## CS 458/658: Introduction to Data Mining

## Data Mining Course Project

Instructor: Lei Yang

Department of Computer Science and Engineering, UNR – Fall 2018

Tue, Thur 12:00PM - 1:15PM, WRB 4050

### Project

- One key goal of this course is to take advantage of your intelligence and (limited) experience (so you're audacious and creative) to expand your knowledge in creating something useful and interesting
- Group project
  - 3 students per group
    - 27 undergraduates —9 groups
    - 9 graduates —3 groups
  - You can apply whatever techniques you learnt from data mining course and other sources

#### **Tasks**

- Tasks
  - Task 1: Document Classification
  - Task 2: Exploring Environmental Data at NRDC
  - Task 3: Exploring Used Auto Purchase Dataset
  - Task 4: Exploring Intent data for B2B marketing and sales (TBD)

## Undergraduates

- **◆**Tasks:
  - Task 1
  - Pick one of the following:
    - Task 2
    - Task 3
- Bonus:
  - One or two of the remaining tasks
  - Other problems in Task 2

### Graduates

- **♦**Tasks
  - Task 1
  - Task 2
  - Task 3
- Bonus:
  - Task 4 ?
  - Other problems in Task 2

#### **Evaluation**

- Final report (due Dec 12, 2018 in Webcampus)
  (35%)
  - Each member need to submit your own report and indicate your contribution in %
- Class presentation and/or demo (5%)
  - Each group will present their work. Each member needs to present.
  - Your presentation will be evaluated by the other groups using an evaluation form.
  - Each presentation is 20 mins with 5 mins for Q&A.
    - Nov. 27, 2018
    - Nov. 29, 2018
    - Dec. 4, 2018
    - Dec. 6, 2018

#### Task 1: Classification

- Provided data
  - The training set and its label information
  - The testing set
- Hidden data
  - The label information of the testing data
  - The data will be used for the purpose of evaluation

#### **Data Format**

- The training set
  - training.txt
  - The first column is the information ID
  - The second column is the feature ID
  - The third column is the value of the feature
  - The default values of features are zeros

1 72 1

#### **Data Format**

- The label information of the training set
  - label\_training.txt
  - Each row represents a data point in the training set
  - 1 is true information while -1 is misinformation

1 -1 -1 1 1 -1 -1

#### **Data Format**

- The testing set
  - testing.txt
  - It has the same format as the training set

```
1 16 1
1 23 1
1 27 1
1 29 2
1 50 1
1 245 1
1 340 1
1 388 1
1 589 1
1 638 1
1 764 1
1 902 1
1 905 1
1 2774 1
1 8066 1
1 10762 2
```

## Model Challenge from Model Selection

- There are so many classifiers
  - Which one is better?

- There may be parameters in classifiers
  - How to determine the optimal values?

#### Evaluation

- Classification accuracy will be used to evaluate the quality of the predicted labels
- Comparing the hidden labels with your predicted labels
- Your final grades will strongly depend on the rankings of the quality of the predicted labels you provide

## Task 2: Exploring Environmental Data at NRDC

#### **Data**

http://sensor.nevada.edu/SENSORData
Search/

### Sample project

- Wind speed prediction
  - Predict the quantitative wind speed at different sites using the historical information in the same sites as well as data in other neighbor sites.
- The data is available from

