiers
- Generalize the interestingness measure for EPs, including the min. support, min growth rate, the subset relationship between EPs and correlation based on common statistical measure (chi-square value)
- Develop tree-based pattern growth algorithms for mining only those interesting EPs

#### Emerging Pattern Mining-Based Methodology ... : T-tree & P-tree structure

- Emerging Patterns (EPs) are contrast patterns between two classes of data whose support changes significantly between the two classes
  - In this study, we extend EPs mining algorithm using TFP algorithm



#### Compressed Patricia FP-Tree for Frequent I temsets Mining

#### Compressed Patricia FP-Tree (CPFP-Tree)

- New data structure
- Can find fast frequent itemsets
- Sparse frequent itemsets
- Dense frequent itemsets

#### CPFP-Tree applications

- Finding frequent itemsets
- Association rules mining & other mining
- System with limited memory
- Sparse & dense frequent itemsets



## **CPFP-Tree:** example







## Conclusion

- In this study,
  - We proposed data mining framework for coronary heart disease perdition
  - Automated diagnosis methodology using ST-segments, HRV features and IMT features
  - Our experiments showed that Function -, Pattern-based classifiers were the most accurate in comparison with all other classifiers for diagnosing heart diseases

#### On going work,

- We proposed emerging pattern mining-based Methodology for Automated Diagnosis of Cardiovascular Diseases
  - Extend emerging pattern based classifier tree-based
  - Introduce the interestingness measure, Chi-square
  - Directly generate EPs and remove uninformative patterns
- We developed a novel fast algorithm for discovering frequent itemsets
  - CPFP-Tree is efficient and scalable
  - Dense and Sparse frequent itemsets mining tasks
  - CPFP-Tree is faster than the Apriori and FP Growth algorithm

# Thank You Any Question?

## Good health is the greatest prize!!

#### This work was supported by

- The Korea Science and Engineering Foundation (KOSEF) grant funded by the Korea government(MOST) (R01-2007-000-10926-0)
- The Korea Research Foundation Grant funded by the Korean Government(MOEHRD)" (The Regional Research Universities Program/Chungbuk BIT Research-Oriented University Consortium)
- a grant from the Personalized Tumor Engineering Research Center(PTERC) and the Korea Science and Engineering Foundation(KOSEF).