

Artificial Intelligence & Data Analytics

Ng Van Duc

L1, Intro

L2, Python

- AI: tools and methods which enable machines to succeed at tasks normally limited to human
- Data Analytics: raw data \Rightarrow insights \Rightarrow $\left\{ \begin{array}{l} \text{actions} \\ \text{decisions} \end{array} \right.$

- Panda.

```
import pandas as pd  
df = pd.read_csv('data.csv')  
df[df.name == "tom"]  
df[df.age > 15]  
df["age"].min()
```

Q & A,

L3, - What are the steps / tasks in Data Preparation?

- What is data exploration? Cleaning? Transformation?

- What do we do in Cleaning Data?

How to deal with Missing Values, Outliers, inconsistencies?

- What do we do in Data Transformation?

What are Normalization, Aggregation, Discretization and their purposes?

- What are the common operation to prepare image data?

- What about Text Data?

What are tokenization, normalization, noise removal here?

What to notice with sensitive data?

L4, - What is the motivation behind the need for Data Integration?

- What is Data Integration, Data Pipeline?

- What are the advantages of Database against conventional file-based system?

- Describe and compare Relational DB and Document-oriented DB.

- Briefly describe SQL.

- What is ETL?

- Describe and compare Data Warehouse & Data Lakes?

L5, - How do you understand feature extraction, selection, dimensionality reduction, feature projection, feature elimination?

- How do you understand the curse of dimensionality?

- Explain Maximum Relevancy Minimum Redundancy!

- Explain briefly PCA, LDA, & compare them.

What is the assumption in PCA?

goal of LDA?

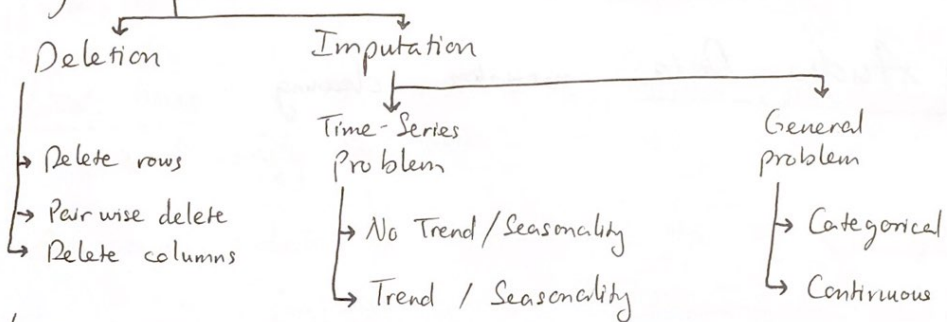
Data Prepagation

$$\text{Data Prepagation} = \text{Exploration} + \text{Cleansing / Transformation} \times \\ = \text{Cleaning} + \text{Transformation} + \text{Integration} + \text{Reduction}$$

⊗ Exploration: Creating initial understanding

⊗ Cleaning: Manual \Rightarrow Automatic

1) - Missing Values



- Outliers:

- Inconsistencies: Ex: Colour: ABB, Fanuc.. \Rightarrow Manually handle
 Quality: Good, Poor, 2, 5.. \Rightarrow Detect manually
 Semantic diff.
 Value diff

⊗ Transformation

- Normalization

- Aggregation: combine >2 attributes into 1 reduce $\left\{ \begin{array}{l} \text{dimension} \\ \text{variability} \end{array} \right.$

- Discretization: ~~Ex~~ Age: $\begin{cases} 1, \dots, 10 \Rightarrow \text{Young} \\ 60 \dots 80 \Rightarrow \text{Old} \end{cases}$ (Int \Rightarrow Word)

convert attributes \Rightarrow 1 discrete value

⊗ Image Data: transform, equalize, segment, augment

- Most common operations: resize, denoise, thresholding, light correction, segmentation ..