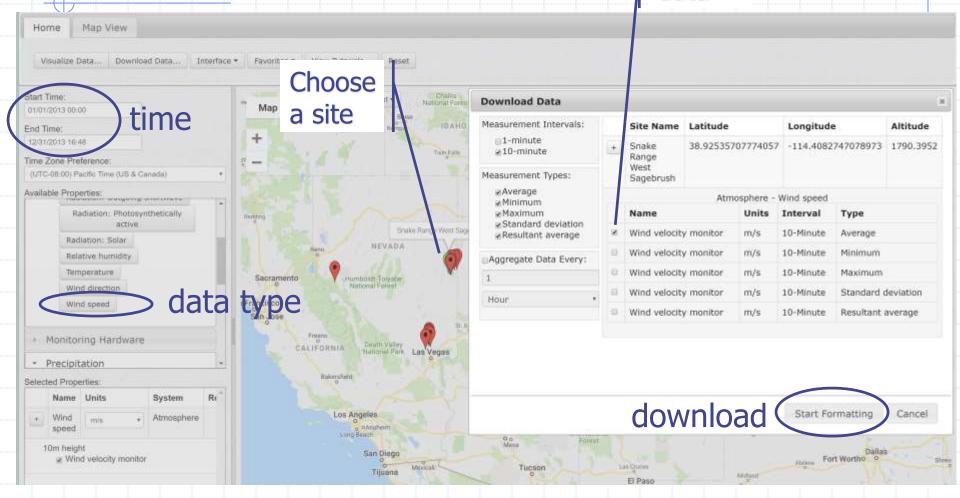
http://sensor.nevada.edu/SENSORDataSearch/

- Snake Range East Sagebrush (EB)
- Snake Range East Subalpine (EA)
- Snake Range East Salt Desert Shrub (ED)
- Snake Range West Subalpine (WA)
- Snake Range West Montane(WM)
- Snake Range West Sagebrush (WB)
- Snake Range West Pinyon-Juniper(WP)

10 minutes wind data

#### Download data

Choose data



### Example data

Site Name:, Snake Range West Sagebrush Deployment: , Wind velocity monitor Monitored System:, Atmosphere Measured Property:, Wind speed Vertical Offset from Surface: , 10m height Units:,m/s Measurement Type:, Average Measurement interval:,00:10:00 Time Stamp ((UTC-08:00) Pacific Time (US & Canada)) 1/1/2013 12:00:00 AM, 0.513648960000000000 1/1/2013 12:10:00 AM, 0.073761600000000000 1/1/2013 12:20:00 AM, 0.348691200000000000 1/1/2013 12:30:00 AM, 0.291023040000000000 1/1/2013 12:40:00 AM, 0.476544640000000000 1/1/2013 12:50:00 AM, 0.435864000000000000 1/1/2013 1:00:00 AM,1.001369600000000000 1/1/2013 1:10:00 AM,0.899444480000000000 1/1/2013 1:20:00 AM,0.206979520000000000 1/1/2013 1:30:00 AM, 0.604398080000000000 1/1/2013 1:40:00 AM,0.710793600000000000 1/1/2013 1:50:00 AM, 0.430052480000000000 1/1/2013 2:00:00 AM, 0.198485760000000000 1/1/2013 2:10:00 AM, 0.175239680000000000 1/1/2013 2:20:00 AM,0.598139520000000000 1/1/2013 2:30:00 AM,1.322791360000000000 1/1/2013 2:40:00 AM, 0.473415360000000000 1/1/2013 2:50:00 AM,0.105948480000000000

1/1/2013 3:00:00 AM, 0.760862080000000000

## Challenges

- Data preprocessing
  - Missing data

What methods to use?

How to tune parameters?



#### **Evaluation**

- For each site, you need to provide prediction accuracy of your proposed approach based on the following measure
  - Mean absolute error (MAE)

$$MAE = \frac{1}{number\ of\ points} \sum |forcast - actual|$$

Root mean squared error (RMSE)

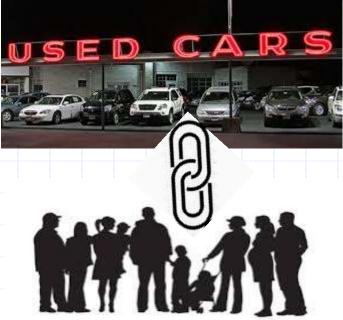
$$RMSE = \sqrt{\frac{1}{number\ of\ points}} \sum |forcast - actual|^2$$

