

# Bra-att-ha-algoritmer...

Identifera grupper av objekt  
Jämföra objekt  
Mönsterigenkänning  
Åskådliggöra fenomen



# Informationssökning

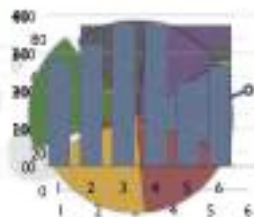
Informationssökning

- Traditioner avse
- Merbaserad indexering (Webb Dokument avse) (glömmer)
- Ordtekniska modeller
  - Vektorer
    - Global (DICE, IDF, Inverse Weight, Word Length, DF)
    - Lokal (Term Frequency (abokännetecken:  $\ln(\text{TF}_i / \text{doc}_i)$ )
  - Värdebaserad sökning
- Hög precision hos informationssökning
  - Recall - andelen av alla relevanta dokument som hittats
  - Precision - andelen av de hittade dokumenten som faktiskt är relevanta
- Problem som ställer sig där
  - Synonymer - Hittas även på andra begrepp (t.ex. Källa, Stäm)
  - Pluraler "resulterade"
- Polysem - flera olika begrepp med samma namn (Blått Ford vs. Forders Ford)

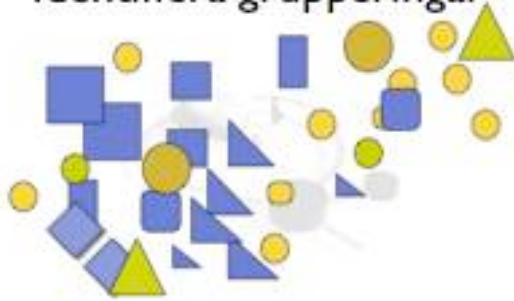


# Åskådliggöra fenomen

	1	2	3	4
1	10	5	20	10
2	5	1	10	10
3	5	5	10	10
4	5	1	10	10
5	10	5	10	10
6	5	5	10	10



## Identifiera grupperingar



## K-means

Identifiera grupperingar

- Dela in datat i ett antal kluster/klasser
- Iterativ algoritim, som jobbar med att uppfylla följande två kriterier:
  1. Varje klass har ett "tyngdpunkt" som är medelpositionen för alla punkter i den klassen.
  2. Varje punkt tillhör den klass vars tyngdpunkt som den är närmast.

## K-means

Identifiera grupperingar

- Initialize
- Loop until termination condition is met:
  1. For each object, assign that object to a class such that the distance from this object to the center of that class is minimized.
  2. For each class, recalculate the means of the class based on the objects that belong to that class.
- End loop;

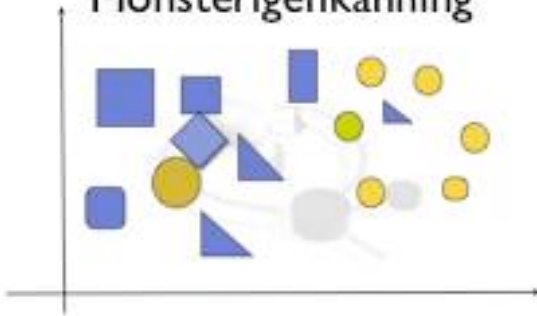
# K-means

Identifiera grupperingar

- Några kritiska punkter i ansatsen
  - Hur många klasser ska man söka efter
  - Hantering av "döda klasser"
  - Initieringen
  - Likhet/avståndsmåttet
  - Termineringsvillkoret
    - Algoritmen terminerar, men kan ta tid...
    - Alternativa termineringsvillkor
      - max ett visst antal punkter som byter klass
      - Kör ett visst antal steg



# Mönsterigenkänning



# KNN

Mönsterigenkänning

- Metod för att avgöra klasstillhörighet
- Hitta dom K närmsta grannarna i träningsmängden
- Undersöka de funna grannarna för att avgöra klasstillhörighet
  - Majoritetsbeslut
  - ???

