

Assignment Project Exam Help  
Other Techniques  
<https://powcoder.com>

Add WeChat powcoder  
Prof. Vibh Abhishek

The Paul Merage School of Business  
University of California, Irvine

# Agenda

- Term Project Presentations next week
- Upload presentation file to Canvas at least 1 hour before class
- Overview of other techniques
- Wiki for contributing final exam questions
  - <https://docs.google.com/document/d/1LFkkveDdileus5zJOg-LfU5siOZT8ObUR0GrSbF3iVE/edit?usp=sharing>

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Attribute Selection

- Weka – Correlation Based Feature (CFS) Selection  
– CfsSubsetEval
- A good feature subset is one that contains features highly correlated with (predictive of) the class, yet uncorrelated with (not predictive of) each other.
- CFS is a fully automatic algorithm -- it does not require the user to specify any thresholds or the number of features to be selected, although both are simple to incorporate if desired

## Other Methods

- Text Mining
- KNN **Assignment Project Exam Help**
- Collaborative filtering **<https://powcoder.com>**
- Logistic Regression
- Support Vector Machines (SVM) **Add WeChat powcoder**
- Neural Nets
- Bagging
- Boosting

# Why Text Mining?

- What can be discovered from text?
- Significant proportion of information of great potential value is stored in documents:
  - News stories pertaining to competition, customers & the business environment at large
  - Technical reports on new technology
  - Email communications with customer, partners, and within the organization
  - Corporate documents embodying corporate knowledge and expertise
  - Legal documents --- automatic reasoning

# Opportunities

## Finding patterns in text:

- Identify and track trends in industry - associations
  - What are my competitors doing?
  - What relevant products are being developed?
  - What are the potential usage of my products?
- Identify emerging themes in collections of documents -cluster
  - Customer communications: cluster messages, each segment identifies a common theme such as complaints about a certain problem, or queries about product features.
- Automated categorization of e-mails (**Spam Filter!**), web pages, and news stories – classification

# Structuring Textual Information

- Many methods designed to analyze structured data
- If documents can be represented by a set of attributes
  - can use existing data mining methods
- How to represent a document ?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Structured  
representation



Apply DM methods  
to find patterns  
among documents

# Text Mining Concepts

- Document
- Token or term
- Corpus
- Bag of Words
- Stop word elimination; Stemming; all lower case
- Term Frequency (TF)
- Inverse Document Frequency (IDF)
- TFIDF
- N-gram sequences
- Named entity extraction
- Topic models

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



# Document Representation

- A document representation aims to capture what the document is about
- One possible approach:
  - Each row in the table represents a document
  - Attribute describes whether or not a term appears in the document

Example

<https://powcoder.com>

	Terms				
	Camera	Digital	Memory	Pixel	...
Document 1	1	1	0	1	
Document 2	1	1	0	0	
...	...	...	...	...	

# Document Representation using TF

- Term Frequency:
  - Attributes represent the frequency in which a term appears in the document
  - $TF(t, d)$

May impose upper and lower limits on TF because the dimensionality is too high

	Terms				
	Camera	Digital	Memory	Print	...
Document 1	3	2	0	1	
Document 2	0	4	0	3	
...	...	...	...	...	

# Inverse Document Frequency (IDF)

- But a term is mentioned more times in longer documents
- Therefore, use relative frequency (% of document):  
$$\text{IDF}(t) = 1 + \log(\text{Total \# of docs} / \# \text{ docs containing } t)$$

<https://powcoder.com>

	Terms				
	Camera	Digital	Memory	Print	...
Document 1	3	2	1	2	
Document 2	1	1.4	1	3	
...	...	...	...	...	