- Compare your approach with the following Benchmark:
  - Persistent forecast: predicted\_wind(t)=actual\_wind(t-1)

## Task 3: Exploring Used Auto Purchase Dataset

## Exploring Used Auto Purchase Dataset (1)

- Dataset: the set of all used auto purchases for the past 5 years in US
  - Number of attributes: 280
    - Vehicle Info (Model, Engine, Drive Type) Home Info (Purchase Price/Date, Value, Year), Address (State, County), Loan Info (Monthly Mortgage), Demographic (Ethnic, Number of Children), Behavior Info (Investment, Interest in Travel/Reading, Presence of Premium Credit Card) .....
    - No Personally identifiable information



# Acknowledgement: Thanks Marketing Evolution for sharing this dataset.

## Dataset(UsedAutoRELEVATEfirst1 0000-noLatLong.csv)

	Attribute 1	Attribute 2	
Data entry 1			
Data entry 2			

	6 A	8	c	D E	F	Ğ	н		J K		M	N O		Q	R	5	t ·
2	Customer	Home Purc La	titude inve	stmem Living Unit C	ounty Carit	ome Lanc	Acctif	Longitude	Home LancCarrier	Roi Vehicle 1	THome ImpiA	ea Code Investn	nem Home Tax 'Ho	ome Year	Transactio: Ho	ime BascHo	me Lan
4 5 6 7	XV6311.HZ	99		1.83E+09	339	0	4.59E+08		O R028	SE	133		2596	1997	20160512	15	9
9 11		(SbZPoB99Iv1	V5wuBoN3T-	WYHPXKkzgTt9qCq0	iAVw		4,46E+08		C022	LT					20161101		
11	UNMATCHE	D					4.11E+06		C002	BASE					20161205		
11 11 11 21	XY6311DIN	1050		1.04E+09	3	0	4.52E+08		0 C002	LX .	845		9999	1988	20170130	O	503
21 22 23 23	XV6311mKn	Dv9rrKvadj9	»AFXNIW2Gm	gZccDrArOuAdwRHi	AcE		3.7E+08		C005	SLT					20160301		
25 25 25 25 25 25	XV6311eBy	CBp112diTWF	cRuDgl9YaYR	3Ny42TJun6Yzk1hD	ti		3.85E+08		C027	SLT					20151229		
	XY63114y5	85		2E+09	25	0	4.36E+08		O R014	BASE	85	443	1372	1996	20161101	30	35
30		0		1,935+09	153	117	4.04E+06		65 C044	BASE	95	515	2994	1971	20161206	0	21
-		RELEVICIESMEN	BBD-your (E)														

