Assignment Project Exam Help

https://eduassistpro.github.io/

Preparin

Add WeChat edu_assist_pro

Prof. Vibs

The Paul Merage School of Business

University of California, Irvine

# Agenda

- Information Theory

- Reminders
  - Assignment 1 posted on Canvas
  - Form grou

- Install Weka

- Working with Data

# From Probability to Information
# Information Theory

- Makes use of the probabilistic relationship between attributes to identify how much information one attribute provides on the other

  - Useful to und                                              ttributes
  - Can also be u ess attributes

- Information = surprise

  - How much surprise is created by
  - Information = expectation – realization

# Logarithm

- $\log_b(X)$ is read as "log of $X$ with base $b$"
  - Microsoft Excel : "=log($X$,$b$)"
  - What does it mean

  - If $Y = \log(X)$, then $X = b^Y$
  - Base 10: lo
    - Microsoft Excel : "=log($X$)"
    - If $Y = \log_{10}(X)$, then $X = 10^Y$

      - If $\log_{10}(1000) = 3$, and $1000 = 10^3$
  - Natural logarithm = $\ln(X) = \log_e(X)$, where e=2.7183
    - Microsoft Excel : "=ln($X$)"
  - Logarithm with base 2
    - Microsoft Excel : "=log($X$,2)"
    - $\log_2(X) = \log_{10}(X)/\log_{10}(2) = 3.3219\ \log_{10}(X)$

UCIrvine | THE PAUL MERAGE SCHOOL OF BUSINESS

# Information Theory

- Entropy of a distribution
  - Let $X$ be a random variable with the probability distribution
  
  $Pr[X=x_i] = p_i, i=1,2,\ldots,n$, where

- Entropy of $X$ (le

  $$H(X) = H(p_1, p_2, \ldots, p_n)$$

- Let $Y$ be another random variable tributed)
  - Knowledge of $Y$ reduces the uncertainty and hence entropy of $X$.
  - Therefore, $Y$ provides the following information about $X$:

- $I(X;Y) = H(X) - H(X|Y)$.
  - Thus $I(X;Y)$ is called Mutual Information

# Properties of Information Measure

- $I(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = I(Y;X)$

- If $X$ and $Y$ are i
  - $H(Y|X) = H($
  - $I(X;Y) = 0$

# Example

- Consider 10 balls in a basket
  - 4 large and red, 1 small and red, 2 large and blue, and 3 small and blue
  - You are to pick o                                        r without looking
- Strategy
  - Check the size of the ball and predict
    - Red if it is large (67% accurate — 4
    - Blue if it is small (75% accurate — 3 out of 4)
- Without the size information, you can only be 50% accurate
- Clearly, size provides information about the color
  - We know that since size and color are not independent
  - Color provides information about the size, as well

UCIrvine | THE PAUL MERAGE SCHOOL OF BUSINESS

# *I*(Color; Size)

- *I*(Color; Size) = *H*(Color) – *H*(Color | Size)
- Without size information:

  *H*(Color) = H()  =1

- With size information:

  *H*(Color | Size = la<span style="color:red">Assignment Project Exam Help</span>

  *H*(Color | Size = s

  *H*(Color | Size) = *H*(Color | Size = large) P(Size = large)

  + *H*(Color | Size = small) P(Size = small)

  =0.918 × 0.6 + 0.811 × 0.4 = 0.875

- Information gain = 1 – 0.875 = 0.125 bit
  - Size, on average, provides 0.125 bit of information on color

# *I*(Size;Color)

- *I*(Size; Color) = *H*(Size) – *H*(Size | Color)
- <span style="color:red">Without color information:</span>

   <span style="color:red">H(Size) (A.69.4)</span> <span style="color:red">Assignment Project Exam Help</span>

- With color i

   *H*(Size | Colo <span style="color:red">https://eduassistpro.github.io/</span>

   *H*(Size | Color = blue) <span style="color:red">Add WeChat edu_assist_pro</span>

   *H*(Size | Color) = *H*(Size | Color = red)× P(Color = red) +
       H(Size | Color = blue) × P(Color = blue)


- Information gain =

# Loan Application Data

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Contingency Table
(Expressing relationship between two attributes)

*Compute H*(Liability)
& *H*(Liability | CR)
& *I* (Liability; CR)

| | | Liability | | |
|---|---|---|---|---|
| | | normal | high | **Total** |
| **CreditRating** | excellent | 3 | 1 | **4** |
| | good | 4 | 2 | **6** |
| | poor | 0 | 4 | **4** |
| | Total | 7 | 7 | **14** |

$H$(Liability) = H

$H$(Liability | CR = excelle

$H$(Liability | CR = good) =

$H$(Liability | CR = poor) = $H$ ( 0

$H$(Liability | CR) = 0.811 × () + 0.918 × () + 0 ×() = 0.625

$I$ (Liability; CR) = 1 − 0.625 = 0.375

⇒ $I$ (CR; Liability) = 0.375.

