



Institut de recherche pour le développement



HỘI NGHỊ KHOA HỌC QUỐC GIA LẦN THỨ VII "NGHIÊN CỨU CƠ BẢN VÀ ỨNG DỤNG CÔNG NGHỆ THÔNG TIN" Trường Đại học Công nghệ thông tin và Truyền thông, Đại học Thái Nguyên Thái Nguyên, ngày 19 – 20/6/2014

### Hierarchical models for the rainfall forecast

# DATA MINING APPROACH

Thanh-Nghi Do

dtnghi@cit.ctu.edu.vn

**June - 2014** 

Introduction

Related works Algorithms

Results





#### Aim

 Statistical downscaling local daily precipitation from large scale weather variables obtained by GCM

Future works

# Introduction

### **Contributions**

- Hierarchical models for the rainfall forecast
- Learning the classification model to categorise the day into the class: no rain, light rain, moderate rain, heavy rain and violent rain;
- Training the regression models to predict the rainfall from data of each class
- Experimental results on datasets collected from SEA START show that the our hierarchical model is the most accurate rainfall prediction compared to linear regression, k neighbor neighbors, decision trees, support vector machines

# **Related works**

### **Dynamic downscaling**

high-resolution regional climate models nested in a GCM to obtain local weather variables

#### Statistical downscaling)

- Statistical relationship between large-scale GCM outputs and local weather variables (low cost)
- Weather typing method
- Stochastic weather generators
- Resampling
- Regression

### **Pasini, 2009**

- Neural network: non-linear regression
- Non-linear relationship between large-scale variables (predictors of GCM) and regional/local variables (predictands)

### Tripathi et al., 2006

- Support vector machine: non-linear regression
- Non-linear relationship between predictors & predictands

### **Chen et al., 2010**

- Support vector machines: non-linear classification/regression
- Classification: dry/wet
- Regression: precipitation of wet days



х

 $y = \alpha + \beta x$ 

6

- Simplest model: regression
- Non-parametric
- Numerical data

#### **Future works**

Х

# **Data mining tools**

#### k nearest neighbors

- Simplest models: classification, regression
- Parameter k=1,2,...
- Error bound: 2 x Bayes risk
- Long time for prediction
- Numerical data

# **Data mining tools**

#### Support vector machine

- Accurate models: classification, regression
- SVM + kernel function = model
- Linear/non-linear tasks
- Global optimum
- Hyper-parameters
- Long training time
- Numerical data

**Future works** 

9

**Results** 

### Support vector classification





Future works

# **Data mining tools**

#### **Decision tree**

- Classification, regression
- Simple, popular, good performance, intuitive model
- Non-parametric
- Any type of data
- Less accurate than SVM

#### **Bagging and random forest**

- Ensemble of decision trees: more accurate than a single one
- Classification, regression
- Simple, popular, good performance
- Non-parametric
- Training time: faster than SVM

Introduction Related works

Algorithms

**Results** 

**Future works** 

### **Decision tree**

- Classification: minimize impurity (Shannon entropy, Gini index)
- Regression: minimize variance



### Bagging







### **Data collection**

### Dataset (predictors)

- Website: http://cc.start.or.th/
- 1980-2006: 26 years
- Can Tho: LON = 105.8 and LAT = 10.2
- Tmax (°C)
- Tmin (°C)
- Precipitation (mm)
- Wind speed (m/s)
- Wind direction (degree from north)
- Solar Radiation (W/m<sup>2</sup>)

s 🚺 Results

**Future works** 

# Software programs

#### Algorithms

- Implemented in R language
- Linear regression (LM)
- k nearest neighbors (kNN)
- Decision trees (DT)
- Bagging of decision trees (BagR)
- Support vector machines (SVM) using a RBF kernel
- Hierarchical approach: BagC-BagR, SVC-SVR



Mean absolute error (MAE)

#### Hold-out protocol

- Repeat 3 times
- Randomly splitting the dataset into the training set (2/3 fullset) and testing set (1/3 fullset)
- Learning prediction models from the training set
- Evaluating the prediction models using the testing set
- Averaging the prediction results

ns Results

Future works

### **Prediction results**

Methods	MSE	MAE
Linear regression (LM)	34.1613	4.4062
k nearest neighbors (kNN with $k=5$ )	19.8154	2.5938
Decision tree regression (DT with <i>minobj</i> =5)	15.5085	2.0525
Bagging (Bag with #trees=100)	<u>10.4347</u>	<u>1.4814</u>
Support vector regression (SVR-RBF, $\gamma = 0.01$ , $\varepsilon = 0.1$ , $C = 10^4$ )	17.0064	2.4554
Hierarchical model (BagC-BagR with #trees=100)	8.9708	1.4337
Hierarchical model (SVC-SVR-RBF)	20.0745	2.4058

Introduction Related works Algorithms

hms

Results

Future works

### **Prediction results**

#### A comparison of rainfall forecast models



20

# **Prediction results**

#### Comments

- Linear regression: not suited for rainfall forecast
- Simple models (kNN, Decision trees): very competitive
- Non-linear models (SVM, Bagging) and hierarchical approach (BagC-BagR, SVC-SVR): most accurate
- Bagging, BagC-BagR: non-parametric, more accurate than SVM



**Future works** 

### To do

Preprocessing: smoothing and cleaning data

Introduction Related works Algorithms

- Improving prediction results of hierarchical models
- Hierarchical models: non-linear regression models for other problems, including water level, number of telephone calls, etc.

### **Thanks for your attention!**