## Combining TF and IDF

- $TFIDF(t, d) = TF(t, d) * IDF(t)$
- Each row represents a document
- Each column is an attribute (term)
- You can use classifier, clustering etc. on this data

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# N-gram sequences

- “The quick brown fox jumps”
- 2-grams or bigrams:
  - {quick, brown, fox, jumps, quick\_brown, brown-Fox, fox\_jumps}
  - You can see that the number of columns can quickly get out of hand

Assignment: Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Named entity extraction

- Example “Silicon Valley”, “LA Lakers”, “Merage School of Business”

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Topic Models

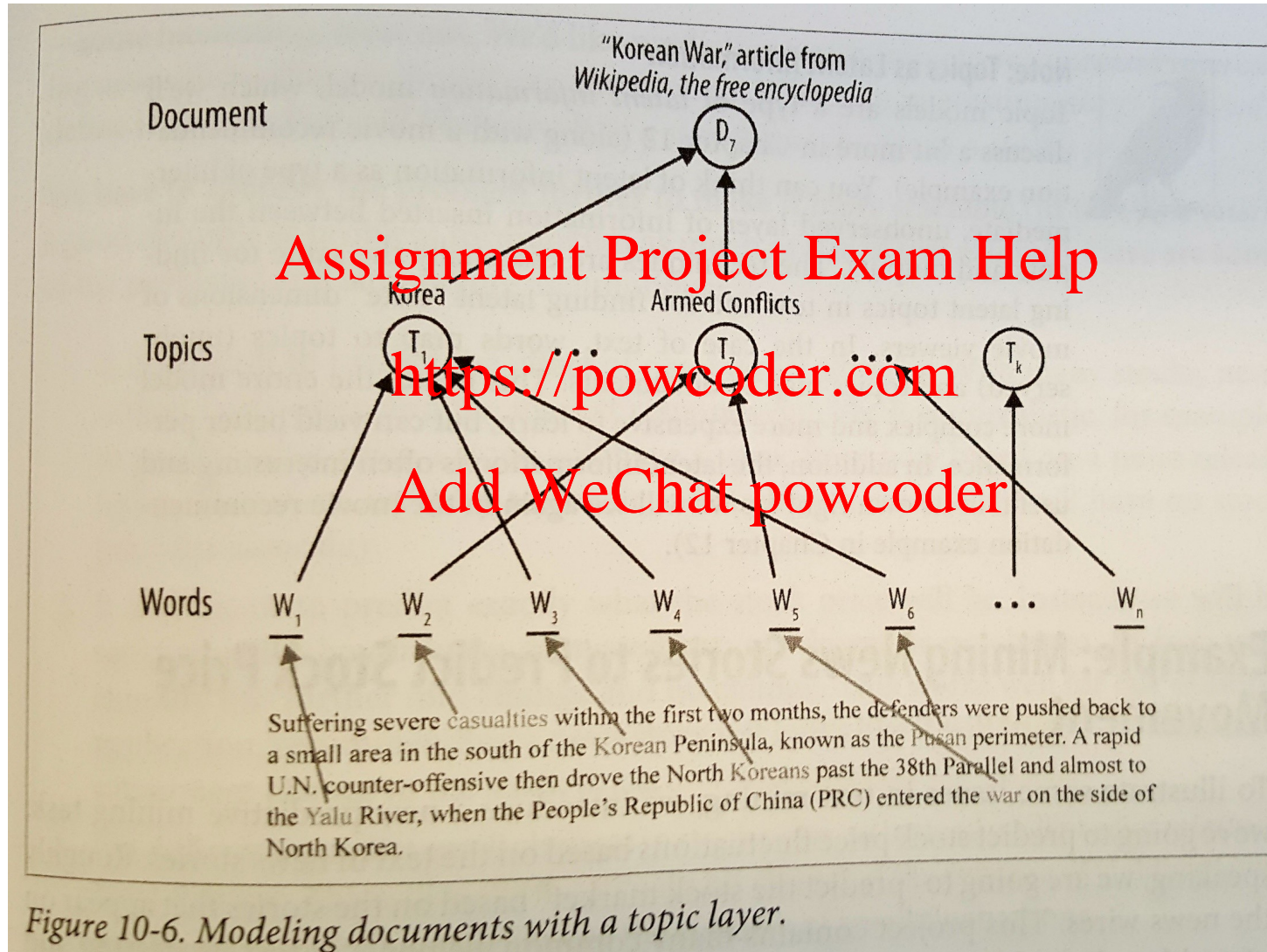


Figure 10-6. Modeling documents with a topic layer.