# Entropy and Gain Ratio

- Even though the mutual information between two random variables is always symmetric, observing or recording one variable may be more difficult than the oth
  - The more ue, higher is the level of this difficulty
  - Entropy measures this diffi

# Gain Ratio

- Gain ratio ($G$) measures the information gain relative to the level of difficulty of coding the attribute
- $G(X; Y) = I$
- $G(Y; X) = I (Y; X) / H(X$
- $G(X; Y) \neq G(Y; X)$

# *G*(CR;Liability) & *G*(Liability;CR)

*I(CR; Liab*

*H(Liability) = H(½, ½) = ½log* ... *) = 1*

*H(CR) = H(4/14, 6/14, 4/14)*

*G(CR; Liability) = I(Liability; CR) / H(Liability) = 0.375*

*G(Liability; CR) = I(Liability; CR) / H(CR) = 0.241*

*G(CR; Liability) ≠ G(Liability; CR)*

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Prof. Vibs

The Paul Merage School of Business

University of California, Irvine

# Steps in Data Mining

1. Develop an understanding of the purpose of the data mining project
2. Obtain the data set to be used in the analysis
   - random sampling from a large database to capture records
   - While data mining deals with very large databases
     - usually the analysis to be done requires only thousands or tens of thousands of records

Assignment Project Exam Help

3. Explore, clean, and
   - This involves verify ___ ition.
   - How should missin https://eduassistpro.github.io/
   - Are the values in a reasonable range, given ___ xpect for each variable?
   - Are there obvious "outliers?" Add WeChat edu_assist_pro
   - The data are reviewed graphically - for exa ___ catterplots showing the relationship of each variable with each other variable

4. Reduce the data, if necessary
   - eliminate unneeded variables
   - transforming variables
   - creating new variables

# Steps in Data Mining

5. Determine the data mining task
   - classification, prediction, clustering, etc.
6. Choose the data mining techniques to be used
   - regression, neural nets, hierarchical clustering, etc.
7. Use algorithms to perform the task
   - This is typically an iterative process
   - Choosing different variables or settings within the algorithm
8. Interpret the results
   - Recall that each ~~validation~~ data for tuning purposes
     - validation data becomes a part ~~process~~
     - likely to underestimate the err ~~yment of the model~~ that is finally chosen
9. Deploy the model in real world
   - For example, the model might be applied to a purchased list of possible customers
   - action might be "include in the mailing if the predicted amount of purchase is > \$10"

UCIrvine | THE PAUL MERAGE SCHOOL OF BUSINESS

# Data Types

- Variable Measures
  - Categorical variables (e.g., CA, AZ, UT…)
  - Ordered variables (e.g., course grades)
  - Numeric variables (e.g., money)
- Dates & Times
- Fixed-Length C                                    Codes)
- IDs and Keys –                                    ata in other tables
- Names (e.g., Company Names)
- Addresses
- Free Text (e.g., annotations, comments, memos, email)
- Unstructured Data (e.g., audio, images)

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Nominal quantities

- Values are distinct symbols
  - Values themselves serve only as labels or names
- Example: att                                        weather data
  - Values: "s                                    ny"
- No relation is implied am                          al values (no ordering or distance measure)
- Only equality tests can be performed

# Ordinal quantities

- Impose order on values

- But: no distance between values defined

- Example: <span style="color:red">Assignment Project Exam Help</span>
  attribute "te                                                    data
  - Values: "ho <span style="color:red">https://eduassistpro.github.io/</span>

- Note: addition <span style="color:red">Add WeChat edu_assist_pro</span> sense

- Distinction between nominal and ordinal not
  always clear (e.g. attribute "outlook")

# The ARFF format

```
%
% ARFF file for weather data with some numeric features
%
@relation weather

@attribute outl                                    ny}
@attribute temp
@attribute humi
@attribute windy {true, false}
@attribute play {yes, no}


@data
sunny, 85, 85, false, no
sunny, 80, 90, true, no
overcast, 83, 86, false, yes
...
```

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

UCIrvine | THE PAUL MERAGE SCHOOL OF BUSINESS

# Additional attribute types

- ARFF supports *string* attributes:

```
@attribute description string
```

  - Similar to ... list of values is not pre-sp

- It also supports *date* attri

```
@attribute today date
```

  - Uses the ISO-8601 combined date and time format *yyyy-MM-dd-THH:mm:ss*

# Sparse data

- In some applications most attribute values in a dataset are zero

  - E.g.: word counts in a text categorization problem

- ARFF supports sparse data

```
0, 26, 0,  0, 0
0,  0, 0, 42, 0, 0,  0, 0, 0,
```

```
{1 26, 6 63, 10 "class A"}
{3 42, 10 "class B"}
```

- More details about ARFF:

  - http://www.cs.waikato.ac.nz/~ml/weka/arff.html

UCIrvine | THE PAUL MERAGE SCHOOL OF BUSINESS

# Sampling Data

- Sampling can be used to create better data sets (training or testing) to build better models.