⊗ Text Data:

- Web page usually has API
- Operations:
 - tokenization (segment into word strings)
 - normalization (not upper/lower case, "one" vs "1")
 - noise removal (headers, footers, tag, ..)
- Sensitive data:
 - identification (name, bank account, personal info..)
 - anonymization / pseudonymization

⊗ Audio Data: normalize, cleaning

Data integration

- Motivation: variety of data sources (machines, data systems..)
 => need a transparent connection, pipeline..
a consolidated system architecture

- Definition: \rightarrow database integration: combining heterogeneous data from various sources, to enable users a unified access

\rightarrow data pipeline: software system (ex: code program)
takes data (from 1 \rightarrow many source) \Rightarrow transform \rightarrow write to output(s)

- ⊗ Data base: - has database-management system (DBMS) to handle / modify the data

- Advantage: - data is searchable, can be cross-gathered
 - failsafe .. ACID (Atomicity Consistency Isolation Durability)

- Types :	<u>Relational DB</u>	<u>Document-oriented DB</u>
	Tables with columns having relation	Documents, similar to JSON
	Predefined schema	Each doc can be complex, different
	Uses SQL	No SQL
	Not as flexible/scalable	High flexibility

schema = structure

- SQL (Structured Query Language)
 Based on Relational Algebra
 Close to Natural English language

Select * from Table;
 Select Column1 from Name;
 Select * from .. where condition?

Insert into — Values —
 Update — set —

+ Approaches to Data Integration

- ETL (~~E~~xtract, Transform, Load) into target system / data warehouse

- Data Warehouse: a central database system which integrates data from all kinds of company-wide op. data sources for subsequent analysis purpose

- Schema-on-write
Define schema before
any data is written to it



Data Warehouse

vs

Schema-on-read

Data Lakes natural format
(raw ..)



Data Lakes

Data
Schema
Quality
Users
Analytics
Price

schema-on-write



Business Analysts

BI ..



vs

Central repository for storing
data in its natural format

schema-on-read
↓ raw data ..

Scientist, developers

ML ..



= Publish / Subscribe mechanism: to transport the data

Data Source publish ⇒ Cloud ..

Consumers subscribe to topic

- MQTT: (Message Queue Telemetry Transport)

a easy to use publish / subscribe real-time protocol

use a broker to address messages to the right channels
enable real-time push

most promising for IoT

Feature

L5,

we have many different representations of data representing the world

- Extract feature: process of deriving features from an real world object; Should be informative
non-redundant !

Bag of word: feature for texts

- Feature Selection

+ Curse of Dimensionality: the exp growths of feature space ..

⇒ Only pick a subset of the features

+ Benefit: moderate the curse of dimensionality
{
reduction of training time
generalization

feature projection + Dimensionality reduction: \neq is different thing represent ^{well} data in a lower dim space

+ Maximum Relevancy Minimum Redundancy Mutual Info?

$$MI(X_i, C)$$

$$MI(X_i, X_j)$$

$$V_I(S) = \frac{1}{|S|} \sum_{X_i \in S} MI(X_i, C); \quad W_I(S) = \frac{1}{|S|^2} \sum_{X_i, X_j \in S} MI(X_i, X_j)$$

$$S^* = \underset{S \subseteq U}{\operatorname{argmax}} (V_I(S) - W_I(S))$$

$$MI(X, Y) = I(X) + I(Y) - I(X, Y) \geq 0$$

$$I(X, Y) = \sum_{x_1} \sum_{x_2} P(x_1, x_2) \log \frac{P(x_1, x_2)}{P(x_1)P(x_2)}$$

⊗ PCA Dimensionality reduction can be divided into
 { feature elimination
 feature extraction
 (create new variables from old ones)

PCA is a technique for feature extraction

Reality implementation

$N < D$

Gram-Schmidt process

$D, N \gg$ Power method

LDA Problem with PCA. Assumptions: "Dim with highest σ^2 are most imp."
 \Rightarrow LDA

⊗ Compare PCA & LDA

	PCA	LDA
Requirements	Unsupervised	Supervised
Optimization goal	Maximize variance	-
Technique		

⊗ Neural network uses part of its layers to perform feature extraction

⊗ These definitions are .. varying

+ Feature selection \approx feature ~~extraction~~ elimination

However, in some algo., f. elimination refers backward algo., start with complete set of feature, then eliminate till criterion are reached
 f. selection ——— forward ———, ——— empty set
 ———, ——— add features ———

+ F. extraction \approx F. projection

+ PCA is just extraction, the elimination step is not officially in PCA

PCA

$$X_{[D \times N]} = \begin{bmatrix} U_K & \bar{U}_K \\ [D \times K] & [D \times (D-K)] \end{bmatrix}_{[D \times D]} \cdot \begin{bmatrix} Z_{[K \times N]} \\ Y_{[(D-K) \times N]} \end{bmatrix}_{[D \times N]}$$

$$X = U_K \cdot Z + \bar{U}_K \cdot Y \Leftrightarrow \begin{cases} Z = U_K^T \cdot X \\ Y = \bar{U}_K^T \cdot X \end{cases}$$

