

L1

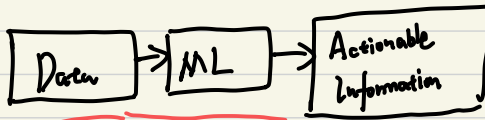
CHINMAY HEGDE (professor)

Statistics
algorithm

ML

optimization
Application

What is ML:



(best used when we don't know the rule of things)

→ Data representation

→ Measure of goodness "loss function"

→ Method to optimize for measure "training algorithm"

Today: Data Representation

Data is a list of attributes that is collected about an object/phenomenon is interest.

Eg. weather data

wind speed, temperature, pressure, humidity, air quality, etc.
 w t p h a

[$w(1)$, $t(1)$, $p(1)$, $h(1)$, $a(1)$]

[$w(2)$, $t(2)$, $p(2)$, $h(2)$, $a(2)$]

⋮

d attributes \longrightarrow tuple of size d
 \longrightarrow vector in d-dimensional space

Vector space:

① Collection of vectors which satisfy 2 properties:

a) addition

$$x = (x_1, \dots, x_d)$$

$$y = (y_1, \dots, y_d)$$

$$x+y = (x_1+y_1, \dots, x_d+y_d)$$

b) Scalar multiplication ~~scalar~~

$$x = (x_1, \dots, x_d)$$

$$\alpha x = (\alpha x_1, \dots, \alpha x_d)$$

Examples:

1) weather resolution \times RGB 个数值变量

2) image $2048 \times 1536 \times 3$ ($=d$)
 $\rightarrow p^d$

3) stocks $d = \text{year}$ p^{year}

② Properties of vector space

1) Dot product / inner products

$$x, y \in p^d$$

$$\langle x, y \rangle = x_1 y_1 + x_2 y_2 + \dots + x_d y_d = \sum_{i=1}^d x_i y_i$$

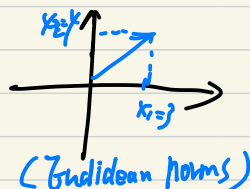
2) Cartesian products / outer products 所有组合可能

$$x \otimes y = \begin{pmatrix} x_1 y_1 & x_2 y_1 & \dots & x_d y_1 \\ x_1 y_2 & x_2 y_2 & \dots & x_d y_2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1 y_d & x_2 y_d & \dots & x_d y_d \end{pmatrix}$$

3) Norms

$$x = (3, 4) = (x_1, x_2)$$

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2} = \sqrt{3^2 + 4^2} = 5$$



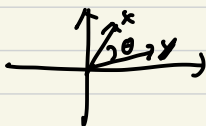
$$\|x\|_1 = |x_1| + |x_2| = 3 + 4 = 7$$

(L1 norms)

(Lp norms)

4) similarity (how we measure similarity?)

$$\frac{\langle x, y \rangle}{\|x\|_2 \cdot \|y\|_2} = \cos \theta$$



Application: NLP (natural language processing)

problem: Given a dataset of n documents

$\{D_1, D_2, \dots, D_n\}$ and a query document D^* , find the closest document to D^* in the dataset.

Solution: Step 1 $d = 80,000$ (# of words in English)

$D_i \rightarrow (x_i) \in \mathbb{R}^d$ (d 维向量) \rightarrow 各种字频数向量

$x_{i,j}$ = # times word j appears in document i
 $D \rightarrow \{x_1, x_2, \dots, x_n\}$ (documents)

Step 2 Define cosine similarity

$$\cos \theta_i = \frac{\langle x_i, x^* \rangle}{\|x_i\|_2 \cdot \|x^*\|_2}$$

\cos 越大, 越相似

Step 3 $i^* = \arg\max \cos \theta_i$

$\cos \theta_i = 1$ - 模一样

["nearest neighbour search"]

夹角为 0