# Topic Models

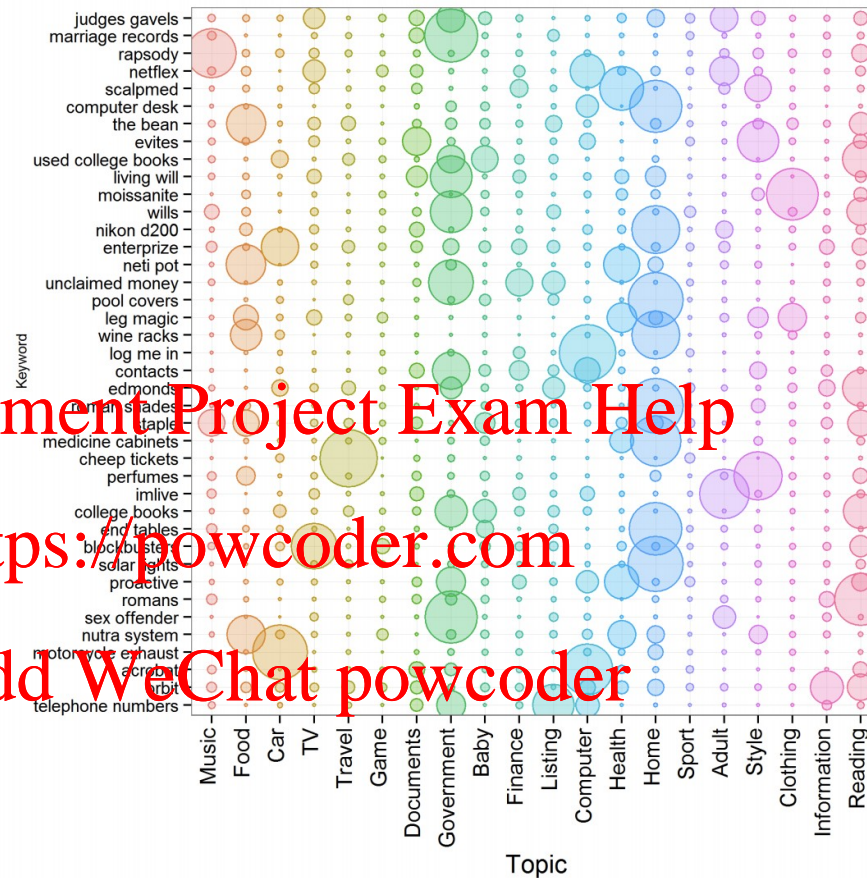


Figure A.1: Topic Distribution of Sample Keywords

Examining the Impact of Keyword Ambiguity on Search Advertising Performance: A Topic Model Approach, Gong, Abhishek and Li (MISQ 2018)



# Text Mining Application 1: Association Rules

After proper representation, data mining techniques can be applied to text, e.g. association rules, clustering, classification.

Keyword-based Association Rules: treat keywords as items.

Microsoft → Antitrust

Assignment Project Exam Help

Document No.	Item 1	Item 2	Item 3	...
100	France	Iraq	US	
101	NASDAQ	NYSE	job	
102	Iraq	US	UK	
103	Microsoft	antitrust	OS	
104	Microsoft	Antitrust	window s	
...				

OR

Doc No.	Microsoft	antitrust	Franc e	...
100	0	0	1	
101	0	0	0	
102	0	0	0	
103	1	1	0	
104	1	1	0	
...				

# Personalized Web Ad Delivery

- Objective:
  - Improve effectiveness of Web ads
  - Customize ad delivery so that ad corresponds to the context user is exploring
- Web content is dynamic → need automated ad placement
  - Example: Gmail
- Solution:
  - Represent each ad as a document with a set of keywords.
  - For example: ad for hybrid car is represented by the following set of keyword: car, electric, environment, etc.
  - Then deliver ads to viewers of pages (i.e., documents) that resemble this description.

# Link Structure Analysis to rank Web pages

- Traditional Information Retrieval methods only examine the appearance of relevant terms, and often fail to account for
  - The quality of the information in the retrieved documents.
  - The reliability of the source
- From the retrieved documents, want to rank authoritative documents higher
- Approach: Mining the Web's link structure to identify authoritative web pages

# Identify Authoritative Web Pages

- The Web includes pages and hyperlinks
- A lot of information is in the structure of web page linkages. **Assignment Project Exam Help**  
Hyperlinks contain rich latent human information  
**<https://powcoder.com>**  
**Add WeChat powcoder**
  - An author creates hyperlink pointing to another page -- can be viewed as endorsement
  - The collective endorsement of a given page by different authors can help discover authoritative pages
- Google uses link structure of the Web to rank documents (PageRank)

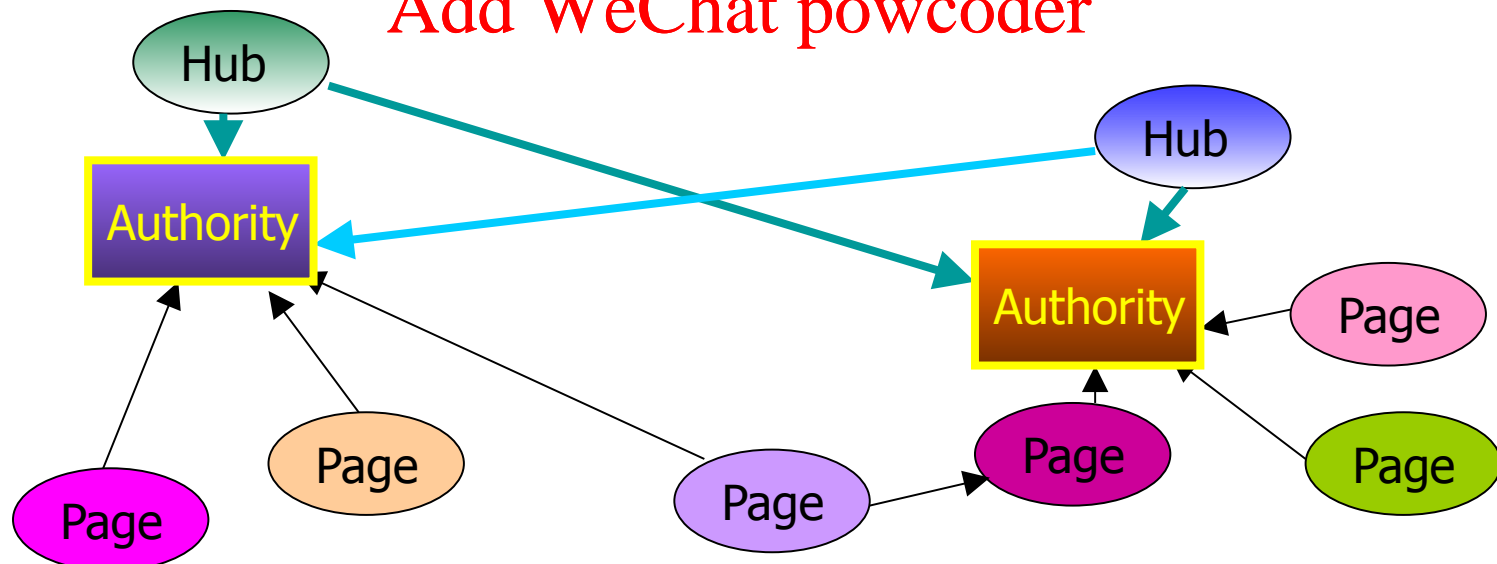
# Using Hubs to identify Authoritative Web Pages

- A hub is a page pointing to many good authorities.
  - E.g., a web page pointing to many good sources of information on business intelligence
- A hub may not be an authority, and have very few links pointing to it.
  - Yet a link from a hub to a page is valued more than a link from a regular page
- An authority is a page pointed to by many good hubs

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



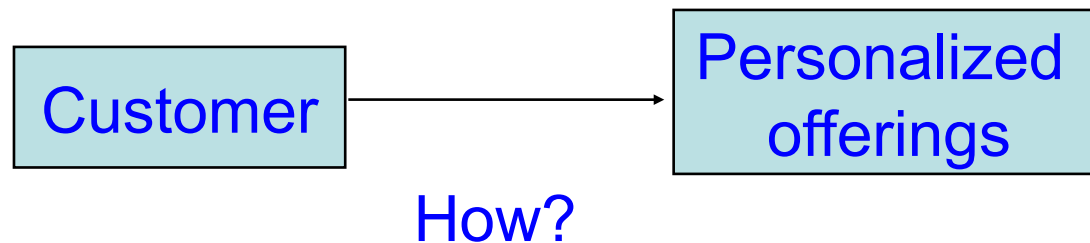
# Personalization

Personalization/customization tailors certain offerings by providers to consumers based on knowledge about them with certain goals in mind.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

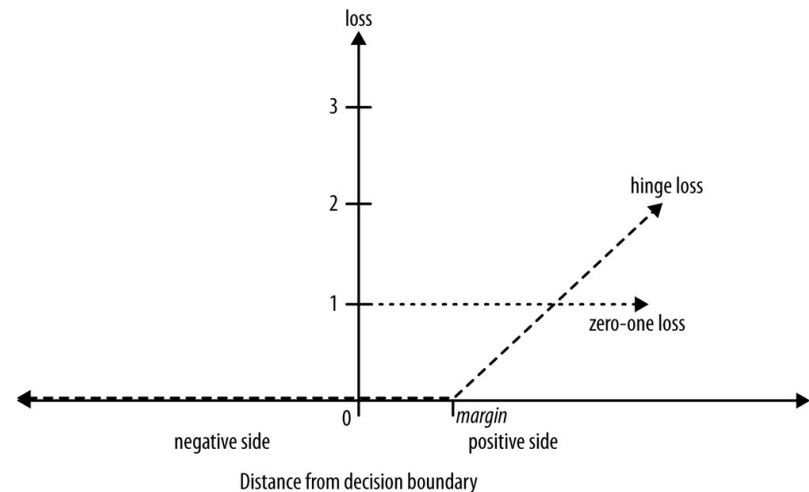
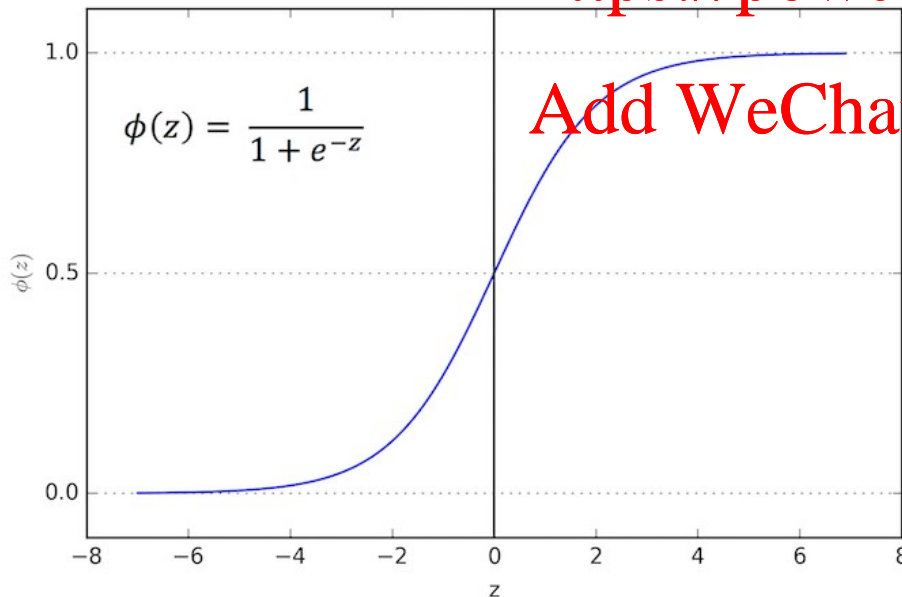