## Data Dictionaries (EXP REL Custom.xls)

ID	Field Name	Description	
Data entry 1			
Data entry 2			

Field Name	Long Description	Start Position	End Position		Field	Туре		Pielé Velues
1415 Address ID	Address ID - Uvique identifier assigned to each address in the Consumer/New repository. The Address ID remains with an address even in the event that the occupants relocate. Values: 10 byte names:	- 1	10	1	0 AN		9509990099	
0337 State Code	State Code	11	11		2 AN		95	D2-ALABAMA, O2-ALASKA, O4-ARIZINA, D3-ARKANSAS, D6-CAUPORNA, D8-COLORADO, D9-CORNO COLUMBIA, (25-ELORIDA, 13-EGORGIA, 35-HAWAU, 15-EDAHO, 17-ELUMOS, 38-INDUMAA, 15-EDAH 13-AMARIE, 24-MARILAND, 25-MASSACHUSETTS, 26-MICHEGAN, 27-MINRESOTA, 26-MISSISSIPPE, A, 32-REYARA, 35-MEW HAMPISHRIT, 34-REW LENSEY, 35-MEW MEXICO, 36-MEW TORK, 27-MICHTE DAXDTA, 36-CHID, 48-CILLAHOMA, 41-OREGION, 43-PENNISPIYANIA, 44-RHODE ISLAND, 43-SOUTH DAXDTA, 47-TENNESSEE, 48-TEXAS, 49-LUTAH, 50-YERMONT, 51-VIRGINIA, 53-WAWANGTON, 54-W YARGINIA, 53-WINCONSIN, 56-WYOMMIN.
10114 State Attenvision	State abbreviation	13	14		2 char			ARTIALASKA, ALI ALABAMA, ARTARKARSAS, AZTARZONA, CATCALIFORMA, COTCO GRADO, CTTCONI COLLIMBRA, DECICLAMARE, FLATLORIDA, SAN-GEORGIA, HITHAMANI LA LOWIN LIDEN AND LIDEN AND ALI CONTROL CHILLINGO ALI CURRIAN, AND THA CAROLINA, AND TRESTRI DANOTA, NETHERBASKA, AND THIN HAMPINES, DITARKAN, AND TROST IN CAROLINA, AND TROST IN CAROLINA, SCHOOL OF A PROPERTY DANOTA, NET REPRESENTANT OF A PERMISSILVANIA CAROLINA, SOUSCITHI DANOTA, TINTENNESSEE, TANTEXAS, UT-UTAH, VAN-FROINIA, VT-VERMONT, WAN-WARDEN NUTON, WINDERINA, WYN-WYSDINA.
10581 Dp Code	Zin Code	15	15	3.5	Scher			200000000000000000000000000000000000000
10579 Zip+4	Zin+4	30			a char			
13272 Delivery Point har code	DirectOPV - Delivery Point barcode / Check digit	24			char			
£1357 Carrier Route	carrier route code	27	- 30		4 cher			
10017 WORKFLOW FIELD Short City Name to Inverted V2	o be apecial 155yte field - field to FCARD for 20 byte field	31	43	1	Skfwr			
10376 City Name	Otyname	44	.71	- 21	9 xbar			
11247 House Number	Primary (hoose) number	72	81	1	0 cher			
\$1249 Pre Offection	Street pre-directional	82	83	1	2 cher			E-East, N+North, NE-Northeast, NW+Northwest, S-South, SE-Southwest, SW-Southwest, W+West,
11003 Street Name	Street neme	84	111	23	S shur			
10633 Street Suffix	Street saffa	112	137	b) 3	i char			ALT-ALLT AND-ANNEX ARC-ANLADE AVE-AVENUE, BOTHERADE, BOHERADE, BE-BLITE, BUTS-BLUI AND-AREAS, BRIG-BRODE, BRIC-BRODO, BRODO, BR

## Exploring Used Auto Purchase Dataset (2)

- Project description:
  - Selection: Due to the size/heterogeneity of the original data, we need to select a target data.
  - Preprocessing: Data exist in many types (continuous, nominal) and forms, and may have missing values.
  - Transformation: To better extract useful patterns from dataset.
  - Data mining: Explore different data mining algorithms
  - Interpretation/Evaluation
- Goal: extraction of useful patterns from dataset
  - What car type will be purchased, given customer's info?
  - What customer's type, given a car purchased?
  - How to divide a market into distinct subsets of customers? 24

### Report Format

- Cover Page
  - Team members and their contribution in %
- Introduction
- Literature review for each task
- Task 1
  - Your approach (e.g., Preprocessing, Model selection, Parameter selection, Your solution)
  - Your conclusion
- Task 2
- Task 3
- Task 4
- List of documents/codes you submitted

## Report requirements

- The report should be as concise as possible while providing all necessary information required to replicate your plots.
- The plots should contain multiple curves and can be formatted so that many plots can fit on one page (so that your report is not longer than it should be).
  - Don't use screenshot!
- In literature review, you need to show your understanding of the literature by reading and comparing the existing work.
  - Cite your references properly. You can use google scholar to download citation.
- Your submitted code should have proper comments.