We want to replace $Y \approx b \cdot \underline{1}_{(1 \times N)}^T$

$$\Rightarrow b = \underset{b}{\operatorname{argmin}} \|Y - b \cdot \underline{1}\|_F^2 = \underset{b}{\operatorname{argmin}} \|\bar{U}_K^T \cdot X - b \cdot \underline{1}\|_F^2$$

$$\text{Set } \frac{\partial f(b)}{\partial b} = 0 \Rightarrow \text{Find } b = \bar{U}_K^T \cdot \bar{X} \quad \text{with } \underline{\bar{X}} \text{ is expectation of } X$$

$$\Rightarrow X \approx \tilde{X} = U_K \cdot Z + \bar{U}_K \bar{U}_K^T \bar{X} \cdot \underline{1}^T$$

$$\text{Loss func: } J = \frac{1}{N} \|X - \tilde{X}\|_F^2 = \frac{1}{N} \|\bar{U}_K \bar{U}_K^T X - \bar{U}_K \bar{U}_K^T \bar{X} \cdot \underline{1}^T\|_F^2$$

If U is orthogonal:

$$\Rightarrow J = \dots \sum_{i=K+1}^D u_i^T S u_i$$

$$L = \sum_{i=1}^D u_i^T S u_i = \sum_i \lambda_i$$

$\Rightarrow L$ doesn't depend on how

U is chosen

$$\Rightarrow J_{\min} \Leftrightarrow F_{\max} = L - J = \sum_{i=1}^K u_i^T S u_i$$

\Rightarrow Choose K maximum eigen values λ_i

PCA algorithm step

1, Calculate expectation/mean: $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

2, Translate to 0: $\hat{x}_n = x_n - \bar{x}$

3, Covariance matrix: $S = \frac{1}{N} \hat{X} \cdot \hat{X}^T$

4, Find λ_i of S , rank/sort from greatest to smallest

5, Find respective eigenvector & normalize $\Rightarrow U_k$ each column is an eigenvector

6, $Z = U_k^T \cdot \hat{X}$

7, $x \approx U_k \cdot Z + \bar{x}$

⊕ Assumption: Dimension with the highest σ is the most important which might not always be true

\Rightarrow LDA

LDA

LDA algorithm steps:

1) Calculate mean μ_i of each class c_i , and of all data

$$\mu_{c_i} = \frac{1}{N_{c_i}} \sum x_{c_i} \quad ; \quad \mu = \frac{1}{N} \sum x_i$$

$[0 \times 1]$

2) Scatter matrix

+ Within-class scatter matrix: $S_W = \sum_i S_i$

$$S_i = \sum_{x \in G_i} (x - \mu_{c_i})(x - \mu_{c_i})^T$$

Or calculate covariance matrix: $\Sigma_{c_i} = \frac{1}{N_{c_i} - 1} \sum_{x \in G_i} (x - \mu_{c_i})(x - \mu_{c_i})^T$

$$\text{then } S_W = (N_i - 1) \cdot \Sigma_i$$

+ Between class

$$S_B = \sum_{i=1}^c N_i (\mu_{c_i} - \mu)(\mu_{c_i} - \mu)^T$$

3) Find eigenvalues/eigenvectors of $S = S_W^{-1} \cdot S_B$

4) Carry on as PCA

Supervised Learning

L6,

⊗ Definition

AI: Any technique which enables a computer to mimic human behavior

General AI (GAI): Transfer knowledge across domains
 ↳ system that can do many tasks

Narrow AI: Perform a single task extremely well

↓ Machine Learning: Algorithms whose performance improve
 as they are exposed to more data

↓ Supervised learning: **labeled data**

↳ Unsupervised learning

↳ Re-inforcement learning

— Regression V/S Classification
 value ← output → category

Linear Regression

Logistic Regression

Decision Tree Regressor

K-nearest neighbors, SVM

Decision Tree

Perceptron ..

CNN ..

— Confusion matrix & metrics: $ACC = \frac{TP + TN}{P + N}$
 $F_1 = \frac{2TP}{2TP + FP + FN}$

— Gradient descent, overfitting, cross validation..

Unsupervised Learning

L7,

Doesn't have label \Rightarrow Model discover relations

1) Finding clusters:

+ Clustering:

Hierarchical clustering

Single-Linkage
Wards method

Partitioning clustering

K-means
Fuzzy C-means
Affinity Propagation
EM-Clustering

Density-based clustering

Mean-shift
DBSCAN
OPTICS

+ Anomaly Detection: objects outside of clusters

2) Association Rule Mining: discover associations (relationship, dependency) between variables

- Support = $\frac{\text{Antecedent} \cap \text{Consequent}}{\text{Total}}$

- Confidence = $\frac{\text{Support}}{\text{Coverage}}$

$\begin{matrix} \textcircled{A} & \textcircled{C} \\ \text{---} & \text{---} \\ \text{A} & \text{C} \end{matrix}$
% confidence

$\begin{matrix} \textcircled{A} \\ \text{---} \\ \text{A} \end{matrix}$
100% Confidence

- from apyori import apriori
transactions = [...]

association-rules = list(apriori(transactions))

3) Dimensionality Reduction:

Feature Selection

Feature filtering
Random forest (Gini...)

Linear Projections

PCA SVD
ICA

Non-linear projections

Transform feature space
while maintaining local / global distance metrics

Reinforcement Learning

28,

- Main idea: goal - feedback : reward
learn from (state, action, reward)

+ MDP (Markov Decision Process): $\langle S, A, P, R, \gamma \rangle$

- Markov Property: $P[S_{t+1} | S_t] = P[S_{t+1} | S_1 \dots S_t]$
- $G_t = R_{t+1} + \gamma \cdot G_{t+1}$
 $V_\pi(S_t) = \mathbb{E}_\pi(G_t | S_t)$; $Q_\pi(s_t, a) = \mathbb{E}_\pi(G_t | S_t, a)$
- Bellman equations: $V^*(s, a) = \max_a Q^*(s, a)$
 $Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \cdot V^*(s')]$
 $V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \cdot V^*(s')]$

+ Classification:

— Offline vs Online
Dynamic Programming
(Value iteration: synchronous / asynchronous)

— Model-based vs Model-free
Learn $P(s' | a, s)$
 $R(s, a)$

on-policy online | actor critic
| Q learning

Learn $V(s)$, $Q(s, a)$

— Exploration vs Exploitation
 $\Rightarrow \epsilon$ -greedy

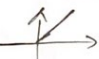
Neural Networks


29,

+ Some functions:

- Sigmoid: $f(s) = \frac{1}{1+e^{-s}}$ \Rightarrow Logistics Regression (2 classes)

- Tanh: $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

- ReLU: 

- Leaky ReLU: 

- One-hot coding

\Downarrow
Softmax Regression
$$a_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

+ Loss functions:

- L2, L4

- Cross entropy: $L = - \sum_i t_i \log(y_i)$

- Hinge loss: $L = \sum_n [1 - t_n \cdot y(x_n)]$

Visual Analytics

L11,

- To interpret data in an efficient way so that it can be understood, used, learned...
- Visualization: a way to communicate concrete, abstract ideas
- Visual Analytics: creation of tools, techniques enable people
 - synthesize info, insight from massive, ambiguous data
 - detect the expected & discover the unexpected
 - provide timely assessment ⇒ communicate ⇒ action

(Statistic Graphs Data mining) Data / Machine + Human (Cognitive perception visual intelligence.)

- Preprocess ⇒ apply algorithmic analysis ⇒ Visualize
⇒ generate insights ⇒ new hypotheses ⇒ updated visualization
- There are available tools:
Tableau, Superset, Qlik, Power BI
- Incomplete problem formulation ⇒ the need for visualizations

1) EDA & visual analytics

Exploratory data analysis

visual representation // of raw data
underlying info

Ex: plot it .. to see if we got the right predict
find outliers?

+ from `pivottablejs` import `pivot-ui`
import `qgrid`..

⊗ Challenges:

SQL UI/PE

Scalability with Data Volumes & Data Dimensionality
Quality of data & graphical representation
Visual representation & Level of Detail
User Interfaces, Interaction Styles & Metaphors
Display Devices
Evaluation & Infrastructure

2) Modeling & visual analytics

visual representation // of model
underlying info and prediction

— Ex: Why given an image, a network classify as a "dog"

+ Target audience:

- Domain experts (doctors, insurance agents...) trust?
- Affected users (...) ⇒ understand, verify
- Regulatory entities (gov..) ⇒ certify, audits..
- Data scientist, developer ⇒ efficiency, research
- Managers, EB ⇒ assess, understand

different audience
different need

⇒ XAI (Explainable AI): ethics, fairness, privacy, trust, transparency
STEPFA

Ex: Activation Atlas

Tensorflow playground