- Random sampling techniques:

  - **Simple random sampling**

  - **Proportionate st**                                         ntative sample.
    Divide data set in                                     'Private' travelers.
    Assuming the pro                                     0% Private, and we needed
    at least 100 Private travelers for our m             andomly select 900
    Business travelers.

  - **Disproportionate stratified sampling**: Select a weighted sample.  Also called 'oversampling': used if a particular group of examples is important but not well represented in the data set.
    e.g. In direct mail response prediction you might select 10 responders in the dataset for every non-responder you select.  For claims analysis, you might weigh the fraudulent claims (which are often naturally rare).

# Missing Value Treatment

- **Reason for missing?**
  - Not recorded
  - Not applicable
  - Customer refused to provide

- **Dealing with m**
  - Delete the records with missing val
  - Add flag fields ('address_missing'
  - Estimate missing value:
    - Use average over entire data set
    - Use average over similar records
    - Use an advanced prediction technique

# Noise in Data

- The biggest challenge with noisy data sets is that it is difficult to identify noise.

- In some spe                                    identified

  - Value out                                    )

  - Meaningless value (e.g., a li

    someone without a license)

  - Mismatched value (e.g., City, State, and PIN not matching against the postal database)

# Attribute Selection

- Smaller attribute sets are simpler to understand, but may produce an overly simplistic model

- Larger attribute sets may lead to overfitting

- Eliminate useless at
  - Related to redundan

- Attribute consolidat
  - Combine a set of binary attributes into one

- Attribute expansion
  - Expand a nominal attribute into a set of binary ones

- Attribute conversion
  - Change the data type of an attribute

# Formal Dimension Reduction

- If you have multiple highly correlated columns, then reduce number of columns
  - e.g. height
- Principal co                                        A)
  - Subsets of numeric (not cat                        riables
    - measured on the same scale
    - highly correlated
  - Come up with few variables (one or two or three)
    - that are weighted linear combinations of original variables
    - retain the explanatory power of the original data

# Principal Components

The line z1 is the direction in which the variability of the points is largest.

**1st principal component**

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Attribute Consolidation

- Example 1: Suppose you have two 0/1 attributes: "Male" and "Female"
  - A row of d                          the attributes
- At the same
  - Create a new nominal attr                r" with two possible values — *male and female*

# Attribute Expansion

- Attribute expansion is the opposite of attribute consolidation
  - A nominal set of 0/1 attribute
- Set-values a
  - Example: Hobby, Genre,
- It can be replaced by a set of binary attributes

# Attribute Conversion

- Ratios
  - e.g. Try income divided by number of employees, to get a measure of productivity per employee

- Derived Values
  - e.g. derive cust                                      *thdate* (or *production-date*), as age m

- Changing the data type of attrib
  - Nominal to numeric or vice versa

# Binning

- Binning (discretization) converts numeric values to discrete categories. e.g. low-income is <= 30, high-income is > 30

- For example:
  - Equal-Interval binning
    - Bin intervals of equal width, irrespective items per bin
  - Equal-Frequency binning
    - Equal number of items per bin, irrespective of bin width

Bins

| | 31-40 | 41-50 |
|---|---|---|
| 26 | 35  37 | 42  45  48 |

| 21-26 | 29-37 | 38-48 |
|---|---|---|
| 25  26  26 | 29  35  37 | 42  45  48 |

# The entropy of a random variable is higher when

A: It has many different states, each of which has low likelihood

B: It has very f ch has high likelihood

C: It has many different stat nly a few states have very high likelihood

D: It has very few states

E: None of the above

The file format used by WEKA is called

A. DOCX

B. XCL

C. WEK

D. ARFF

E. TXT

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# When to Normalize Data?

- Rescale attributes to the range of 0 to 1
  - Subtract the min, and divide by (max – min)
- Results in all variables getting equal importance
- Not advisable
  - When the uni~~ts~~ for the variables (e.g. dollars), and _____portance
    - e.g. sales of jet fuel, sales of hea~~ting oil~~
- Advisable
  - if the variables are measured in quite differing units
    - unclear how to compare the variability of different variables
    - e.g. dollars for some, parts per million for others
  - or if variables measured in the same units, but scale does not reflect importance
    - e.g. earnings per share, gross revenues

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Data Preprocessing using Weka

- Download file 4bank-data.csv from Canvas

- Follow steps on the following page.

- http://facwe ... her/classes/ect5 84/WEKA/p

# RFM, Pivot Tables and London Jets Data

- [http://www.dbmarketing.com/articles/Art149.htm](http://www.dbmarketing.com/articles/Art149.htm)

- London Jets Data in Excel format posted on Canvas for RFM analysis and Pivot tables.

  Assignment Project Exam Help

  - Do RFM anal
    https://eduassistpro.github.io/
  - Think about strategies that Lo                    uld use to revive their fortunes

    Add WeChat edu_assist_pro

- Go to http://office.microsoft.com/en-us/
  - Search for "Pivot Table" and read up on creating and using them

# Next Session

- Classification using Exact Bayes & Naïve Bayes

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro