

# Datalogi intro

Kristoffer Klokke

2021

# Contents

<b>1</b>	<b>Repræsentation af tal</b>	<b>4</b>
1.1	Konvertering fra 2 tals system til 10 tals system . . . . .	4
1.2	Konvertering fra 3 tals system til 10 tals system . . . . .	4
1.3	Konvertering fra hex tals system til 10 tals system . . . . .	4
1.4	Konvertering fra hex tals system til 2 tals system . . . . .	5
1.5	Konvertering fra 10 tals system till 2 tals system via algoritme	5
1.6	Two's complement . . . . .	5
1.7	Kommatal i binær . . . . .	6
1.7.1	Fast decimalpunkt . . . . .	6
1.7.2	Flydende decimalpunkt . . . . .	6
<b>2</b>	<b>Gates og memory</b>	<b>7</b>
2.1	Gates . . . . .	7
2.2	Memory . . . . .	8
<b>3</b>	<b>CPU</b>	<b>9</b>
3.1	Opbygning . . . . .	9
<b>4</b>	<b>Algoritmer</b>	<b>9</b>
4.1	Binary search . . . . .	9
4.2	Køretid . . . . .	10
4.3	Invariant . . . . .	10
<b>5</b>	<b>Merging and hashing</b>	<b>10</b>
5.1	Datastrukture . . . . .	10
5.1.1	Sekventiel tilgang . . . . .	11
5.1.2	Random access . . . . .	11
5.1.3	Merging . . . . .	11
5.2	Hashing . . . . .	11
5.2.1	Sandsynlighed for overlap . . . . .	12
<b>6</b>	<b>Machine Learning</b>	<b>12</b>
6.1	Forms of machine learning . . . . .	12
6.2	Supervised Learning . . . . .	12
6.2.1	Regression problem . . . . .	12
6.2.2	Classification problem . . . . .	13
6.3	General framework . . . . .	13
6.4	Overfitting . . . . .	13
6.5	Training and Assessment . . . . .	13
6.6	McCulloch-pitts "units" . . . . .	14

6.7	Network structures . . . . .	14
<b>7</b>	<b>Feature spaces for representation of images</b>	<b>15</b>
7.1	Color histogram . . . . .	15
7.2	Distance of color histograms . . . . .	15
<b>8</b>	<b>Clustering</b>	<b>16</b>
8.1	Lloyd/forgy . . . . .	16
8.2	MacQueen/ $k$ -means . . . . .	16
<b>9</b>	<b>Models of Computation, languages and Recursion</b>	<b>17</b>
9.1	Models of Computation . . . . .	17
9.2	Deterministic Finite Automaton (DFA) . . . . .	17
9.3	Context-Free Grammar . . . . .	18
<b>10</b>	<b>Online algorithms</b>	<b>18</b>
<b>11</b>	<b>Satisfiability</b>	<b>19</b>
<b>12</b>	<b>Databases</b>	<b>21</b>
12.1	Conceptual . . . . .	21
12.2	Logical . . . . .	22
12.3	Physical . . . . .	22
12.3.1	Relational Algebra . . . . .	22
12.4	Integrity Constraints . . . . .	22
12.5	Primary key . . . . .	22
<b>13</b>	<b>Cryptography</b>	<b>22</b>
13.1	Encryption . . . . .	23
13.2	Symmetric key systems . . . . .	23
13.3	Asymmetric keys . . . . .	24
13.4	RSA . . . . .	24
13.4.1	Why it works . . . . .	24
13.4.2	Why RSA is secure . . . . .	25
13.4.3	RSA implementation . . . . .	26
13.5	Symmetric and public key system . . . . .	26
<b>14</b>	<b>Graph theory</b>	<b>27</b>
14.1	Matrix representation . . . . .	27

# 1 Repræsentation af tal

Bit - 0 / 1

Byte = 8 bits - 0 - 255

## 1.1 Konvetering fra 2 tals system til 10 tals system

$$1011_2 \rightarrow 11_{10}$$

1	0	1	1
$1 \cdot 2^3$	$0 \cdot 2^2$	$1 \cdot 2^1$	$1 \cdot 2^0$
8	0	2	1
$= 11$			

## 1.2 Konvetering fra 3 tals system til 10 tals system

$$1202_3 \rightarrow 47_{10}$$

1	2	0	2
$1 \cdot 3^3$	$2 \cdot 3^2$	$0 \cdot 3^1$	$2 \cdot 3^0$
27	18	0	2
$= 47$			

## 1.3 Konvetering fra hex tals system til 10 tals system

Grund tal for hex:

0123456789ABCDEF

$$A32B_3 \rightarrow 47_{10}$$

A	3	2	B
$10 \cdot 16^3$	$3 \cdot 16^2$	$2 \cdot 16^1$	$11 \cdot 16^0$
40960	768	32	11
$= 41771$			

## 1.4 Konvetering fra hex tals system til 2 tals system

Da hex = 4 bits kan hver samling af 4 bits direkte oversættes til hex  $2F7E_3 \rightarrow 0010111101111110_2$

2	<i>F</i>	7	<i>E</i>
2	15	7	14
0010	1111	0111	1110
= 0010111101111110			

## 1.5 Konvetering fra 10 tals system till 2 tals system via algoritme

Følgende algoritme finder cifrene fra højre til venstre i den binære representation af et positivt heltal  $N$

$X \leftarrow N$

While  $X > 0$

    Next digit = rest ved heltalsdivision  $X:2$

$X =$  kvotient ved heltalsdivision  $X:2$

Eksempel:

Heltalsdivision	Kvotient	Rest	
25:2	12	1	$25 = 11001_2$
12:2	6	0	
6:2	3	0	
3:2	1	1	
1:2	0	1	

Hvis man rammer 0 skal der tilføjes 1 og husk at vende den rigtige retning!

## 1.6 Two's complement

Bruges til at beskrive positive og negative tal.

Det første bit vil negativ og dermed hvis et binært tal starter med 1 vil det være negativt i denne repræsentationsform.

For at skifte fortegn kan alle bits flippes og 1 adderes. Eks.

$$0110_2 \rightarrow 6_{10}$$

$$1001_2 + 1 = 1010 \rightarrow -6_{10}$$

## 1.7 Kommatal i binær

### 1.7.1 Fast decimalpunkt

Den binær eksponent vil blive negativ efter komma med ens regne metode som normalt.

Eks.  $1.101_2 \rightarrow 1.875_{10}$

1.	1	0	1
$1 \cdot 2^0$	$1 \cdot 2^{-1}$	$0 \cdot 2^{-2}$	$1 \cdot 2^{-3}$
1	$\frac{1}{2}$	0	$\frac{1}{8}$
$= 1\frac{5}{8} = 1.625$			

### 1.7.2 Flydende decimalpunkt

For det flydende decimalpunkt  $-9.87 \cdot 10^{-2}$

Er eksponenten -2 og mantissen er 9.87 og den har et negativ fortegn.

For 8 bits er der afsat følgende bits til at beskrive et flydende decima ud fra pensum: 1, 3 og 4

Første bit er til fortegnet hvor 1 er negativt og 0 er positivt.

De næste 3 er eksponent i binært

De sidste 4 er til mantissen, hvor den fyldes op med 0, hvis mantissen ikke udfylder alle 4. Her vil det sidste bit i tallet også undlades da det ikke er nødvendigt og da det altid vil starte med 1.

Da i det givende eksempel mantissen er for stor til 4 bits kan half-precision bruges som tager udgangspunkt i 16 bit med fordelingen: 1,5 og 10.

Dermed bliver  $-9.87 \cdot 10^{-2}$  til 1|10010|1110110110

I tilfælde hvor mantissen ikke fyldte alle bits ville der bare fyldes 0 på i slutningen, da det blot giver en større præcision og ikke påvirker tallet.

Dette kan også bruges fra fast decimal i binært til flydende decimal, eksempelvis  $-0.01101_2$

Her er fortegnet negativt og dermed 1

Eksponenten vil laves til -2 dermed 110 i two's complement for at det bliver

negativt

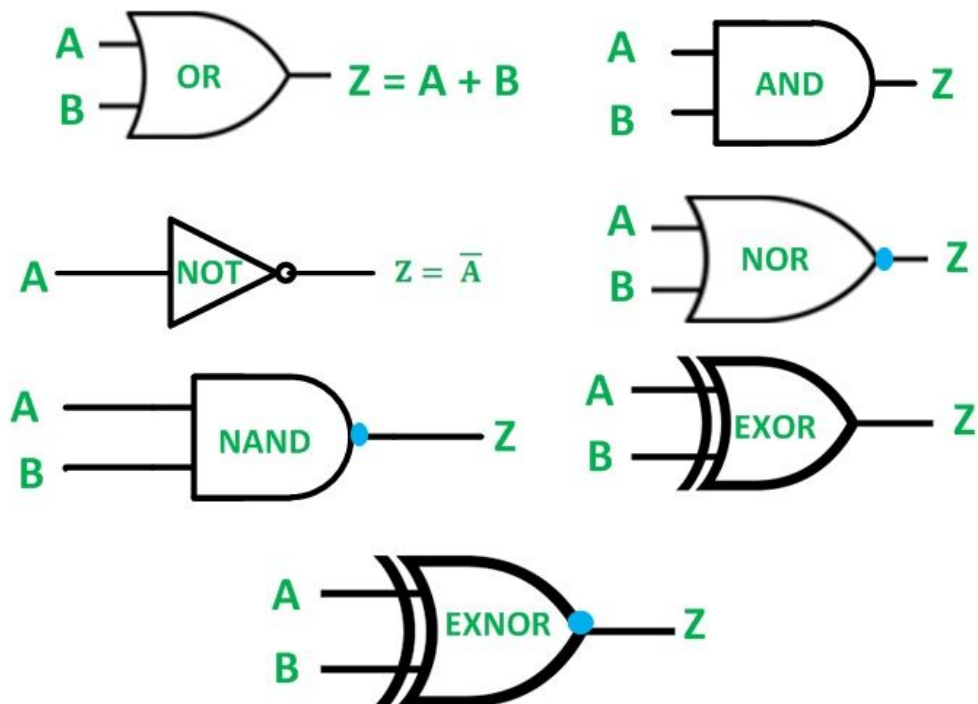
Mantissen er 1101, men da det første cifre altid vil være 1 undlades den

Dermed bliver det 1|110|1010

## 2 Gates og memory

### 2.1 Gates

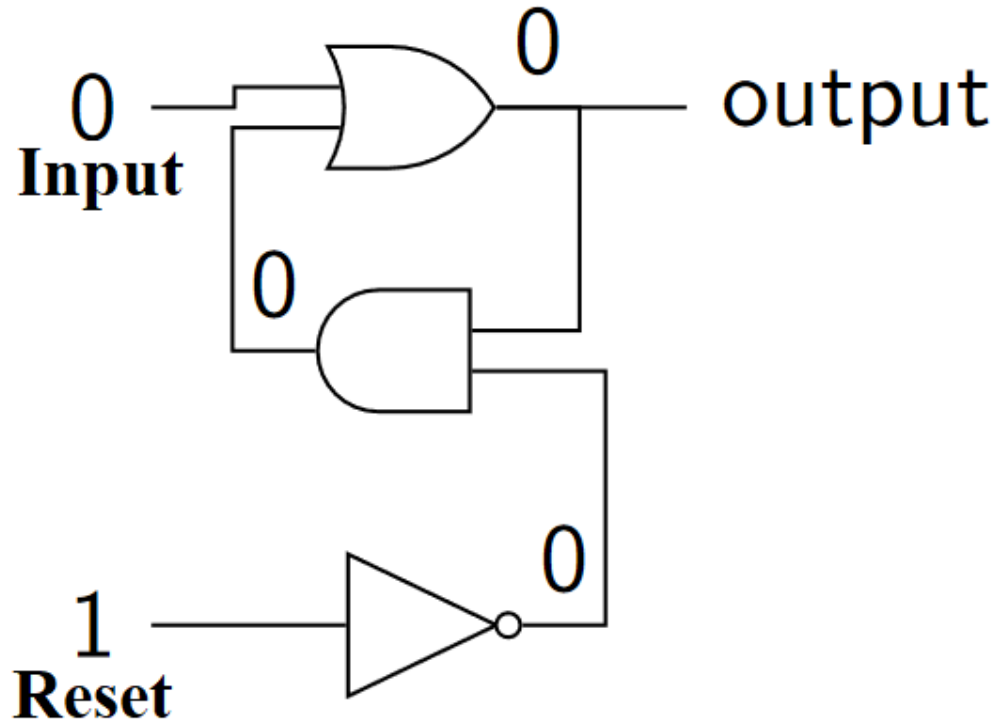
#### Various Logic Gates



Circuit Globe

## 2.2 Memory

A memory can be made with gates, called flip-flop

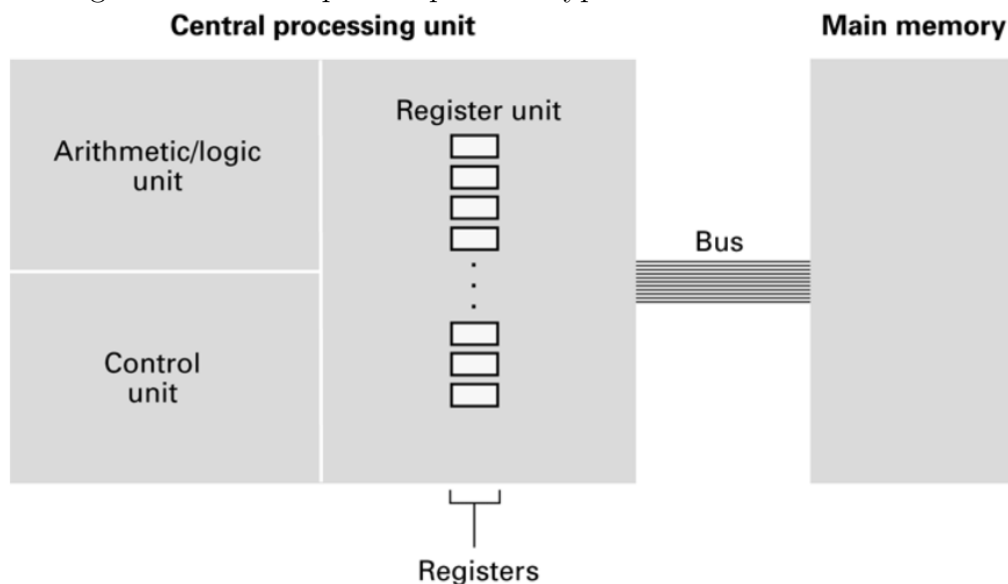




## 3 CPU

### 3.1 Opbygning

ordlængde + antal bits pr manipulation typisk 32 eller 64.



Kommando typer - flytning og manipulation af memory og ændring af program flow

CPU cycle - Fetch (get command), decode (interpretate command), execute

Program flow - styret af counter (nuværende command) og instruction pointer (næste command) som er en register

En commands opbygning er først en opcode, som er et antal bits som repræsenterer en command, herefter er de resterende bits argumenter.

## 4 Algoritmer

An algorithm is an ordered set of unambiguas executable steps that define a terminatiry process.

Ofte beskrevet med pseudokode

### 4.1 Binary search

An algorithn which search thorough an array.

```

BinarySearch(L, x)
l:=1, r:=L.lenght, m:=(l+r)/2
while l <= r and L[m]!=x
    if x < L[m]
        r:= m-1
    else
        l:=m+1
m:=(l+r)/2
if l <= r
    Return m
else
    Return "Not found"

```

Efter 1 check vil der være  $\leq \frac{n}{2^i}$  elementer tilbage. Værste tilfældet vil der være 1 tilbage i listen.

$\frac{n}{2^i} < i \iff n < 2^i \iff \log_2(n) < i$   
 Dermed er køretiden  $O(\log n)$

## 4.2 Køretid

Køretid findes ud fra den karakteriserende operation som er propationel med køretiden.

Til at beskrive køretid bruges  $O(n)$ ,  $O(n^2)$  og lign. hvor at den enkelte algoritmes køretid vil beksrives som  $T(n)$ ,  $T(n^2)$

Her bruges  $T$ , da den er defineret som den tilsvarende  $O$  multipliceret med en konstant.

$T(n) \in O(f(n)) \iff \exists k, m \in \mathbb{Z} : \forall n \in \mathbb{Z}, n \geq m : T(n) \leq k \cdot f(n)$

## 4.3 Invariant

Invariant er en måde hvorpå man kan skabe et logisk overblik over en algoritme ved brug af matematiske udtryk.

Eksempelvis ved en binary search vil, invarianten være "Hvis  $x$  findes i  $L$ , findes den i  $L[l..r]$ "

# 5 Merging and hashing

## 5.1 Datastruktur

Datastruktur er metoder til at hente og gemme data.

Interface er en række af metoder til at håndtere dataen

Der er to grundtyper af interfaces Sekventiel og random access.

### 5.1.1 Sekventiel tilgang

Sekventiel interface håndtere data i en lang stream som et dokument. Her er metoder som, `writeNext()`, `open()`, `close()`, `endOfFile()`, `readNext()`, osv.

### 5.1.2 Random access

Data er gemt i et område, med en tilhørende ID. Metoder: `findElm(ID)`, `sertELM(id, data)`, `deleteELM(ID)`, `open()`, `close()`, osv.

### 5.1.3 Merging

Merging er en teknik, hvor to sorteret lister bliver merged. Her bliver det første element fra begge altid sammenlignet og det mindste bliver taget ud og sat i den nye sorteret list.

Køretiden er her  $O(\log(n)n)$  grundet hvor  $n$  er størrelsen på samtlige elementer og da der bliver færre pr. iteration vil den ikke køre  $O(n^2)$  men  $O(\log(n)n)$ .

Ved ikke allerede sorteret lister, vil alle elementer blive skilt ad, til lister med 1 element, hvorefter de bliver parret op.

## 5.2 Hashing

Hashing er en metode brugt til at validere filer og lign. Her kan et hashing metode bruges, hvor dataen bliver komprimeret ned, til en mindre mængde af data, hvor den kan bruges som validering til nye filer der bliver komprimeret ned.

Den mest simple hashing metode er, at tage en array med en given størrelse, og på alle keys i den nye data, vil de blive udregnet `key%hash size`. Her vil hvert element så blive indsat på den nye plads. Dette vil dog resultere i overlap, som kan ordnet på forskellige måder med den mest simple chaining som er at lave en array med element værdi og adresse til en ny array med element og adresse til den næste osv. En hashing metode vil ikke kunne regnes tilbage og kan derfor kun bruges til validering og ikke komprimering alene.

### 5.2.1 Sandsynlighed for overlap

$S_n = s_{n-1} \cdot \frac{L-(n-1)}{L}$  Hvor L er hash længden.

## 6 Machine Learning

Agent a software.

An agent is learning if it improves its performance by observation.

Machine learning is usefull for

- Non anticipated possible solutions
- Non anticipated solution over time
- No programmable solution

### 6.1 Forms of machine learning

Supervised learning - Learns from problems with solutions and apply them to future problems. Ex. housing pricing.

Unsupervised learning - Problems have no real solution but rather categorize data. Ex. Categorizing Amazon users in buying habits.

Reinforcement learning - Rule judges if agent is right or wrong from which it learns. Ex. AI learns games

### 6.2 Supervised Learning

Input = features

Output = responses

Learning is based on  $m$  input vector with all features where it is a set with  $\{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_m, y_m)\}$

From learning it will generate a function  $g(\vec{x})$  from which  $g(\vec{x}_m) \approx y_m$

$\hat{y} = g(\vec{x})$

#### 6.2.1 Regression problem

Used to predict a quantitative solution

Nearest neighbor - Takes the nearest points to the given  $x$  values and returns the average  $\hat{y} = \frac{1}{k} \sum_{\vec{x}_i \in N_k(\vec{x})} y_i = g(\vec{x})$

To find the best function in a 2 dimensional regression a loss/error function

can be defined as  $\hat{L}(a, b) = \sum_{i=1}^m (y_i - ax_i - b)^2$

The more general loss for the poly function is

$$\hat{L}(\vec{\theta}) = \sum_{i=1}^m (y_i - h(\vec{\theta}, \vec{x}_i))^2 = \sum_{i=1}^m (y_i - \sum_{j=0}^p \theta_j x_{ij})^2$$

Where  $h(x) = \text{poly}(\vec{\theta}, x) = \theta_0 + \theta_1 x + \dots + \theta_k x^k$

The functions distance can be squared or absolute to counter negative variation and squared to make higher values more important.

From this partial differential can be used where one value is constant and the differentiation is found and set equal to zero to find one constant.

From this the other constant can be found through differentiation and finding it when equal zero. This is done until all variables are found.

With this technique a local optimum for the lowest error.

### 6.2.2 Classification problem

used to predict qualitative (binary) solution Nearest Neighbor - Just like regression but instead of average, it takes the majority

## 6.3 General framework

Learning = Representation + Evaluation + Optimization

Representation: How input is represented. Ex. the hypothesis set for regression problems

Evaluation: Defining the loss function

Optimization: Finding the best learner in the language for minimizing the loss

## 6.4 Overfitting

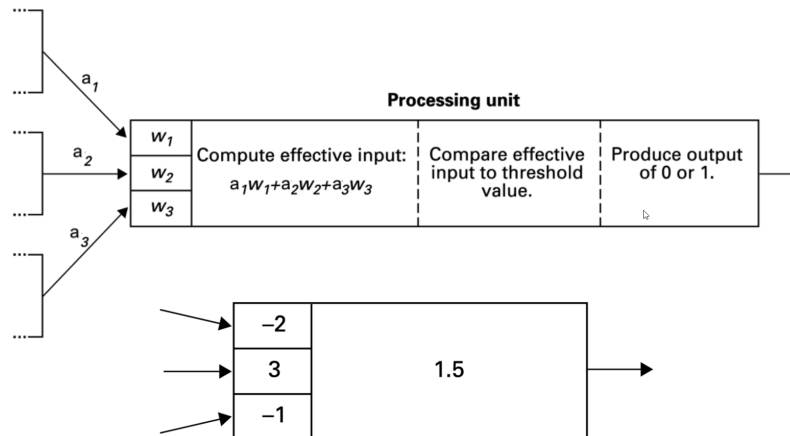
Overfitting is the act of using a too high  $k$  value in the poly regression. Here a too high value will end up hitting all data, but as a result most likely not represent the actual underlying function.

## 6.5 Training and Assessment

To ensure not overfitting the training methods can be used. Here there are two techniques.

- Holdout method - Using only 75% ish for training and 25% for testing
- K-fold cross validation - Used in small data, where samples of the training data is used for testing

## 6.6 McCulloch-pitts "units"



$a$  is input

$w$  is weights

The comparing can also be made an input is  $a_0$  and instead it compares to 0.

Instead of comparing to 0 an activation function can be used. This will make it derivable to find the minimum.

The most common is the sigmoid function  $\frac{1}{1+e^{-x_i}}$

Neurons are able to replicate boolean functions and therefore opens neurons possibilities a lot.

Though the XOR is not directly possible due to it not being possible to create a activation function which includes (0,1) and (1,0) but not (0,0) and not (1,1).

## 6.7 Network structures

Structure(/Architecture): definition of number of nodes, interconnections and activation functions  $g$  (but not weights). Two types of structures Feed forward networks and Recurrent networks

Feed-forward network has no cycles and can be single-layer (no hidden layer) or multi-layer perceptrons (hidden layers). Used in classification and regression.

Recurrent networks has cycles between nodes which result in possible memory in the network

Multiplayered perceptrons have hidden units typically choosen by hand

With more layers the threshold function will combine meaning a more precise threshold

## 7 Feature spaces for representation of images

Features are ways to describe an image such as color distribution, shapes and so on.

### 7.1 Color histogram

Color histograms are used to represent the distrubution of colors.

here each color is distributed into intensity and then normalised with the total pixels used.

Here the amount of distribution determine the amount of compression.

### 7.2 Distance of color histograms

Euclidean distance - representring the histogram as a vector and using calculating the distance by subtracting every points and multiplying them together and taking square root. Ex

$$\text{dist}(P, Q) = \sqrt{(h_P - h_Q) \cdot (h_P - h_Q)}$$

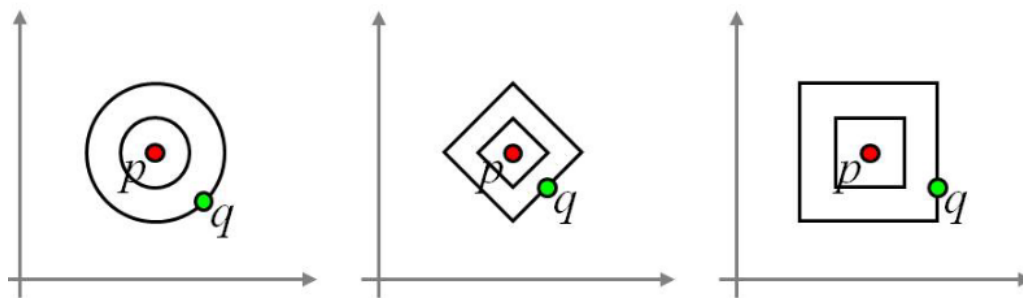
Psychologic distance - takes in measures such as pink and cyan being close together as colors and weights that higher.

Distance can also be measured in different ways, for example let it be a classification problem the distance can be

$$\text{dist}_P(p, q) = w_1(|p_1 - q_1|^P \cdot w_2|p_2 - q_2|)^{\frac{1}{P}}$$

Where  $w$  is the different weights.

Here if  $P = 2$  we get the distance formula from before. Different  $P$  values will here create different shapes of classification. The different shapes can be seen on the picture with 1 distance unit between the two black lines.



Left -  $L_2$   
Middle -  $L_1$   
Right -  $L_\infty$

As seen the higher  $P$  values results in the corners being included.

## 8 Clustering

Clustering is used to categorize data.

Good clustering will have an amount of clusters which will result in furthest distance between clusters and shortest distance between point and clusters.

There are algorithms which will help clustering data points.

### 8.1 Lloyd/forgy

Assign all points random clusters.

Calculate the mean point of all clusters.

Assign each point to closest cluster.

Recalcluate mean cluster points.

Repeat until no change.

### 8.2 MacQueen/ $k$ -means

Assign all points random clusters

Calculate mean cluster points

In a given order assign an element to closest cluster point and recalculate cluster mean points.

Repaeat last step until no change.



## 9 Models of Computation, languages and Recursion

### 9.1 Models of Computation

Models of computation are ways of working with problems the most effective way

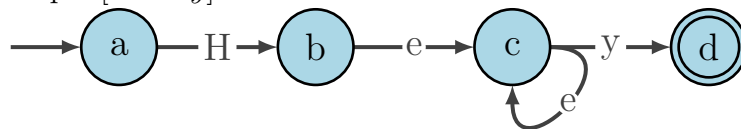
The models focus on important parts of problems and sort out unimportant parts

Example compilers how the first took centuries to develop and today is a simple project.

### 9.2 Deterministic Finite Automaton (DFA)

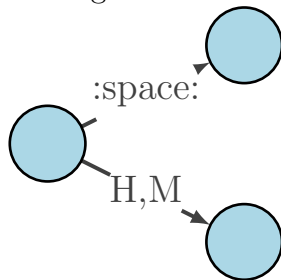
DFA is a way of illustrating regular expressions via graphs.

Forexample  $[He + y]$  will be:



The double ring at d illustrates a acceptance state, inwhich a DFA ends.

An edge can have multiple poosibilites but not the same.



The : notation is used to describe groups of characters or not visible.

Concatenation is combining multiple DFAs

DFAs are limited, one example is the langauge made of ( and ) which always has to matched.

Here the DFA has to be endless, to compensate for an endless string of ( and ).

### 9.3 Context-Free Grammar

CFG is a rule based regex. Here characters is terminals and regex applications (possibilities) are called derivations

Example:

$S \rightarrow \text{Hello } T$

$S \rightarrow \text{Hey } T$

$T \rightarrow \text{there}$

It will here insert rules if given in the string of a rule.

The derivations are: Hey there og Hello there

CFG is able to store memory using recursion. To describe the language of matched parenthesis like the DFA not could it would be.

$S \rightarrow (S)$

$S \rightarrow SS$

$S \rightarrow ()$

By recursion it will be able to go to the closing parenthesis after an opening.

## 10 Online algorithms

Algorithms based on unexpected data. Unlike offline data where data is known and the algorithm can be optimized based on data online algorithms has to optimize for the unknown.

An algorithms competitive is the ratio between offline and online performance. From which the c-competitive ratio is  $\forall I : \frac{ALG(I)}{OPT(I)} \leq C$ .

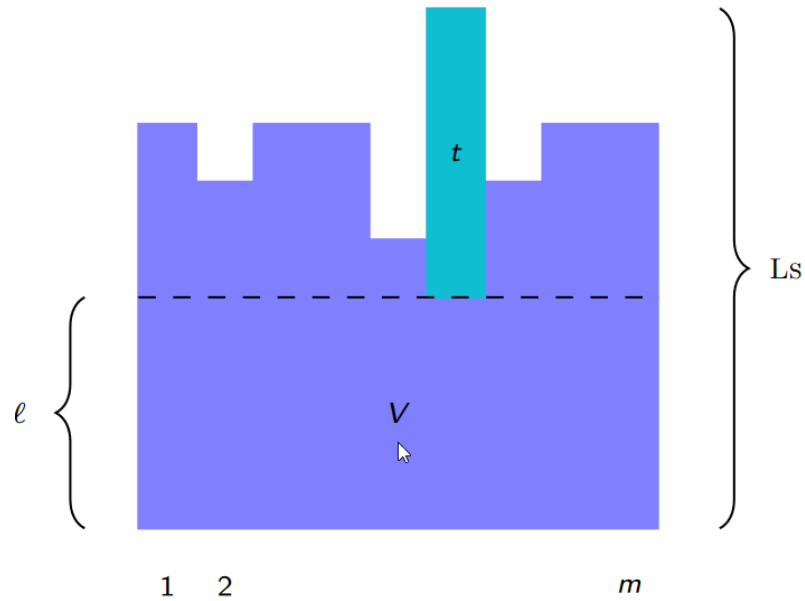
Therefore an algorithms ratio is the ratio between the online and offline guarantee worst case for an algorithm.

An example is computation queue algorithms which gives computations to  $m$  computers.

The offline version can see all the data and optimize for all the data, and the online version will only see the next data.

The online algorithm can here be given the computer with the least amount of work the next computation.

The ratio can here be proofed by the following:



$t$  is the last job

$l$  is the length of the lowest job before  $t$

$T$  is the total length of all job

$V = l \cdot m$

$\text{OPT} \geq T/m$  and  $\text{OPT} \geq t$

$T \geq v + t$

$$\begin{aligned}
 L_S &= l + t \\
 &\geq \frac{T - t}{m} + t && \text{since } l = \frac{V}{m} \text{ and } V \geq T - t \\
 &= \frac{T}{m} + (1 - \frac{1}{m})t \\
 &\geq \text{OPT} + (1 - \frac{1}{m})\text{opt} && \text{since } \text{OPT} \geq T/m \text{ and } \text{OPT} \geq t \\
 &= (2 - \frac{1}{m})\text{OPT}
 \end{aligned}$$

## 11 Satisfiability

A propositional formula is satisfiable if variables can be assigned such it returns true.

By describing problems via propositional formulas, it can be solved by SATs.

A SAT is a program which can find solutions for propositional formulas.

To use a SAT the formula has to be CNF by:

- $F \rightarrow G = \neg F \vee G$
- $\neg(F \wedge G) = \neg F \vee \neg G$
- $\neg(F \vee G) = \neg F \wedge \neg G$
- $\neg(\neg F) = F$
- $F \vee (G \wedge H) = F \vee G) \wedge (F \vee H)$

For a SAT input lines beginnings are compiled as:

- p - Problem formatted as *p end #variable #cause*
- c - comment
- s - output for solution
- v - variable assignment from output

Lines are OR'ed and rows are AND'ed.

SAT solving is based on trying every possible solution

This is optimized by:

- Incremental assignments
- Backtracking solver
- Pruning the search
- Backjumping
- Conflict-driven learning
- Restarts
- Forgetting

## 12 Databaser

Databases are a software tool to store and organize large amount of data.  
Data can be accessed by multiple users and data are enforced by constraints.

DBMS - database management systemt

Modern DBMS - helps with backup, permission, consistency, searching sorting filtering, multi-users

There are to types of DBMS

Data entries are called tuple

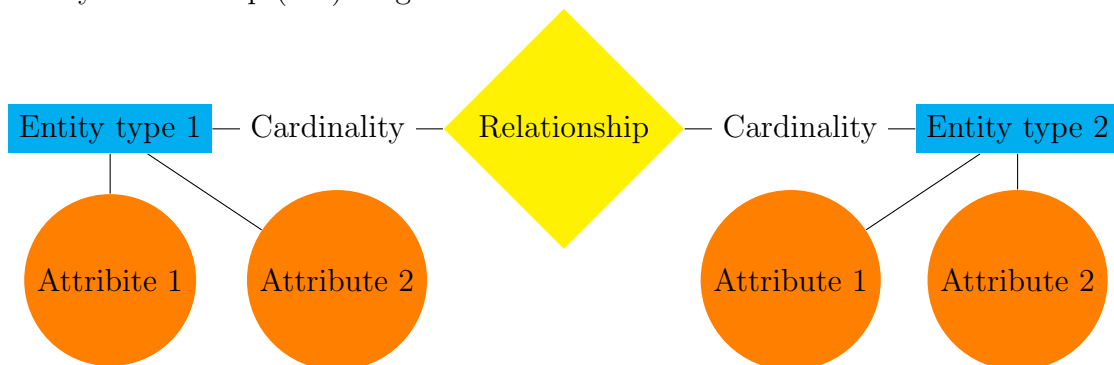
- Relational DBMS - Table system with data
- Graph DBMS - Directional graph good to see patterns

A database is structured as:

- Conceptual - ER diagram
- Logical - Relation diagram
- Physical - SQL language

### 12.1 Conceptual

Entity-relationship (ER) diagrams



Cardinality describes the amount of relation. Ex. 1..n mean there must be at least 1 to n amount in the given relation.

Relation - describes the relation between two entities: Ex. owns

Relation can also have attributes such as date of ownership.

## 12.2 Logical

Here the tables are somewhat established

This is through:

$relation_{name}(attribute1, attribute2, \dots)$  or

$relation_{name}(attribute1 : type1, \dots)$

## 12.3 Physical

The SQL command level

Remember string must be with ' and not "

The usual commands *CREATE*, *INSERT*, *DELETE*, *SELECT*, *FROM*, *WHERE*

### 12.3.1 Relatio Algebra

$\sigma_c(R)$  Selection - gather rows with the given  $c$  condition

$\pi_{A_1, \dots, A_k}$  projection - gather columns  $A_1, \dots, A_k$  All of the usual set operators

$\cup, \cap, \times \dots$

## 12.4 Integrity Constraints

Constraints is in form of data types and size limit

This is things such as name begin *CHAR*(20)

The ensure data integrity

## 12.5 Primary key

Column with data which is unique data.

Can not have duplicate to ensure data integrity

## 13 Cryptology

The science of cryptology is in two genres:

Cryptography is written something hidden

Cryptanalysis is breaking cryptography

The goal of cryptography is securing

- confidentiality to secure encryption
- Integrity is often done through hash keys which verify the message

- Authenticy is also a part of integrity
- Non-repudiation is important so its validated that a message actually was sent

### 13.1 Encryption

Encryptions goal is to have a key which maps the messages and is bijections so decryption is possible

Message space is the readable message and cipher space is the encrypted format

The most simple is ceasar sipher which mapping shift all charachters by the key number

The problems is bruteforcing is pretty easy with only 29 keys.

A good number is  $2^{128}$  for amount of keys so its not bruteforceable

Kerckhoffs's Principle - The security should rely in the key and not the system

Vigenère Cipher - Ceasar cipher but with a rotating key such as 314

One-Time-Pad - using a non repeating key but an infinite long key

The problem is integrity is hard to ensure no bit flip and sharing the key is hard

A good way to analyse different encryption methods is images, which is easy to find patterns in.

### 13.2 Symmetric key systems

Symmetric - both sides has the same key.

- Caesar cipher
- Enigma
- DES - had only  $2^{56}$
- Tiple DES - had  $2^{112}$  but it was a combination of DES which resulted in flaws
- AES -  $2^{128}$  or  $2^{256}$  keys and currently not bruteforceable, worked since 2001 and standard today

### 13.3 Assymtric keys

To counter the problem with two people having the same key a private key and public key is used

People sending a message to a person they encrypt it with the persons public key

Then it is decrypted with the private key

Here the private key has to be hard to find out both from bruteforce or reverse engineering the public key.

### 13.4 RSA

$N_A = p_A \cdot q_A$  where the right side is primes

$\gcd(e_A, (p_A - 1)(q_A - 1)) = 1$

$e_A \cdot d_A \equiv 1 \pmod{(p_A - 1)(q_A - 1)}$ .

Public key  $PK_A = (N_A, e_A)$

Private key  $SK_A = (N_A, d_A)$

To encrypt:

$c = E(m, PK_A) = m^{e_A} \pmod{N_A}$

To Decrypt:

$r = D(c, SK_A) = c^{d_A} \pmod{N_A}$ .

#### 13.4.1 Why it works

We want to proof that  $r = (d(E(m, PK_A)SK_A)$

We know that  $e_A d_A = 1 + k(p_A - 1)(q_A - 1)$

$$r \equiv (m^{e_A} \pmod{N_A})^{d_A} \pmod{N_A} \quad (1)$$

$$\equiv m^{e_A d_A} \pmod{N_A} \quad (2)$$

$$\equiv m^{e_A \cdot d_A} \pmod{N_A} \quad (3)$$

$$\equiv m^{1+k(p_A-1)(q_A-1)} \pmod{N_A} \quad (4)$$

1. first we insert the definition on encryption and decryption.
2. This is then reduced
4. Then the definiiton of  $e_a \cdot d_A$  is inserted



$$r \equiv m^{1+k(p_A-1)(q_A-1)} \pmod{p_A} \quad (5)$$

$$\equiv m \cdot (m^{(p_A-1)})^{k(q_A-1)} \pmod{p_A} \quad (6)$$

$$\equiv m \cdot 1^{k(q_A-1)} \pmod{p_A} \quad (7)$$

$$\equiv m \pmod{p_A} \quad (8)$$

5. We then use  $p_A$  instead of  $N_A$

6. We then rewrite

7. By Fermat's little rule we know that  $m^{(p_A-1)} \pmod{p_A}$  is equivalent to 1

8. We then reduce by everything to the power of 1 is 1

We see in case we use  $q_A$  instead of  $N_A$  like with  $p_A$  it will result in the same result  $r \equiv m \pmod{q_A}$ . We can then use Chinese remainder theorem due to  $\gcd(p_A, q_A) = 1$

$$r \equiv m \pmod{p}$$

$$r \equiv m \pmod{q}$$

$$N = pq$$

$$r \equiv m \equiv m \pmod{N}$$

$$r \equiv m \pmod{N}$$

This means  $m$  can be everything due to  $m$  always will be equivalence to  $m$ . And since  $r$  is equivalent to  $m \pmod{p, q, N}$  and  $r$  and  $m$  can not be larger than  $N - 1$  must therefore be equal  $r = m$

### 13.4.2 Why RSA is secure

Well to decrypt a text given  $c$  and  $PK_A = (N_A, e_A)$ , it has to find

$$c = m^{e_A} \pmod{N_A}$$

Which there is no simple algorithm to solve and brute force is needed.

To ensure a private key is not found the factoring of  $N_A$  must be hard. To brute force  $N_A$  if a single prime is known it will only take  $O(\sqrt{N})$  time. This means if  $N$  is 3072 bits it will take almost forever.

If the length is  $n = \lceil \log_2(N + 1) \rceil$  the run time will be  $O(2^{n/2})$

But it may not be true due to quantum computers can have more efficient algorithms

### 13.4.3 RSA implementation

To implement RSA for the encryption part the expression  $a^k \pmod n$  is used. If done as said there would be  $k - 1$  modular multiplications, which is too many

This can be optimized by a recursive function where if the power is odd its possible to minus 1 from the exponent and multiply, and if even it can be halved and multiplied. RSA program implements this for more detail.

To find  $e_A$  and  $d_A$  the following expression has to be used

$$\begin{aligned} \gcd(e_A, (p_A - 1)(q_A - 1)) &= 1 \\ e_A \cdot d_A &\equiv 1 \pmod{(p_A - 1)(q_A - 1)} \end{aligned}$$

This can be done with the extended Euclidean algorithm

This can find  $d_A$  through the gcd. The algorithm is implemented in the RSA program.

To find prime numbers there are a number of different ways

- Sieve of Eratosthenes - Take all numbers the first in the list is a prime and after removing it remove all multiples
- Bruteforce - Take all numbers from 0 to  $\sqrt{n}$  and test if it divides with the prime
- Miller-Rabin test - Choose a random number from which it can be tested if Fermat test works

The most effective is the Miller-Rabin which is implemented in the RSA program.

Carmichael numbers are numbers which are hard to determine if composit, ex 561 where only 241 out of the 560 determine it as a composit.

## 13.5 Symmetric and public key system

To speed up RSA due to it slow nature, a symmetric session key is generated from the RSA.

Digital signatures is also implemented with RSA, by encrypting the fixed length:q signature with the private key and public keys can then validate if it can be decrypted.

## 14 Graph theory

Vertices points in graph

Edges lines connecting vertices

Degree number of edges on a vertice

Categories of graphs

- labelled graph - vertices has names
- Directed graph - vertices has a direction
- Weighted graph - edges has an weight for moving past

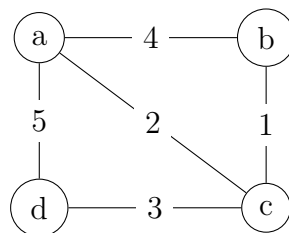
Walks - a route on a graph which goes between two vertices

Path - walk without any repetition on an edge

Cycles - Closed walks

Connectedness - If a number of graph is in more parts

### 14.1 Matrix representation



Incidence Matrix

Edges has labels

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

Adjacency Matrix

Edges has weight

$$\begin{bmatrix} 0 & 4 & 2 & 5 \\ 4 & 0 & 1 & 0 \\ 2 & 1 & 0 & 3 \\ 5 & 0 & 3 & 0 \end{bmatrix}$$

To

find paths between points the adjacent matrix can be multiplied with itself  $k$  ( $A^k$ ) times where if there is a 1 there exists a path with  $k$  length.

For the shortest path between two points the adjecenst matrice can be multiplied but instead of multiply plus is used and instead of plus the minimum of the sums is choosen.

It then has to be multiplied the amount of vertices -1.

$$w = \begin{bmatrix} 0 & 4 & 2 & 5 \\ 4 & 0 & 1 & 0 \\ 2 & 1 & 0 & 3 \\ 5 & 0 & 3 & 0 \end{bmatrix} \quad w^2 = \begin{bmatrix} 0 & 3 & 2 & 4 \\ 3 & 0 & 1 & 0 \\ 2 & 1 & 0 & 1 \\ 4 & 0 & 1 & 0 \end{bmatrix} \quad w^3 = \begin{bmatrix} 0 & 3 & 2 & 3 \\ 3 & 0 & 1 & 0 \\ 2 & 1 & 0 & 1 \\ 3 & 0 & 1 & 0 \end{bmatrix}$$

To optimize the matrix can be multiplied by itself so for a graph with 121 vertices instead of 121 computations the matrix can be  $a^{2 \cdot 2 \cdot 2 \cdot 2 \cdot 2 \cdot 2}$