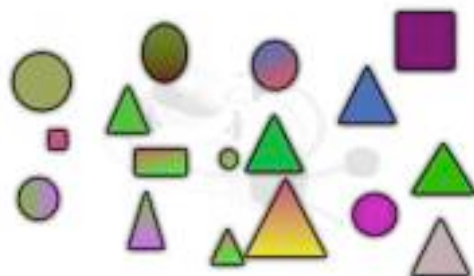


# KNN

- Några kritiska punkter med KNN
  - Likhetsmått
  - Proceduren för att utvärdera grannarna
  - Antalet grannar
  - Storleken på träningsmängden



# Reducering av egenskapsrymden



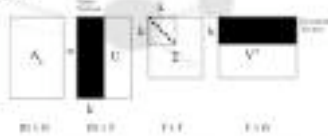
# LSA/LSI

- LSI utvecklades slutet av 80-talet början av 90-talet av Susan Dumais & Michael Berry
- Patenterad metod Bell Communications Research (Bellcore)
  - Computer information retrieval using latent semantic structure U.S. Patent No. 4,839,853, June 13, 1989.
- Samma grundprincip som traditionell approach
  - Matrisbaserad (ordfrekvensmatris)
  - Vektorbaserad sökning
- Stora skillnaden är indexeringsstrukturen



## LSA/LSI

- K-dimensionellt semantiskt rum konstrueras där termer och dokument placeras in i samma sökrum
- Jämföra och hitta klasser av dokument
- LSI använder en trunkerad SVD för att konstruera indexstrukturen



## LSA/LSI

- För informationsöppning är första steget att representera frågan i samma indexstruktur, dvs det k-dimensionella semantiska rummet
- Användan tillser ord, ordlistor, dokument eller dokumentanläggningar att vara frågor
- Frågor är en samling av ord
- En fråga q representeras som en termvektor eller ett pseudo-dokument

$$V^T = \Sigma^{-1} U^T A = V^T_{:,j} = \Sigma^{-1} U^T A_{:,j}$$

$$V_{:,j} = (A_{:,j})^T U \Sigma^{-1}$$

$$\hat{q} = q^T U_k \Sigma_k^{-1}$$

## LSA/LSI

- Vidare, detta pseudo-dokument jämförs mot alla dokument som ingår i indexstrukturen, och de dokument som uppfyller ett likhets/närhetskriterie returneras som ett söksvar



## LSA/LSI

- Förändra Indexstrukturen
  - Beräknas om hela (ingen uppdatering)
    - Reservskrivande
    - Bibehåller "korrektheten" i strukturen
  - Uppdatering är en process där man förändrar en existerande LSI-genererad indexstruktur
    - Inkorporera "pseudo dokument" och "pseudo termer"
    - Snabb och enkel
    - Garanterar inte ortogonalitet i det sekundärskanummet
  - SVD-uppdatering (G.W. O'Brien (1996))
    - Addera nya dokument & nya termer
    - Korrigera för förändringar i "termviktarna"
    - Bibehåller "korrektheten" i strukturen
    - Snabbare än att räkna om hela

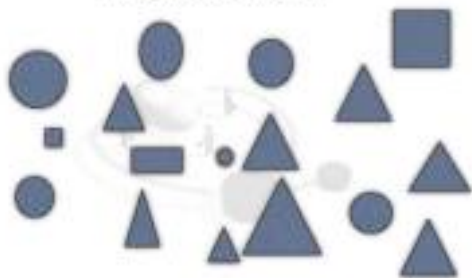


## LSA/LSI

- Kritiska punkter för LSI
  - Kräver mycket beräkningsresurser
  - Minne
  - Beräkningskraft
- Välja dimensionen på indexstrukturen ( $k$ )
- Termviktning mm
- Mätning av närhet



## Likhetsmått



## Olika sätt att jämföra

- Vinkeln mellan vektorer
- Avstånd mellan objekt
  - Ex. Euclidiskt avstånd
- Överlappning
- Area
- Form
- 



## Lite grann om komplexitet

- Rumskomplexitet
  - Term<sup>n</sup>dokument matriser
- Tidskomplexitet
  - LSA/LSI =  $O(n^2)$
  - K-means =  $O(ndcT)$
  - KNN =  $O(n^2m)$