# Classifier: Logistic Regression

- This is not a regression
- Uses logistic function and hinge loss function

<https://powcoder.com>

Add WeChat powcoder



# K Nearest Neighbor (KNN)

K-Nearest Neighbor can be used for classification/prediction tasks.

**Step 1:** Using a chosen distance metric, compute the distance between the new example and all past examples.

**Step 2:** Choose the  $k$  past examples that are closest to the new example.

<https://powcoder.com>

**Step 3:** Work out the predominant class of those  $k$  nearest neighbors - the predominant class is your prediction for the new example. i.e. classification is done by *majority vote* of the  $k$  nearest neighbors. For prediction problem with numeric target variable, the (weighted) average of the  $k$  nearest neighbors is used as the predicted target value.



# How do we determine our neighbors?

- Each example is represented with a set of numerical attributes



John:

Age=35

Income=95K

No. of credit cards=3



Rachel:

Age=41

Income=215K

No. of credit cards=2

Assignment Project Exam Help  
<https://powcoder.com>

- “Closeness” is defined in terms of the *Euclidean* distance between two examples.
  - The Euclidean distance between  $X=(x_1, x_2, x_3, \dots, x_n)$  and  $Y=(y_1, y_2, y_3, \dots, y_n)$  is defined as:

$$D(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

# K-Nearest Neighbor Classifier

## Example : 3-Nearest Neighbors

Customer	Age	Income	No. credit cards	Response
John	35	35K	3	No
Rachel	22	50K	2	Yes
Hannah	63	200K	1	No
Tom	59	170K	1	No
Nellie	25	40K	4	Yes
David	37	50K	2	?

# Collaborative Filtering: Finding like-minded people

- One seeks recommendations about movies, restaurants, books etc. from people with similar tastes
- Automate the process of "word-of-mouth" by which people recommend products or services to one another.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Collaborative Filtering

- Starts with a history of people's personal preferences
- Uses a **distance function** – people who like the same things are “close”
- Determine a neighborhood size (say  $k$  closest data points). We will examine recommendations from this neighborhood only.
  - Typically  $k$  is between 20 and 50
- Uses “**votes**” which are weighted by distances, so close neighbor votes count more

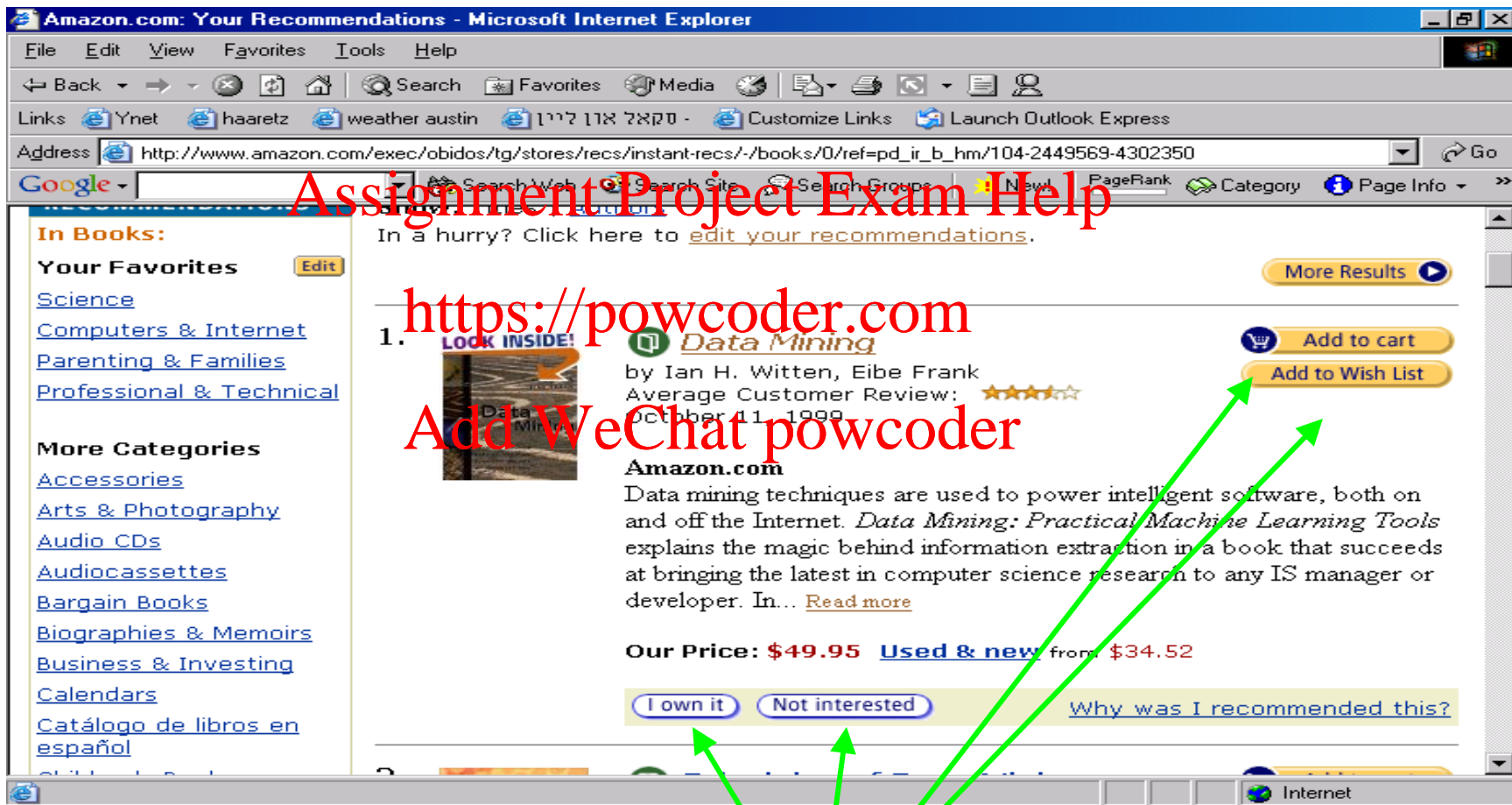
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Example:

amazon.com.



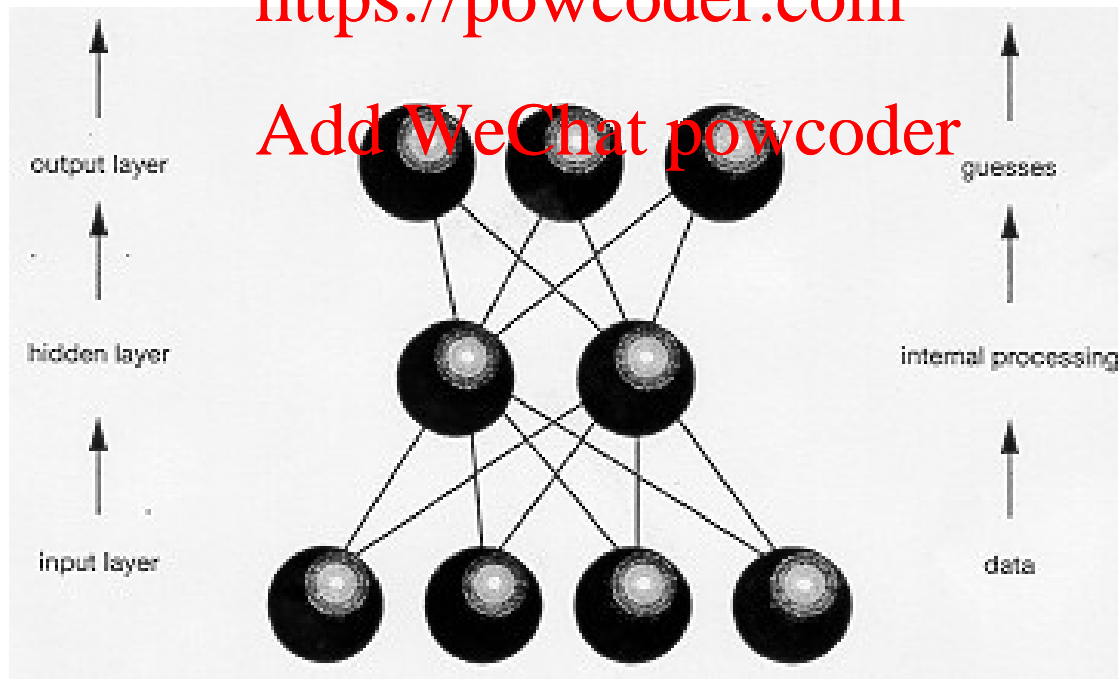
Implicit rating

# Artificial Neural Networks

- An artificial neural network (ANN), usually called neural network (NN), is a mathematical model or computational model that is inspired by the structure and/or functional aspects of biological neural networks. -- Wikipedia
- A neural network consists of an interconnected group of artificial neurons, and it processes information using a connectionist approach to computation. -- Wikipedia
- Neural Nets learn complex functions  $Y=f(X)$  from data.
- ANN can approximate any function (e.g. logistic regression, linear regression).

# Components of Neural Nets

- Neural Nets are composed of
  - Nodes, and
  - Arcs
- Each arc specifies a weight.
- Each node (other than the input nodes) contains a Transfer Function which converts its inputs to outputs. The input to a node is the weighted sum of the inputs from its arcs.



# Recommender Systems

- Collaborative Filtering

Assignment Project Exam Help

- Content Based Recommendation

<https://powcoder.com>

- Use document content to create a description (tags)

- Create user profile with weights for different tags

- Example Books: Genre, Author, Length, Pictures etc.

- Knowledge Based Recommendation

- When we do not have history of purchases (Camera)

- Examine customer needs and match to product features



# Bagging

- Combining predictions by voting/averaging
  - Each model receives equal weight
- “Idealized” version
  - Sample several training sets of size  $n$  (instead of just having one training set of size  $n$ )
  - Build a classifier for each training set
  - Combine the classifiers’ predictions

# Bagging classifiers

## Model generation

```
Let  $n$  be the number of instances in the training data
For each of  $t$  iterations:
    Sample  $n$  instances from training set
    (with replacement)
    Apply learning algorithm to the sample
    Store resulting model
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## Classification

```
For each of the  $t$  models:
    Predict class of instance using model
Return class that is predicted most often
```

# Boosting

- Also uses voting/averaging
- Weights assigned according to performance
- Several variants
  - Read text for AdaBoost

<https://powcoder.com>

Add WeChat powcoder

Link Analysis is used for ...

- A: Identifying similar consumers for product recommendations
- B: Highly non-linear classification
- C: Replicating logistic regression
- D: Determining which web sites or documents are more authoritative and credible.
- E: None of the above

# Next Session

- Project Presentations
  - All Students must attend
  - Please upload slides/files to drop-box on Canvas at least 1 hour before class.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder