

CO-VL

Rückfing:

Man erzählt was zur Projektaufgabe. Ca. zur Mitte des Semesters legt man sich fest.

Dazu erzählen und zur Beziehung zur VL setzen!

#3 Python + pandas → man kann so ähnlich wie in R arbeiten (Tabelle einer DB)

Exercises:

1. 4^n

2.

$\begin{matrix} AA \\ TT \end{matrix}$	$\begin{matrix} AC \\ TG \end{matrix}$	$\begin{matrix} AG \\ TC \end{matrix}$	$\begin{matrix} AT \\ TA \end{matrix}$	$\frac{12}{2} + 4 = 10$
$\begin{matrix} CA \\ GT \end{matrix}$	$\begin{matrix} CG \\ GC \end{matrix}$	$\begin{matrix} CC \\ GG \end{matrix}$		
$\begin{matrix} TC \\ AG \end{matrix}$	$\begin{matrix} TA \\ AT \end{matrix}$	$\begin{matrix} CC \\ CG \end{matrix}$		

bei ungerade $\frac{\#RNA}{2}$

bei gerade komplizierter

⇒ deswegen arbeitet man eigl. mit ungerade Nucleotid-Anzahl

Bei Genome-Assembler kann man mit ungerades k au-geben!

A

Allgemeine Formel für gerade aufschreiben

$$\frac{4^n - 4^{\frac{n}{2}}}{2} + 4^{\frac{n}{2}}$$

bei n gerade

$\frac{n}{2}$ geht auch, da n gerade sein muss.

$$4^{\frac{n}{2}}$$

bei n ungerade

Kleine sequenzierte Fragmente DNA

20.04.17

CO-U

N : # Fragmente
 F : Fragmentlänge (100-1000)

G : Genomlänge ($10^6 - 10^9$)

(Säugetiere, Pflanzen
[Viren: $10^3 - 10^4$])

$$C = \frac{N \cdot F}{L}$$

durchschnittliche Coverage



je größer C , desto eher eine Normalverteilung dran

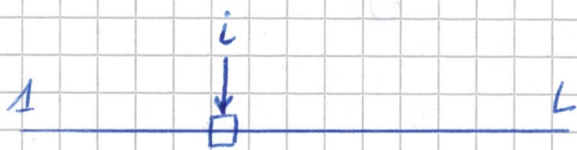
Normalverteilung stimmt nicht ganz

- durch 0 beschränkt
- in Realität wir ganze Zahlen

Idee:

geworfene Bälle

NF Bälle zufällig auf L Positionen



Jeder Ball landet mit Wkt $\frac{1}{L}$ bei i .

Wenn die Bälle unabhängig sind:

(Was nicht ist, aber obs ignorieren wir.)

Binomial-Verteilung!

Binomial-Verteilung!

X: Anzahl Bälle bei i

$$P(X=k) = \binom{NF}{k} \left(\frac{1}{L}\right)^k \left(1 - \frac{1}{L}\right)^{NF-k} \quad k=0,1,\dots,NF$$

$n \rightarrow \infty$
 $p \rightarrow 0$
 $np \rightarrow C$
 wie wir uns das angucken

$NF = 10^{11}$
 $k = 30$
 $L = 3 \cdot 3 \cdot 10^9$

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

$$= \frac{n^k}{k!}$$

Allg.: $\binom{n}{k} p^k (1-p)^{n-k}$

← "n hoch k fallend"

$$= \frac{n(n-1)(n-2)\dots [k \text{ Faktoren}]}{k!}$$

Zahlen sind sehr groß, bzw. sehr klein.
 => beliebig viele Probleme beim numerischen Rechnen
 => approximieren

n sehr groß (10^{11})
 p sehr klein ($1/(3 \cdot 3 \cdot 10^9)$)
 $p = 1/L$
 $np = C = 2 \cdot B \cdot 30$

weder sehr groß, noch sehr klein, sondern ein unwichtiges Wert

Approximation:

$$\log \left[\frac{n^k}{k!} p^k (1-p)^{n-k} \right]$$

Logarithmus wandelt aus Produkten Summen

$$= \log n^k - \log k! + k \cdot \log p + (n-k) \cdot \log(1-p)$$

Was darf man approximieren und was nicht?
 => nicht "zu falsche" Approximation, bzw. nicht zu große Fehler

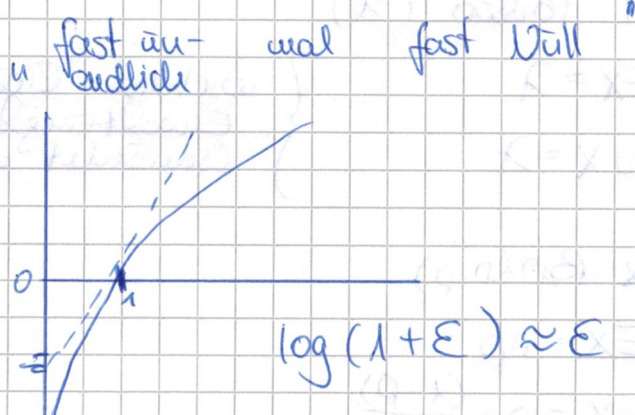
z.B. $n^k \approx n^k$
 \downarrow
 n^k

Unterschied nicht schließen, wenn n sehr groß und k eher klein!
 \downarrow
 $n \cdot (-p)$

= "und wie passiert ein Wurf?"

$$\approx k \cdot \log n - \log k! + k \log p - np$$

$$= k \log(np) - \log k! - np$$



np ist endlich und eine konst./konstante Zahl

$$= k \log C - \log k! - C$$

$$\approx (k \log k)$$

es gibt genauere Darstellungen von log k!

$$\exp[-c] = \frac{c^k}{k!} \cdot \exp(-c)$$

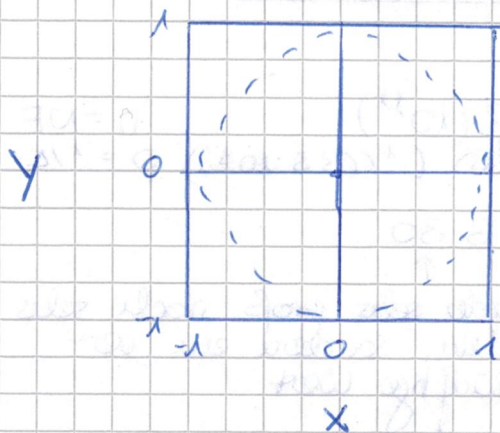
$$P(X=k) = \binom{NF}{k} \left(\frac{c}{NF}\right)^k \left(1 - \frac{c}{NF}\right)^{NF-k}$$

$$= e^{-c} \cdot \frac{c^k}{k!}$$

Poisson-Verteilung mit Erwartungswert c

(Approximation der Binomial-Verteilung mit entspr. Werten)

Wo kommt Poisson-Verteilung noch vor?



Kreis Fläche $\pi \approx 3.141$

Quadrat Fläche 4

Zufällig Punkte "reinschießen"

↳ mit gleichem Wert 2 Koordinaten (je x, y ziehen)

$x^2 + y^2 \leq 1^2$ für jeden Punkt testen, ob im Kreis oder nicht

$$p = \frac{\pi}{4}$$

X : Punkte im Kreis

$EX: np \approx X$ Erwartungswert X

ungefähr Beobachtung X

$$\pi = 4p \approx \frac{4x}{n} \quad \text{genauer je größer } n \text{ ist}$$

mit Poisson-Verteilung approximieren!

X Poisson (λ)

$$EX = \lambda$$

$$\text{Var } X = \lambda$$

wichtige Eigenschaft:
Erwartungswert und Varianz gleich
(Schranke aber nicht ein)

X Bin(n, p)

$$EX = np$$

$$\text{Var } X = np \frac{1-p}{1}$$

"Gefülltes Palindrom"

Ensemble Genomes

> Ensemble Bacteria

E.coli

> Download genes

> DNA

sm = softmasks

rm

sonst

• fa.gz mit Gzip FASTA-Datei

→ repeats mit Kleinbuchstaben

→ "-" ersetzt mit x o.ä.

→ normal

oplevel

nonchromosomal

→ alles

→ nicht auf Chromosomen

⇒ Softmask am angenehmsten, da man noch selber
entscheiden kann

Script:

→ nur privat von Sven, stattdessen olnopy

```
from dnair import fasta_read
```

(offiziell)

```
with open
```

```
for entry in fasta_reads(fname):
```

```
    (header, seq)
```

```
    (:
```

S. Elias

⇒ C+G und A+T sind "gleiches" zueinander: als zu den
anderen von der Anzahl her.

CO-VL

(...)



DNA-Sequenzierung

gleiche Länge
"die Zyklen" da
billig, hat sich
aber die Fehlerrate
nicht durchgesetzt

fluoreszierende
Stoffe

pH-Wert → sehr ungenau

Minima → viel teurer und viel genauer als mit pH
i.H. ca. 80-90% aller Sequenzierungen
(Fehlerrate < 1%)

durch Nano-Pore → neue Technik (Einzelmolekültechnik)
(höhere Fehlerrate)

⇒ die meisten können mit einem Molekül nichts anfangen

→ PCR-Amplifizierung, um ausreichende Signalstärke zu bekommen

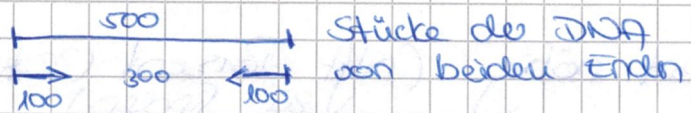


PCR

- zum Glück recht geringe Fehlerrate
- Fehler am Anfang schwierig zu unterscheiden

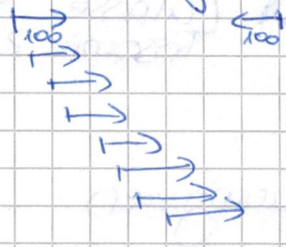


Paired-end-sequencing



⇒ 2 × 100 bp Stücke sequenziert

Man hat gewisse Fehlerraten (PCR, Sequ., ...)



durch Überdecken kann man Fehler korrigieren

Fragmente: 100-200bp

↑ mehrere Millionen davon
⇒ meist mit Qualitätsmetriken (Stichwort: FASTQ)

ERROR CORRECTION PROBLEM (#3)



DE-DUPLICATION PROBLEM (#4)

Ziel: unabhängige Sequenzen bekommen
(PCR-Dubletten rausfiltern)

(#5) k-mer \equiv q-gram

- es ist in der Tat schwierig
- immer noch aktuelles Thema (wird immer noch viel publiziert)
- z.B. 16 Byte (oder mehr) pro k-mer und Zahl und man will schneller Zugriff haben \rightarrow so lange klein und ins RAM möglich, sonst schwierig

- was macht man?

- k-mer als Integer (128 Bit \equiv ~~63~~ 63-mer)

- * k bis 31 gehen, da es in 64 Bit passt \uparrow wir betrachten wir ungeade
- * revers komplementäre Molekül! \rightarrow 2 versch. Integers \rightarrow es wird nur das "kanonische Integer" betrachtet \rightarrow man wählt den kleineren Integer-Wert

- Datenstruktur

- * kleines k \rightarrow Integer-Array
- * bei großen k geht das nicht 2^{62} (mehrere Ex-Byte) \Rightarrow k-mer codes hashen
- * Bloom-Filter

- k-mer Poisson verteilt (mit Mittelwert λ)

solide vs. schwache k-mer

↑
vernünftliche
korrelier

↑ wenn in PCR früh Fehler o.ä.

\Rightarrow guter Schwellenwert nötig (hilft Poisson-Verteilung zu finden)

- biologische Effekte:

- lokale Abschnitte sind dupliziert (krasser Fall: Trisomie 21)

- sequenzielle Effekte

\Rightarrow ggf. nicht wegschneifen, sondern korrigieren

(#7) Methode für OIen ggf. in Ordnung (nicht pauschal "weschneifen") \rightarrow dann am einfachsten

(#8) versch. Phylosofien: Hash-Tabellen sollten zw. 50% und 70% gefüllt werden

(Hashing ist sehr weites und sehr interessantes Gebiet)

- Strategien um Kollisionen aufzulösen (z.B. lineares sondieren)

↳ ~~maine einfachste Frage ist die Adresse des Arrays~~

▷ Bloom-Filter (#9-11)

Schlüssel drin?

nein \rightarrow bin sicher das nicht drin
ja \rightarrow könnte drin sein, muss aber nicht

Wüt f. falsch positiv ist bekannt

Schlüssel können aus Bloom-Filtern nicht wieder raus genommen werden.

\rightarrow man kann bis 1 zählen

\leadsto Es gibt Varianten von Bloom-Filtern, mit denen man auch Werte zählen kann

• Bloom-Filter ist Bit-Array

Bloom Filter B

1010101 1101 m Bits

Universum X : "beliebig" groß, alle Schlüssel,
z.B. $|X| = 4^{71}$

k Hash funkt: $h_i: X \rightarrow \{0, \dots, m-1\}$

B.add(x) : Setze $B[h_i(x)] := 1 \quad \forall i=1, \dots, k$

x in B ?

$$\bigwedge_{i=1..k} B[h_i(x)]$$

Was kann passieren?



Man muss versuchen, dass die Wüt, dass man zufällig nur 1en trifft sehr gering ist

\Rightarrow größeres $k \Rightarrow$ mehr Nullen

\Rightarrow Nachteil: dauert länger, mehr Speicherplatz

Man kann das ausrechnen!

(Vgl. "lage" in ADBS)

s. #11

- wenn man zu viele Werte einfügt, sind alle Werte auf 1 \Rightarrow abschätzen wie viele Werte man einfügt
- false-positive-Rate auf ϵ beschränkt
 \hookrightarrow eigentliche ϵ

Design-Aufgabe:

Wähle (m, k) so, dass $FPR \leq \epsilon$, wenn $\leq n$ Schlüssel eingefügt werden.

Man reduziert approximativ viele Annahmen die nicht wirklich so eintreffen

n Elemente rein ich gucke mir ein zufälliges Bit an. wie groß 'wilt' das es gerade 10 ist.

Wie ich hier (#11) reduziere ist falsch, aber wie es ist ähnlich wie in Approximation der Poisson-Verteilung. Es ist nur "ein bisschen falsch", aber nicht "schlimm falsch".

bestimmte Anzahl 1-Bits.

$$\text{Anteil 1-Bits: } \approx \frac{kn}{m}$$

$$\text{Anteil 0-Bits: } \approx 1 - \frac{kn}{m}$$

kn klein gegenüber m ist $1 - \frac{kn}{m}$ also richtig?

$$\text{wenn } kn \text{ größer wird: } 1 - \frac{kn}{m} + \frac{\left(\frac{kn}{m}\right)^2}{2!} \approx \dots$$

$$\approx \exp\left(-\frac{kn}{m}\right)$$

Auf Folie ist nur Näherungsweise. So stimmen würde es dann, wenn alles unabhängig ist.

\hookrightarrow in Praxis unvollständig

Approximation ähnlich wie letztes Mal mit Log.

Warum gibt es ein k das minimiert? Warum wird es nicht immer besser, je mehr Hash-Fkt. ich nehme?
 \rightarrow es gibt eine feste Größe und dann würde es mehr Ten geben

Bloom-Filtere sind eine schöne Sache, aber man kann z.B. nicht zählen.

Aufgabe Paper "lighter" lesen → muss nicht ganz gelesen werden, aber "Trick finden" wie man k-mere trennen kann ohne echt zu zählen!

Ausblick Übung:

- Aktueller Stand Hashing
- Bloom-Filter

27.04.17

co-ü

in Arbeitsgruppe Paketmanager Conda

hier: Miniconda

<https://conda.io/miniconda.html>

Umgebungsname

```
conda create --name omics python=3 jupyter numpy  
scipy matplotlib ib seaborn pandas hspy numba
```

managt Paketabhängigkeiten

Mit Conda kann ich beliebig viele von solchen Umgebungen auflegen mit speziellen Konstellationen

⇒ beeinflusst nicht normales Systempython!

Wie verwende ich das?

`source activate omics`

(bis `source deactivate` legt er Pfade um)

jupyter notebook

> Python 3 auswählen

zeigt Wert letzten Rückwerts
mit Shift + Return



Bloomfilter:

Fehlerrate ist ϵ in Abhängigkeit von k , m und n und ich kann k optimieren

10^{-6} ist kleine Falsch-Positiv-Rate

Näiv überlegt:

$|X| = 4^{71}$, davon kommen 4^{16} vor

4^{16} versch. Objekte repräsentieren

\Rightarrow wir benötigen Vielfaches von 4^{16}

\leadsto für solche Spieldereien Notebooks benutzen.

Da kann ich auch selb. HTML-Seite machen und verschicken.

[Notebook s. Elias]

$\Rightarrow x, y$ auf log. Skala \Rightarrow eine Gerade

ganz am Anfang gegen 1, aber ab da eine Gerade

Bsp.: mit 10^{-9} Fehlerquote $\rightarrow x=40 \Rightarrow 5$ Byte
 $\Rightarrow 40$ GB Platz f. Bloom-Filter

Plotten:

`%matplotlib inline` \leftarrow "magic" Kommando f. Notebook

`plt.plot (x-Achse, y-Achse)`

`plt.yscale ('log')`

`plt.show()`

\triangleright Fehlerquote von 10^{-6} oder 10^{-9} klingt erst mal gut

4^{16} Schlüssel

$$(4^{71} - 4^{16}) \cdot 10^{-9}$$

$$= 4^{71} \cdot 4^{-15} (\approx 4^{-15})$$

$$= 4^{56}$$

\Rightarrow Man bekommt die richtigen Antworten und einen kleinen Teil des Restes falscher Antworten.

Da aber der Rest so groß ist, bekommt man entsprechend viele falsche Antworten (hier 4^{56})

schon bei 3 Ebenen
rekursiv wird es sehr
unübersichtlich

Bloomfilter

Ausnahmemenge

Bloomfilter

Ausnahmemenge

etc.



alles außerhalb davon
kann ich ignorieren

⇒ sehr aufwändig / rechenintensiv!

Beim Erstellen gucke ich noch mal alle Mengen durch!

Wie kann ich den Bloomfilter exakt (exakte?) machen, ohne dass ich meine Fehlerrate so extrem niedrig zu setzen.

⇒ clevere Methode

Ausnahmemenge komplett zu speichern ist Quatsch. Die Ausnahmemenge besser so speichern!

Dadurch kann man wieder mit kleineren Fehlerraten (z.B. 10^{-3}) arbeiten.

Aufg. 1

mit Bloomfilter rümpfen

wie groß werden Ausnahmemengen?
Lohnt sich dann noch ein Bloomfilter?

Mit Bloomfilter kann ich mir bis 1 zählen. Wie kann ich das trotzdem lösen (l-mer von häufig und selten unterscheiden)?

Mehr listen bräudchen mehr Spielere. Geht es auch mit nur einem Bloom-Filter?

⇒ Sampeln / Ausdünnen der Menge

↳ nur ein Teil nehmen (5% → bei ~~1000~~ $n=30$)

dann mit Bloomfilter
⇒ Rest überprüfen und ausdünnen

⇒ Mit Formeln spielen

Auf Webseite wird Quellcode und weitere Fragen
reingestellt.

CO-VL

Lighter-Paper

→ da keine gelesen hat → keine Besprechung

Genom Assemblierung

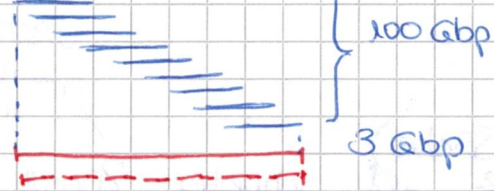
Seq.: viele DNA Reads (Länge 100-500 bp) (oder 1000)

Ges.: Genom, aus dem die Reads stammen

(Genom: mehrere Chromosomen, linear oder zirkulär)



Idee: hohe Abdeckung (coverage), überlappende reads



DZW bei Chrom.:

Zahlen bsp für Mensch

Reads sind in der Datei natürlich nicht sortiert!

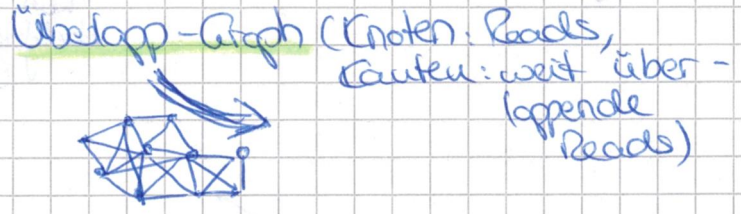
Viele reads, sind lang und haben Fehler.

Heute kein Problem was Speicherplatz etc. angeht. Humangenom ca. 1998-2000 assembliert. Zwar waren die reads länger (S.O.) aber keine so große Rechner, Festplatten etc. Damals riesen-Problem!

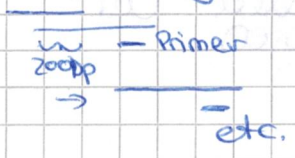
Ideensammlung:



Klassische 90'er



Um Lücken zu füllen:
Prime degenet



- ~ 100 Mio. Knoten
- ~ vielfaches an Kanten (> 1000x)
- 100 Mrd Kanten
- ~ (100 Mio)² / 2 Überlapp-Tests (naiv)

=> früher wenig Sequenzen, die länger waren, darunter geht das

"modern" 2000er

De Bruijn Graph (Knoten: ^{solide} k-mer
Kanten: (k+1)-mer)



~ 3 Mrd. Knoten
Kanten implizit

s. Lighter-Paper

in Vorlesung
näher betrachten

neue Ideen 2015+

Überlapp mit Hashing

Fluiddichtebeibehaltendes Hashing
(engl. locality sensitive hashing (LSH))

Datenstruktur für k-mer:

- sortierte Liste (binäre Suche)
- Hashtabelle (mehr Platz, ggf. schneller Zugriff)
viel Speicher
- Bloom Filter + Ausnahmeliste (Bloom Filter...)
(Zugriff ok, wenig Speicher)

wenn man es
wirklich gut
macht

↳ es gibt Bloom Filter, die codefreundlicher sind;
im Lighter-Paper unter Supplement

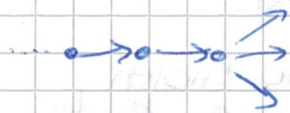
Variante: Blockweise Bloom Filter
(schneller, aber größer)



Idealfall, es gibt nur einen Weg
(die Krl. das das Chromosom
zu Ende) ist "vergessen" wir jetzt einfach)

kein
Nachfolger

↳ Lücke in Sequenzdaten
oder Chromosom zu Ende
(lokal gesehen zu Ende)



vernünftige Ausfälle im Genom
mit versch. Nachfolgen

im 1. Schritt an allen Ver-
zweigungspunkten aufhören

3 Schritte: \rightarrow v. contiguous (?) $\hat{=}$ zusammenhängend

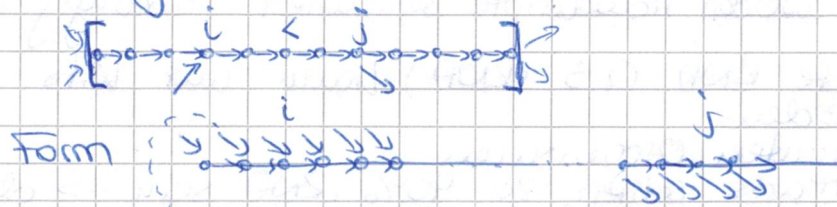
1. Contig-Assemblierung
2. Scaffolding (Verbinden von Contigs)
3. Gap filling

1

Contig: Sequenz, die sicher im Genom $\geq 1 \times$ vorkommt
(Es gibt versch. Klassen von Contigs. Hier werden 2 erwähnt.)

a) Unitigs: maximale Pfade mit inneren Knoten mit Ein- und Ausgrad 1.
(eher kurz)

b) Overlappungs



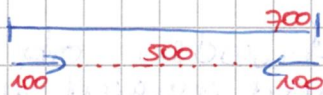
maximal lang, bis es $i < j$,

i eingehend verzweigt
j ausgehend verzweigt gibt,
mit Pfad $j \rightarrow i$

Overlappungs sind länger als Unitigs

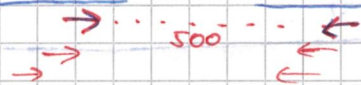
2

Fragmente = Paare von Reads



Contig 1

Contig 2



} nicht nur ein Read-Paar, sondern ganz viele, die es konsistent halten

Contigs dürfen mehrfach verwendet werden!

\Rightarrow Davaus sog. Scaffolds bauen

„Puzzlespiel“ auf höherer Ebene, wobei die Lage-Information genutzt wird.
Es können Lücken bleiben!

Ausgabe: Scaffolds (mit Glück ganze Chromosomen, aber eher kleiner) (17)

→ https sichern, ssh-link klappt nicht
 ↑
 man braucht ssh-key

Email mit Gitlab-Benutzername schicken

Review-Artikel nicht frei zugänglich → v.a. Bible
 ⇒ muss man nicht lesen; gibt einen guten Eindruck

Vorschlag für Projekt: RNA-Assemblierung

- nach dem Splicing meist 3-10 kleine Teile (alles Exons)
- mature mRNA → cDNA → doppelsträngige DNA
 ⇒ Sequenzierung
 ↳ es werden Transkripte sequenziert (abhängig v. Anzahl)
- ~~un~~ unerwünschte RNA (z.B. rRNA) kann mit Kits reise-geprüft werden.
- bei unbekanntem Organismus
 → solide Reads sollten zu 90% rRNA sein → damit Kit bauen lassen, Rest wie oben

~60.000 Gene insg.
 25.000 davon Proteinkodierend → führt zu mehreren Proteinen

Aus RNA Omnitigs machen. (Kann man mit „Hausmittelchen“ machen. Man braucht keine fancy Datenstruktur.)

→ Datensatz v. Pilz o. Weisbüschelaffen (rel. groß) können zur Verfügung gestellt werden.

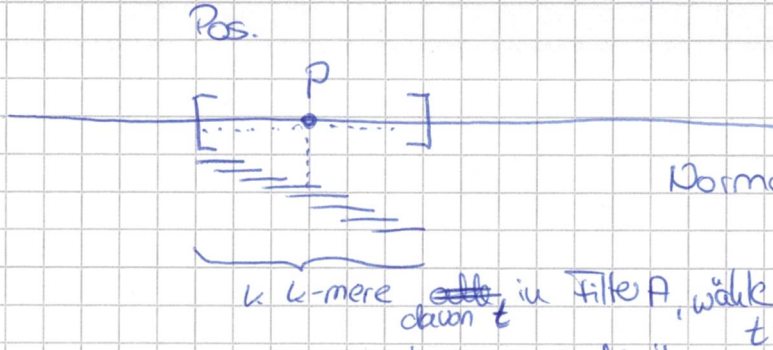
Jedes Contig müsste untersucht werden, ob ähnlich zu bekannten. (Müsste hier nicht gemacht werden.)

⇒ Sinnvolles Projekt mit vernünftiger Dateigröße

Lighter-Paper:

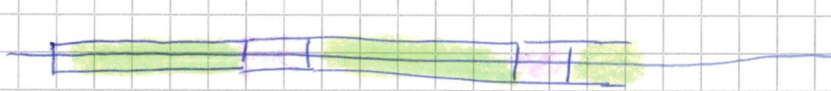
Schritt 1: k-mere anzählen

Schritt 2: Positionen im Read anzählen



Normalfall: wenn in Umgebung alles i.O. ist, sollte alle k -mer im Bloomfilter A sein

$t = \lceil \frac{k}{z} \rceil$
 $\rightarrow p$ "vertrauenswürdig",
 sonst "nicht vertrauenswürdig"



$\overset{grün}{\text{}} \stackrel{!}{=} \text{vertrauenswürdig}$
 $\overset{rosa}{\text{}} \stackrel{!}{=} \text{nicht}$

Nimm wir k -mer, bei denen alle Positionen vertrauenswürdig sind.

\Rightarrow Filter B (oder hier rauskommt) ist ziemlich gut.
 Es ist möglich das in Filter B k -mer sind, die nicht in A sind.

Auch wenn ich zähle, kann es sinnvoll sein. Wenn ich sample, dann brauche ich das.

Schritt 1 ist I/O-dominant. Rechenarbeit im Vgl. zum I/O sehr wenig. Durch Parallelisieren wird es nicht (wirklich) schneller.

Schritt 2 wie geht es sich zu parallelisieren
 Problem hier: schreiben des Bloom-filters B. Hier muss synchronisiert werden.

Auf jeden Fall kann man Positionsklassifizierung parallelisieren.

Es wird kein normales Bloomfilter, sondern Pattern-blocked Bloom filter genommen (Supplementäre Note 1)
 - Bitmuster setzen

Projekt-Idee (?): vlt. andere Bloom-filter ausprobieren und vergleichen
 \rightarrow analysieren, Eigenschaften nachweisen, zeigen das es in Praxis geht

H Hausaufgabe: Tomascu-Folien oder -Paper angucken
 \rightarrow Omnitags \rightarrow auf jeden Fall

Review-Paper optional

\rightarrow eher nicht formal (eher für Biologen)
 \rightarrow kann man schneller lesen (Text)

$\Rightarrow \log_2 k$ an Sven schicken

Genom Assemblierung

(vom biologischen zum formalen Problem)

1. Gegeben (Eingabe): Strings über $\{A, C, G, T\} =: \Sigma$
 $\{s_1, \dots, s_n\}$ IUPAC Alphabet (15 Zeichen) $=: \Gamma$
 $\Sigma \times \{0, \dots, 40\}$ Strings mit Qualitäts-
 $\Gamma \times \{0, \dots, 40\}$ werten

Man könnte noch was über Anzahl und Länge sagen

2a. Gerüst (Ausgabe): String G über $\{A, C, G, T\}$, so dass
 gilt: \forall oder rev.comp (rs) Teilstring von G $\forall i = 1, \dots, n$
 und G ist minimal lang mit dieser Eigenschaft
 "shortest common superstring"
 im Wesentlichen
 \rightarrow i.w. Genome ohne lange Wdh
 \rightarrow Genom enthält auch fehlerhafte Reads



So was das lange gerüstet!

NEU:

1. Gegeben (Eingabe): a) Strings über $\{A, C, G, T\} =: \Sigma$,
 Parameter $k \in \mathbb{N}$
 b) Menge von k -meren über $\{A, C, G, T\}$
 ["solide k -mere"]
 a) enthält noch etwas mehr Kontext

2a. Gerüst (Ausgabe): De Bruijn Graph (V, E) , $V: k$ -mere aus $\{s_i\}$
 $E: (k+1)$ -mere aus $\{s_i\}$
 Pfad "G", der jede Kante ≥ 1 mal enthält
 ? Pfadmenge Π , so dass jede Kante in ≥ 1 Pfad enthalten
 minimaler Cardinalität

Es gibt mehrere (mind. 20) sinnvolle Formulierungen für Genom Assemblierung
 \rightarrow lernen sodass formulieren! (Klassische hat \mathcal{O} Faktorisierung einer Polynomele
 Laufzeit) bzw. Approximation

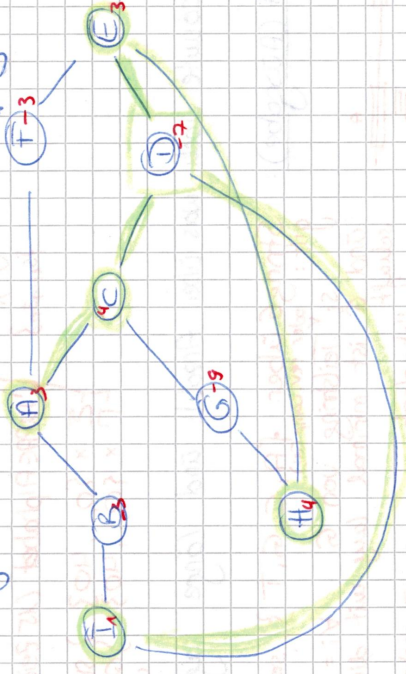
Protein-Interaktionsnetzwerk

Knoten: Proteine
mögliche

Kanten: physikalische Interaktionen

Knotengewichte $\in \mathbb{Z}$ $\begin{cases} > 0: \text{erhöhtig} \\ < 0: \text{erniedrigend} \end{cases}$

$$y = (y, v)$$



⊗ Was dieses Protein in einer Klasse von kranken Patienten auffällig im Vgl. zu einer Klasse von gesunden Patienten

Gesucht:

relevante Netzwerkmodule

- Steinerbaum zwischen allen auffälligen Knoten

↳ wenigsten Kanten

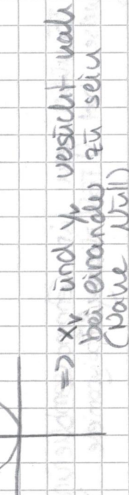
- maximum scoring connected subgraph (+ alle pos. Knoten enthalten)

- thresholdgraph (> 0)
connected components $\geq m$

$$F(x) = \sum_{i \in V} (x_i - y_i)^2 + \lambda \sum_{e \in E} |x_i - x_j|$$

→ lokal konstante Lösung x

Erklärung in Abstraktion



→ x_i und y_i versucht nah beieinander zu sein (Naher Null)

NP-schwer

NP-schwer

relativ populär geworden

benutzbare Knoten

konvexes Optimierungsproblem in x

(ähnlich mit vektorielles wie lineares Problem wird schwieriger)

⇒ immer polynomial lösbar

Metagenomik

Analyse / Assemblierung mehrerer Genome in einer Probe

- z.B. Haut-Bakterien
- Stuhlproben-Bakterien (Darmflora)
- Umweltproben - z.B. Gewässer-Organismen

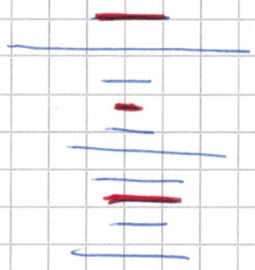


Für jeden Read:
 - finde Spezies
 - finde Genfamilie

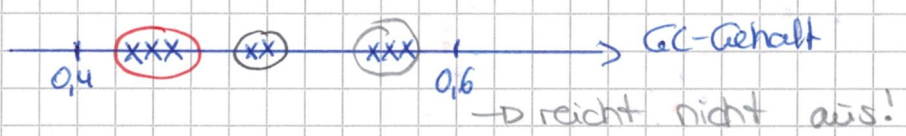
Gleichzeitige Assemblierung aller Genome
 Zuordnung Read \rightarrow Spezies nicht bekannt

1. Assembliere Contigs (z.B. unitigs) wie bei der einfachen Genom-Assemblierung
 (Contigs $\hat{=}$ 1000-10000 bp)

Metagenomikspez. Problem \rightarrow 2. Gruppierung Contigs \rightarrow Genome



die roten gehören zu einer Spezies!
 das ist aber nicht bekannt!



Contigs \rightarrow Genome

\triangleright Eigenschaften der Contigs:

a) Tetranukleotid-Häufigkeit

Anzahl ~~4~~ 4-merer: $4^4 = 256$
 $4^2 = 16$ selbstkomp.

Gleiche Organismen \rightarrow gleiche Tetranukleotidvektor

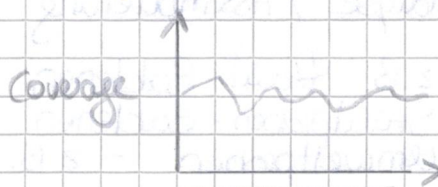
$$\frac{256-16}{2} + 16 = 136$$

Contigs $\rightarrow TR^{136}$

$C \mapsto x = (x_1, \dots, x_{136})$ mit $x_i \geq 0 \quad \forall i$

$$\sum_i x_i = 1$$

b) Durchschnittliche Abdeckung aller k-merer, aus denen das Contig assembliert würde

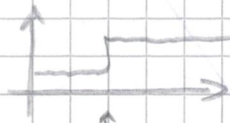


Abdeckung der Häufigkeit der k-merer in einem Contig

Idee:

Contigs aus dem gleichen Genom haben approximativ

- a) gleiche Tetranukleotid-Häufigkeit
- b) gleiche Abdeckung ("coverage")



würde man hier trennen, da man vermutet, dass es 2 versch. Organismen sind

Contig $\hat{=}$ Punkt in \mathbb{R}^{136+1}

haben eine andere Dimension als

deswegen steht da auch "nicht 137!"

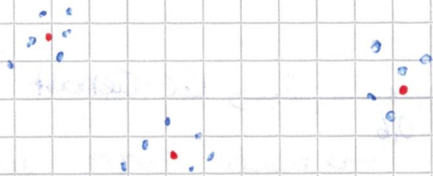
Clustering-Verfahren:

- Lloyd's Algorithmus ("k-means")
- DBScan (density based) (Dichte-basiert)

•• EM-basierte Verfahren (probabilistisches Clustering, expectation maximization)

k-means

Punkte in \mathbb{R}^d , x_1, \dots, x_n



Finde k Clusterzentren c_1, \dots, c_k

② Zuordnung $x_i \rightarrow c_j(i)$, so dass Summe der Distanzen $d(x_i, c_j(i))$ minimal ist

- 2 Probleme: ① Zentren
② Zuordnung

np-schweres Problem (beide Probleme gleichzeitig)

→ Einzelproblem gegeben das andere ist einfach/trivial

Zentren gegeben → nächster Nachbar
Zuordnung gegeben → in der Mitte

Lloyd's Algorithmus ist "schlau" und löst die Probleme immer abwechselnd. Der garantiert aber nicht das globale Optimum zu finden. Es kann sein, dass er in einem lokalen Optimum hängen bleibt.

Man muss initiale Werte vorgeben.

Wie macht man das? Man macht das ganze 100x ~~mal~~ mit versch. initialen Werten und wählt das beste Ergebnis.

Lloyd's Algorithmus ist für die Praxis wirklich gut und löst das Problem in polynomieller Zeit, er löst aber nicht das k-means-Problem da er nicht garantieren kann das globale Optimum zu finden.

• Allgemeiner Ansatz für probabilistisches EM-Clustering

- Struktur wie bei k-means:

Initiale Zentren

→ probabilistische Zuordnung der Punkte

→ neue Zentren

← spezifisch für unser Modell

$$P(x_i | c_j) = P_{\mu}(\dots) \cdot P_{\sigma}(\dots)$$

mit d-dimensionalen Normalvert. = $c \cdot \exp\left(-\frac{\|x_i^I - c_j^I\|_2^2}{2\sigma^2}\right) \cdot \text{Poisson}_{c_j^II}(x_i^{II})$

Konstante, da Sven grad nicht weiß, was der richtige Normalisierungsfaktor für d-dimensionales Clustering

mit $x_i = (\underbrace{x_{11}, \dots, x_{136}}_{x_i^I} | \underbrace{x_{137}}_{x_i^{II}})$ → Poissonvert. für Coverage

Berechne $\log P(x_i | c_j)$ für alle c_j für alle x_i

Berechne Zugehörigkeiten:

$$w_{ij} = \frac{P(x_i | c_j)}{\sum_j P(x_i | c_j)}$$

→ $\sum_j w_{ij} = 1$ für alle i .

(relative) Es gibt mir für alle i die Verteilung der Zentren an

Man kann den Alg. auch begreifen als (offensichtliche) Verallgemeinerung des Lloyd's Algos.

⊗ man könnte auch σ_j^2 nehmen, aber das macht es auch schwerer.

⇒ Wenn man das alles so macht, dann funktioniert Meta-genomik relativ gut. v.A. auf simulierten Daten (da kann man es überprüfen)

Metagenomik

alle Referenz-
genome bekannt

Referenz-Genome
unbekannt

xenografts

Spezialfall

Menschl. Tumor
in Maus

- > contigs
- > Eigenschaften
- > Clustering

Sequenzierung: Mischung von
Reads aus Mensch
+ Maus

- Einteilung:
- 1) Mensch
 - 2) Maus
 - 3) beides gleich gut
 - 4) kein oder beidem (wenig)
 - 5) widersprüchlich (teilw. 1, teilw. 2)

schnell:
Hashing,
z.B.
Cuckoo hashing

z.B. Mensch hat ca. 1,4 Mrd 25-mer in nicht repetitiven
Bereichen.

Wie viele Bits sind minimal nötig, um Mengen der Größe
1,4 Mrd von 25-meren darzustellen?

Nair: 25-mer $\hat{=}$ 50 Bits
-> 1,4 Mrd \cdot 50 Bits

Es gibt $\binom{4^{25}}{1,4 \cdot 10^9}$ verschiedene mögliche Mengen.

Bräuche $\geq \log_2 \binom{4^{25}}{1,4 \cdot 10^9}$ Bits. $\log_2 \binom{n}{k} = \log \frac{n!}{k!(n-k)!}$

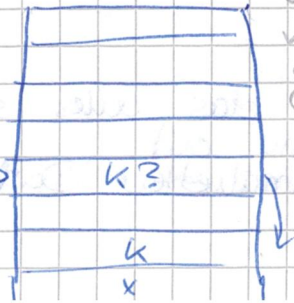
$$= \log n! - \log k! - \log [(n-k)!]$$

$$\log 1 + \log 2 + \dots + \log n \leq n \log(n)$$

Hash:

Schlüssel k 50 Bits

$h(k)$



alle Schlüssel müssen
vorkommen und es
gibt leere Felder

$$k \mapsto (h(k), r(k))$$

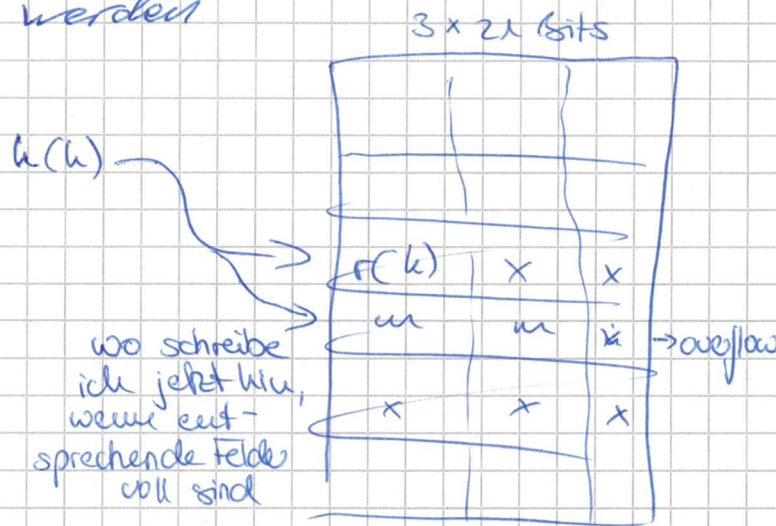
Adresse/
Seite i.d.
Hash-Tabelle

"Rest"
fingerprint

so dass $\exists f$ gibt mit $f(h(k), r(k)) = k$

beispielsweise Zerlegung: $h(k)$: erste 28 Bits
 $r(k)$: restl. 27 Bits

→ gemeinsame Anfänge müssen nur einmal gespeichert werden



Das was überläuft sammel ich (wenn gut, dann) 5% der alten Menge)

→ 2. Tabelle, ggf. 3. Ebene

⇒ breitere Tabelle:

weniger überläuft, mehr verschwendeter Platz

⇒ schmalere Tabelle:

genau umgekehrt

→ passende Einstellung finden

Variantedetektion → in 2 Wochen



Hauptproblem:

techn. Probleme von Varianten unterscheiden

⇒ statistisches Thema

dadurch: RNA-Sequenzierung (quantitativ) ⇒ statistisches Thema (27)

CO-VL

von Johannes:

DEFECTINA GENOMIC VARIANTS

(Folien von 2013)

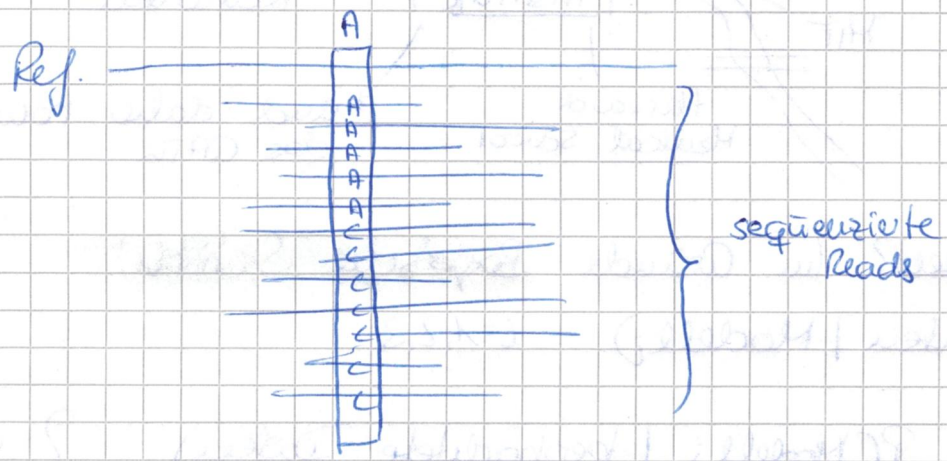
- DNA Referenzgenom: Durchschnitt aus zufällig ausgewählten Personen
- SNV $\hat{=}$ single nucleotide variant
- Indel $\hat{=}$ insertion or deletion
- SNP $\hat{=}$ single nucleotide polymorphism (in a population)
 \Rightarrow muss beobachtbar sein, also nicht im Einzelfall
- Variante: homozygot: gleiche Allele (\neq Ref.)
heterozygot: verschiedene Allele ($=$ Ref, \neq Ref.)



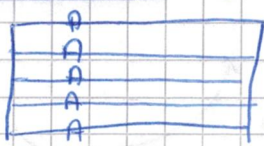
m $\hat{=}$ maternal
p $\hat{=}$ paternal

homozygot (aber keine Variante):	beide Allele = Ref	0
heterozygote Variante	1 Allel = Ref, 1 Allel \neq Ref	1
homozygote Variante	beide Allele \neq Ref (aber gleich)	2
		↑ kodierung

Genotyp Bsp.: 0001112200...



\wedge zu erwarten: max. 2 verschiedene Basen

Probleme: - korrelat mappen
- PCR-Dublikate:  nicht unabhängig!
eigtl. 1 Read, 4 Kopien

Warum PCR-Dublikat, wenn zufällig gleiche lang
- Fehlerrate heute eher $< 1\%$, dennoch vorhanden

A und O: tief genug sequenzieren!

- > teuer (Cost nahezu proportional zur Sequenzierhöhe)
- > zum sparen nicht alles sequenzieren, sondern Exom (bzw. etwas drum herum) (s. Bild #7) Coverage Histogramm

#9 GATK = genome analysis tool kit

* aus typischer Exom-Analyse ~ 50.000 Varianten

Qualitätswerte in Reads:

Q	10	20	30
Feehler - w'kt	10^{-1}	10^{-2}	10^{-3}

(PHRED-Skala)

$10^{-Q/10}$

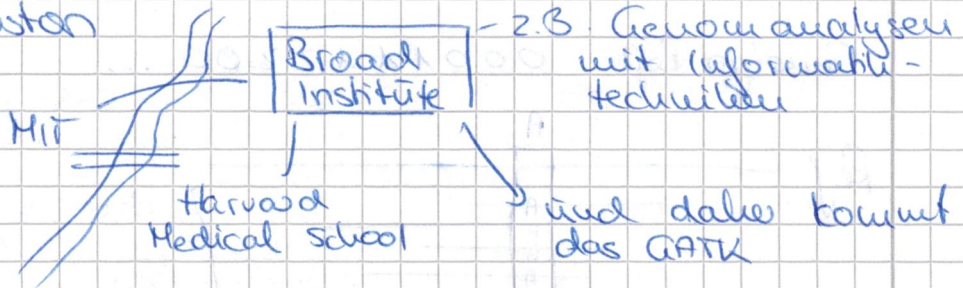
98% ab bp

Qualitätswerte kommen aus dem Sequenziergerät

→ alles ab hier ist zuverlässig (erleichtert die Prozedur, da man nicht mehr die Qualitätswerte beachten muss)

▷ GATK (ab # 18?)

kommt aus Boston



Was die wollen? (in Gründe Bayes'sche Statistik!)

(leicht: $P(\text{Daten} | \text{Modell}_i)$ $i=1,2,3$)

(für jedes Modell)

Wir wollen aber: $P(\text{Modell}_i | \text{beobachtete Daten})$

$$= \frac{P(\text{Daten} | \text{Modell}_i) \cdot P(\text{Modell}_i)}{P(\text{Daten})}$$

Satz von Bayes

dann statt = \propto (proportional)

Konstante: durch 1 ersetzen

=> wir wollen ja nur das Modell bestimmen

#19

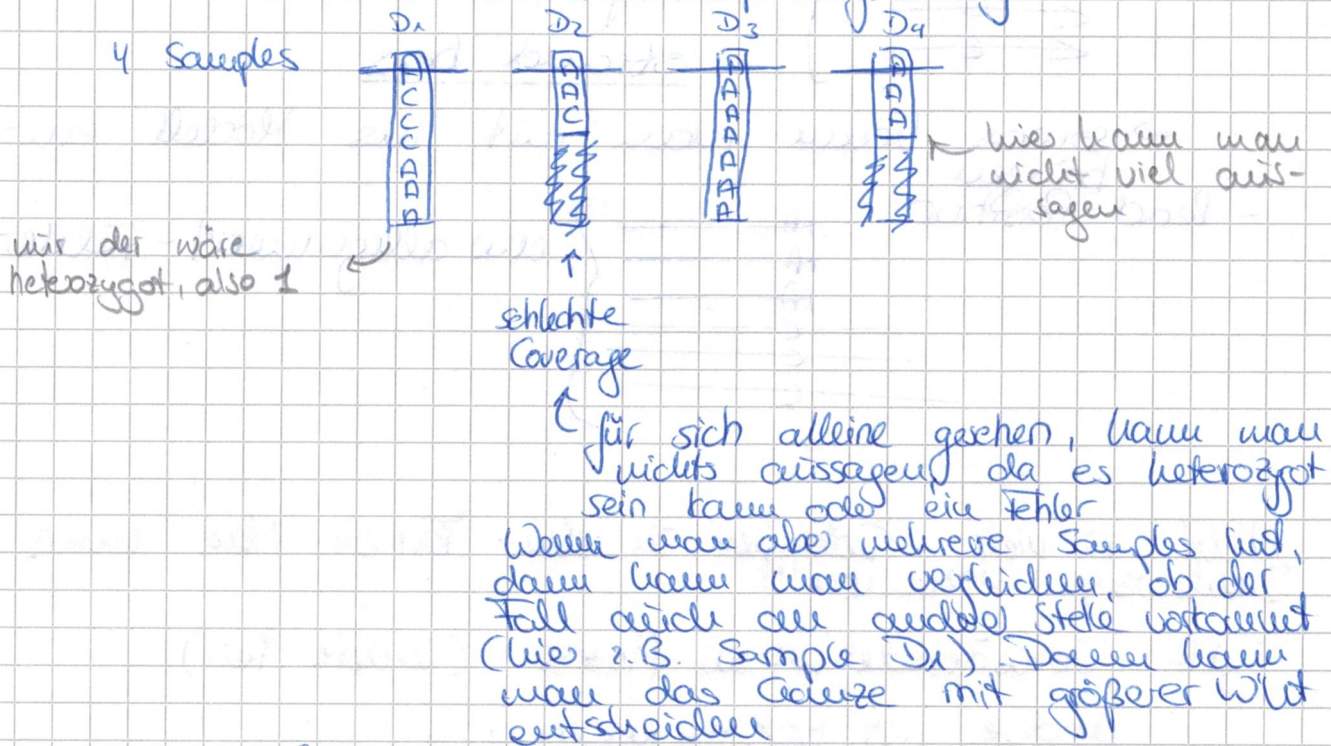
- mehrere Patienten i
- alliierte Liste von Basen D_i ; z.B. ACCCTA...
- W'kt in Liste i die j . Base X ist
- Tabelle ist empirisch ermittelt
"Wenn ein C steht und ich einen Fehler mache, wie Wahrscheinlich ist dann ein A, T oder G?"
⇒ Sehr findet das "nicht so wahnsinnig wichtig"

#20

- Formel: $\frac{1}{2}$ W'kt linker Allel + $\frac{1}{2}$ W'kt rechter Allel

#21

- GATK macht das über alle Samples gleichzeitig



mehrere Samples gleichzeitig

#22

- alle Samples übereinander legen und sagen, ob es überhaupt eine Variante gibt
- für Gesamtw'kt muss man über alle gehen
- sehr viele Möglichkeiten, aber man kann es ausrechnen
- für jedes m eine vernünftige prior-W'kt, die sich aus evolutionstheoretischen Modellen zusammensetzt

$$P(\text{krank} | +) \propto P(+ | \text{krank}) \cdot P(\text{krank})$$

$$P(\text{gesund} | +) \propto P(+ | \text{gesund}) \cdot P(\text{gesund})$$

$$P(\text{krank} | +) = \frac{P(+ | \text{krank}) \cdot P(\text{krank})}{\sum_{m=\text{gesund, krank}} P(+ | m) \cdot P(m)} \quad (= P(+))$$

$$P(\text{krank} | +) \propto 0,999 \cdot 10^{-6} \approx 10^{-6}$$

$$P(\text{gesund} | +) \propto 0,02 \cdot (1 - 10^{-6}) \approx 0,02$$

$$\frac{10^{-6}}{10^{-6} + 0,02} \approx 0 \quad \frac{0,02}{10^{-6} + 0,02} \approx 1$$

dieses Verhältnis bestimmt Ergebnis

$$\text{also } P(\text{krank} | +) \approx 0 \\ P(\text{gesund} | +) \approx 1$$

Selbst bei positivem Test, bin ich ziemlich sicher gesund!

Bayes schaut in diese Richtung →

1 Mrd. Bevölkerung	krank	gesund
+	999	20 Mio
-	1	980 Mio
	1000	~ 1 Mrd

absolute Zahlen am Bsp. von 1 Mrd. Menschen

⇒ es lohnt sich (dem Schema (Bayes) zu folgen, da sehr intuitiv!

- Zufällige Meersch:

Was ist wahrscheinlicher? Informatiker oder Informatiker, der Linux benutzt?

$$P(A \cap B) \leq P(A)$$



Bücherepfehlung: Thinking Fast and Slow

Ich kann jetzt Modelle bauen wie z.B. wenn da ein C steht, geht das Rad mit hoher W'it nach rechts. Eigentlich müsste das unabhängig sein. Wenn aber so ein Modell gewinnt, dann weiß ich, dass da etwas nicht stimmt!

$P_1 (C=c | (X,Y) = (x,y)) \rightsquigarrow$ unsere Tabelle

$P_2 (D=d | C=c, Z=z)$

erweitern: $P(C=c, D=d, E=e | (X,Y,Z) = (x,y,z))$

Z ist Bayes??? Richtung-Bayes?

Z: -1: eher \leftarrow bei kleinem c

0: 0,5

+1: eher \rightarrow bei kleinem c

: usw.

Ich habe viele verschiedene Modelle und für alle kann ich in der Vorwärtsrichtung meine Daten auswerten. Und dann kommt Bayes und bewertet meine Modelle. Und da ist dann das Z? Wenn Richtung-Bayes mehr bringt als ohne, dann stimmt da etwas nicht!

\rightsquigarrow Diese Modelle kann man beliebig kompliziert machen!

\Rightarrow großer Rechenaufwand

Ober die Parallelisierung ist es sehr gut. (Trivial Parallelisierbar) Man kann z.B. pro Kon eine Stelle gleichzeitig betrachten.

HA Folien im Git Modelle noch mal angucken!

Varianten-Detektion

Modell: $P(\text{Beobachtung} | \text{Genotyp})$

a-Prion: $P(\text{Genotyp})$ (Ref. am wahrscheinlichsten)

Inferenz mit Bayes: $P(\text{Genotyp} | \text{Beobachtung})$

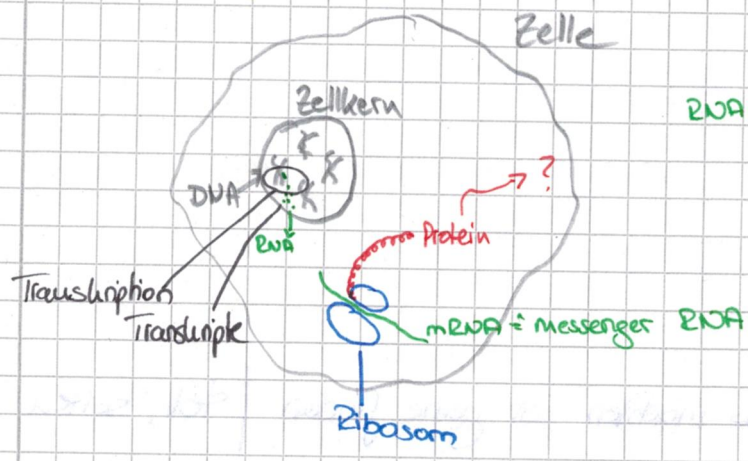
Probleme: sehr komplexe Komponenten
sehr große Zustandsräume

→ keine exakten Rechnungen → Approximationen
Markov Chain Monte Carlo (MCMC)

→ dient i.d.R. hier etwas auszurechnen indem die Werten in einer Simulation berechnet werden, da man sie nicht mehr exakt ausrechnen kann

von letzten Mal
Zusammenfassung

▷ Transkript-Expression



RNA kann Zellkern verlassen

~ 20.000 - 25.000 Protein-Codierende Gene

~ 65.000 Gene

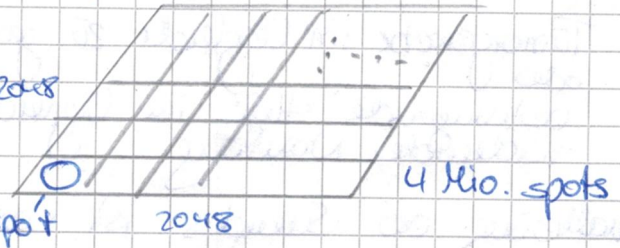
Man misst wie häufig jedes Gen abgelesen wird (transkribiert wird) und bestimmt dann in welchem Zustand die Zelle ist.

Technologien: (Transkripte zählen, Gene zählen)

- DNA Microarrays
- RNA-seq (S.S. 38)

(Farbstoff, radioaktiv)
Label
einzelsträngige DNA
(25 nt oder 60 nt)
RNA über Chip greifen

Komplementäre Sequenz



viele Mio. Kopien desselben Moleküls

Luftigkeit des Ländchens entspricht der Menge des Traces-
 Wipps, also starkes Ländchen: es gibt viel was leben bleiben
 kann, etc.

Firma: Affymetrix 25-mer
 Agilent 60-mer
 (50-70-mer)
 ~60 €

} spezifisches, aber wenn
 braucht mehr DNA,
 es kann sich mehr ver-
 wedern

=> sie können nicht die volle Bandbreite abdecken (ggf. Spot
 übersättigt), bzw. sind sehr ungenau

=> vor 157. Interessantes Bioinfo-Problem: Designen der Arrays
 (eindeutige DNA-Sequenzen, die gleiche thermodyn. Eigenschaften
 haben (Gc-Gehalt).

RNA-seq

sequenzieren alle RNA-Moleküle (Umschreiben in cDNA),
 Gen zuordnen,
 Gene zählen

=> Beides hat ähnliche Probleme

Ausgabe: Matrix

x Samples y

Objekte (Chromosom, Gene)	16-bit Array-Werte	
	Zähler	
	9999...	999 000 ... 00
	1000-1500	2000-3000
100	100000	
4 Mio	Klasse ⊕	⊖

← was möchten wir gerne finden / toll, selten

interessant

↳ z.B. Tumorgewebe im Vergleich zu gesundem Gewebe
 oder
 verkümmerte Hefe im Vergleich zu Hefe mit aus-
 reichender Nahrung

Problem: Skalierung der Samples ist verschieden

Normalisierung:

- "house keeping genes", relativ überbust (Gene mit stets konstanter Expression)

- Division durch Summe / (je Sample Spalte)
Mittelwert

a	1,5a
b	1,5b

- Division durch Median
↳ ignoriert Ausreißer

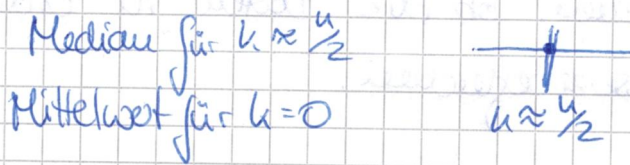
c	⋮
d	⋮
⋮	⋮
x	⋮
y	⋮
z	10k · z

- "..." gewichteter Mittelwert: ⊗

Gez. X, sortiere X, $X_{sorted} = (x_{(1)}, x_{(2)}, \dots, x_{(n)})$

Σ 1	2
-----	---

$$\mu_{(k)} = \frac{1}{n-2k} \sum_{i=k+1}^{n-k} x_{(i)}$$

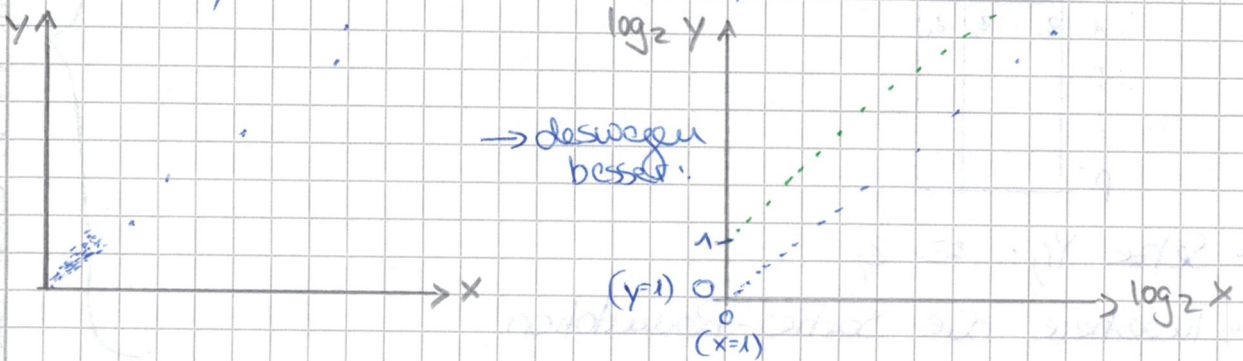


⊗ weis nicht symmetrisch sein, sondern man kann von 70-80%



gesehen, da sehr viele Gene gar nicht exprimiert werden und damit bei 0 liegen. (ca. 50%). Damit hat man die Gefahr, dass als Median 0 berechnet.

x und y s. S. 38 in der Matrix



→ deswegen besser!

in der Praxis natürlich nicht so ideal

$$y = 2x$$

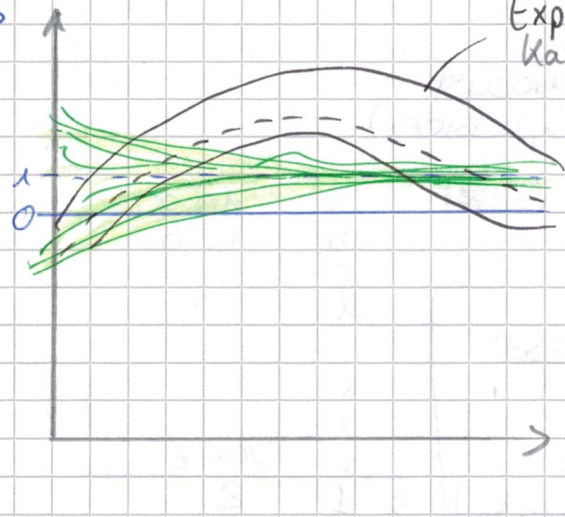
$$\begin{aligned} \lg y &= \lg 2 + \lg x \\ &= 1 + \lg x \end{aligned}$$

grüne Kurve ⇒ macht einen Offset von 1

⇒ jedes konstante Faktor macht einen Offset

M
minus
log y -
log x

Expression ist verschieden verteilt.
Kann auf techn. Probleme hinweisen!



=> wenn man Glück hat, dann sieht man diese Wolke
=> wenn man Pech hat, dann sieht man eine solche Bauform
=> Parabel, dann „platt machen“, also zur Gerade transformieren

MA-Plots sind diagnostisches Tool, um 2 Spalten zu vergleichen.
Damit kann man sehen, ob die Daten in Ordnung sind.
=> Je gerade, desto besser zu vergleichen.

hart Normalisierungsmethode:

Geg.: Matrix X p Objekte \times n Samples
 $n \ll p$

- Sortiere jede Spalte für sich
→ Y
Jede Zeile in Y ist ein Quantil der Datensätze in X .
- Berechne durchschnittliche Quantile

$$z_i = \frac{1}{n} \sum_{j=1}^n Y_{ij}, \quad i = 1 \dots p$$

	Y		Z
1	z_{11}	$\dots z_{1n}$	z_1
\vdots			
p			

- Setze $Y_{ij} := z_j$
- Invertiere die Sortier-Permutation

→ das k-te Quantil ist in allen Samples gleich $\forall k$
Quantilnormalisierung

=> härteste, sinnvolle Normalisierungsmethode (die man noch verwenden kann)

→ $\frac{7}{p}$ -Quantil $\hat{=}$ der 7-leinste Wert jeder Spalte ist gleich
✓ (Median, Mittelwert, etc. ist auch überall gleich) (?)

=> ist sehr konservativ. Nach dem Normalisieren sehe ich vermutlich wenig Unterschiede, aber die die ich sehe sind relativ sicher

→ Problem einer Zellen-T...

(0-ü)

Mit Notebook (Python)

conda create --name weech8 python=3 numpy scipy (...)

nachinstallieren: ^{panda?} source activate weech8
conda install jupyter

=> s. Elias

▷ Mitschrift vom Programm hat Elias. Hier vier Notizen:

- df = data frame
- 11_Cf2 etc. sind sample Nomen
- describe liefert zu jeder Spalte Statistiken

→ mean = Mittelwert
→ std = Standardabweichung
→ 50% = Median
→ x% = Quantile

- "Magic-Commands" %matplotlib inline ~~###~~
↳ Notebook-spezifisch (läuft im Document)

=> Programmcode im Git

- Q-Q-Plot: Quantil gegen Quantil-Plot
d.h. mit sortierten Werten

- pandas, um nicht auf rohen Arrays zu arbeiten

Normalisierung über alle Spalten ist sinnvoll

1. $\frac{1}{x^2} = x^{-2}$
 $\frac{d}{dx} x^{-2} = -2x^{-3} = -\frac{2}{x^3}$

2. $\frac{d}{dx} \ln(x) = \frac{1}{x}$
 $\frac{d}{dx} \ln(x^2) = \frac{1}{x^2} \cdot 2x = \frac{2}{x}$

3. $\frac{d}{dx} e^x = e^x$
 $\frac{d}{dx} e^{2x} = e^{2x} \cdot 2 = 2e^{2x}$
 $\frac{d}{dx} e^{-x} = e^{-x} \cdot (-1) = -e^{-x}$

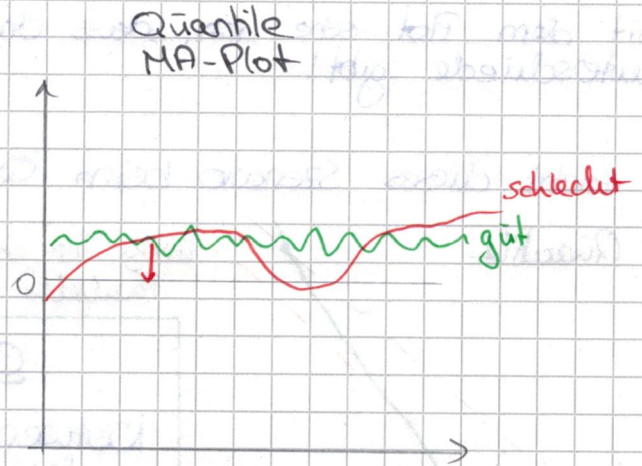
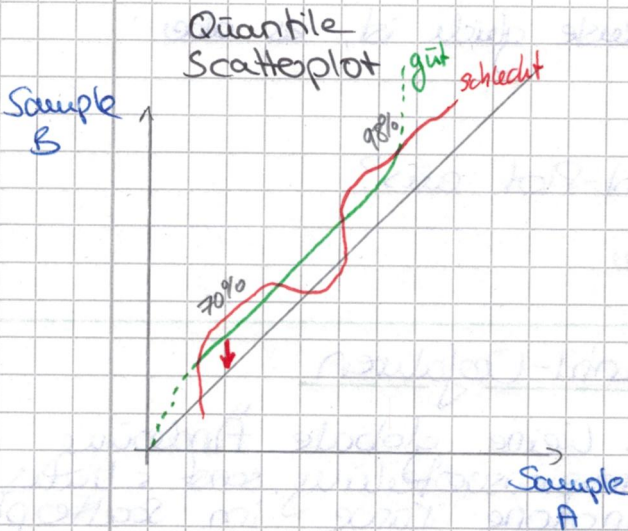
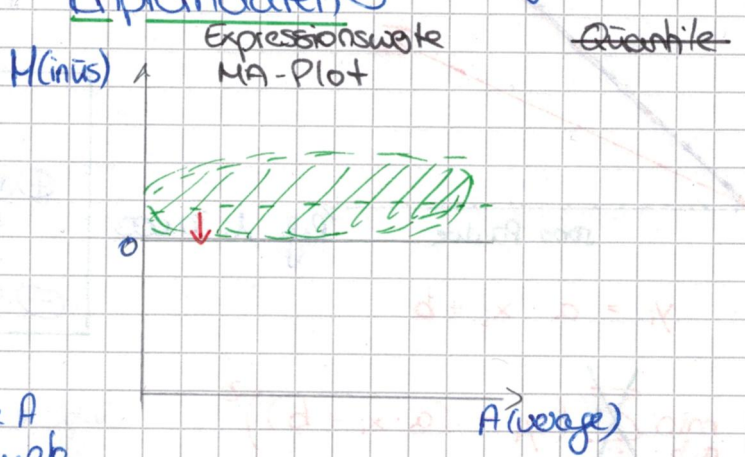
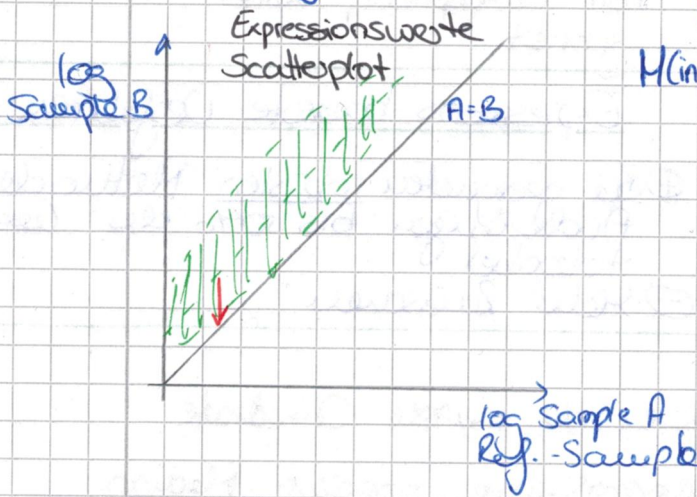
4. $\frac{d}{dx} \sin(x) = \cos(x)$
 $\frac{d}{dx} \cos(x) = -\sin(x)$
 $\frac{d}{dx} \tan(x) = \sec^2(x)$

5. $\frac{d}{dx} \arcsin(x) = \frac{1}{\sqrt{1-x^2}}$
 $\frac{d}{dx} \arccos(x) = -\frac{1}{\sqrt{1-x^2}}$

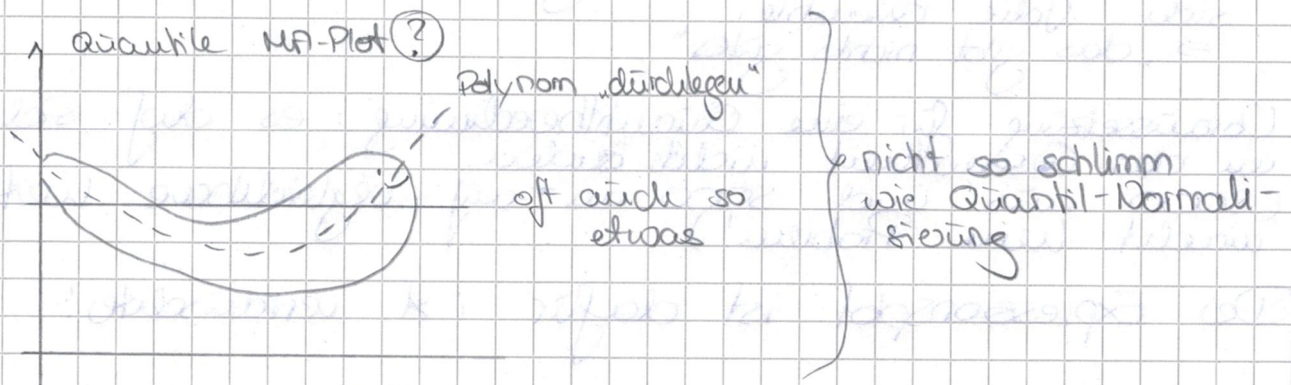
6. $\frac{d}{dx} \arctan(x) = \frac{1}{1+x^2}$
 $\frac{d}{dx} \ln(x^2 + 1) = \frac{1}{x^2 + 1} \cdot 2x = \frac{2x}{x^2 + 1}$

Wiederholung:

Diagnostische Plots / Normalisierungsplots für Transkriptomdaten

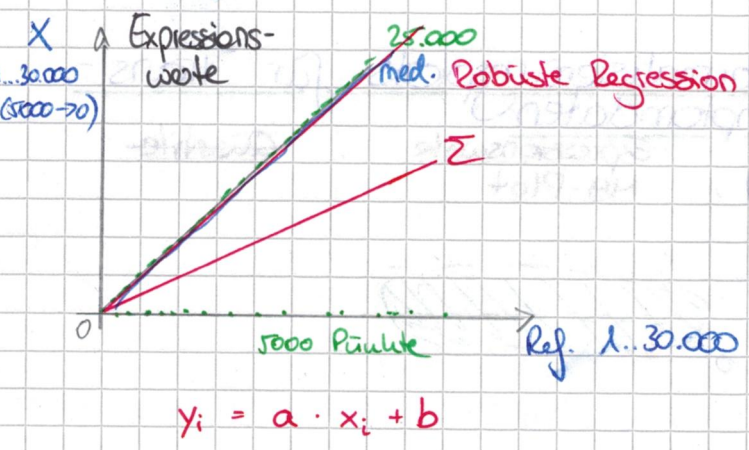


- zweckmäßig sich Referenzsample zu beschaffen
 - z.B. beste Qualität, meiste Sequenzierungen (= Beste)
 - Durchschnitt



- 70-98% der Quantile sollten schön gerade sein

Ich schalte 5.000 Gene ab (Sample X). Wie sieht Scatterplot aus?



→ Minderheit, die der Mehrheit widerspricht, wird ignoriert

Expressions-basierte Verfahren:

- ⊕ Bei geeigneten robusten Methoden Ausbügeln bis 50% des Gens tolerabel
- ⊖ Fehler rauschen

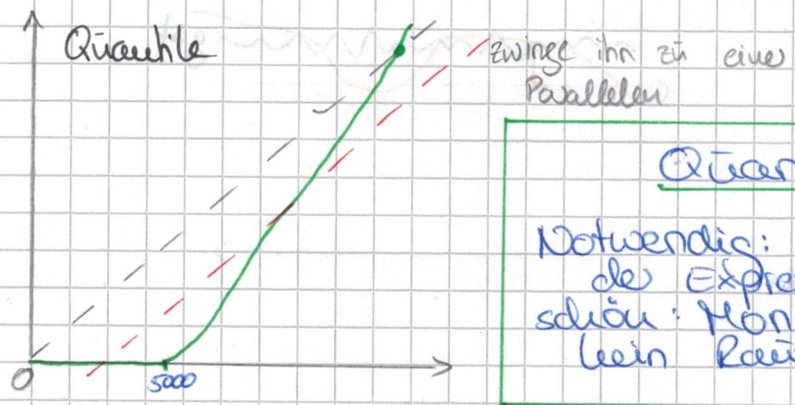
→ $\min_{a,b} \sum_i (y_i - (a \cdot x_i + b))^2$
~~median~~

Verfahren der kleinsten Quadrate

→ Nicht Quadratsumme, sondern Median
 ⇒ viel rechenaufwändiger

⇒ mit dem Plot sehe ich, dass das Meiste gleiche ist, es aber Unterschiede gibt!

Wie sieht dieses Szenario beim Quantil-Plot aus?



Quantil-Verfahren

Notwendig: Keine globale Änderung der Expressionsverteilung, sonst s. links
 schön: Monotone Kurve im Scatterplot, kein Rauschen

⇒ d.h. wenn ein Großteil des Gens stark ändert (z.B. 0 → hoch, oder auf 0), der Rest aber gleich bleibt, verschieben sich die Quantile
 ⇒ „das gibt nichts gutes“

Voraussetzung für eine Quantilberechnung: es darf sich an der Geschwindigkeit nichts ändern.
 (5000 aus und 5000 an + auf vergleichbare Werte macht kein Problem!)

Der Expressionsplot ist dafür i.d. verräuselt!

→ Beide Verfahren nutzen und Plots anschauen
 (Bevor man mit einer automatischen Pipeline draus geht, sollte man sich das anschauen. Deswegen macht man ja auch Plots.)

M z.B. verdoppeln. Aber Problem des kleinen Zahlen, deswegen addieren wir was dazu.

Sei $M := \max(F) + 10$

Betrachte $\delta := \frac{\Delta}{M} \geq 1$

⇒ Wenn die Lücke groß genug ist, dann könnte man auch die erste Bedingung wieder aufweichen, z.B. ein Wert δ auf in der Lücke liegend o.ä.

⇒ 1. Teil: Normalisierung → wichtig! Wenn man die verbodet ist der Rest sinnlos

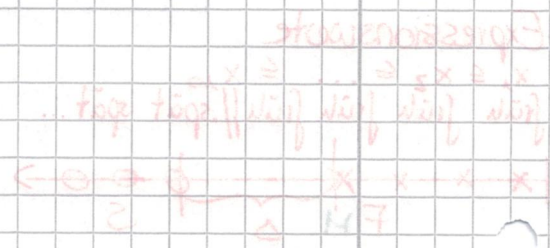
2. Teil: Testen auf Unterschiede → statistische Mittel

z.B. t-Test (nicht so gut)

Ränge (wird erst schrittweise, wenn man genügend samples hat)

bei 4-6 Samples kann man nicht über $\frac{1}{210}$ kommen. Mehr Samples werden aussagekräftiger

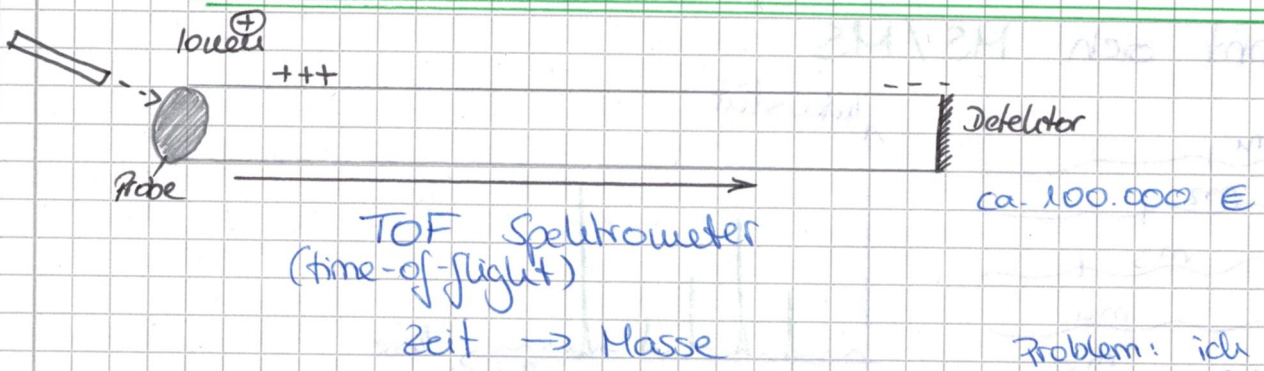
Anm.: wenn man zu strikt ist, findet man häufig 0 Gene, da es oft streut!



$$\frac{10}{10} = 1$$

$$200.0 = \frac{10}{0.05} = 200$$

PROTEOMIK / MASENSPEKTROMETRIE



Problem: ich brauche ein Vakuum



Proteine sind relativ große, schwere Objekte. So ein ganzes Protein behalvere ich überhaupt gar nicht überschreitet.
Ich kann Proteine zerbrechen.

Idee

Protein



Aucare Kette von Aminosäuren:

enzymatischer Verdau

viele Peptide (5-20 AS)

→ MS → Massen $\{m_1, m_2, m_3, \dots\}$

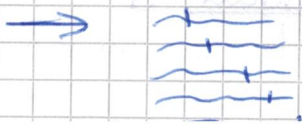
Fingerpint

„Ist so 90er Jahre“

- Man kann keine Reihenfolgen unterscheiden, d.h. 2 Peptide mit exakt denselben AS in unterschiedlicher Reihenfolge haben die gleiche Masse
- die Massen sind sich auch sonst oft sehr ähnlich

Was hat man sich überlegt?

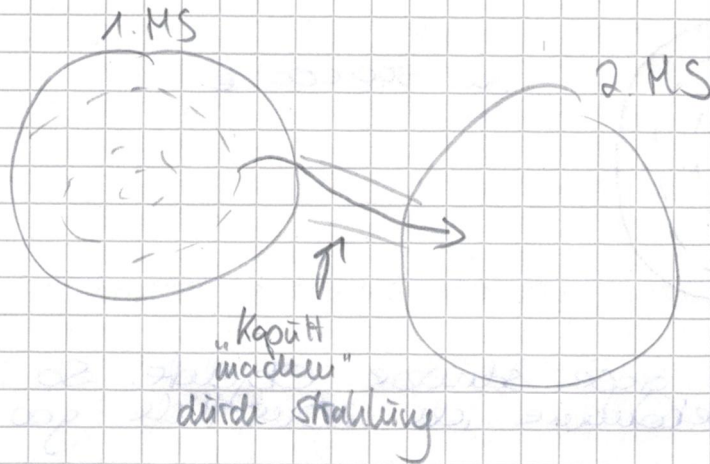
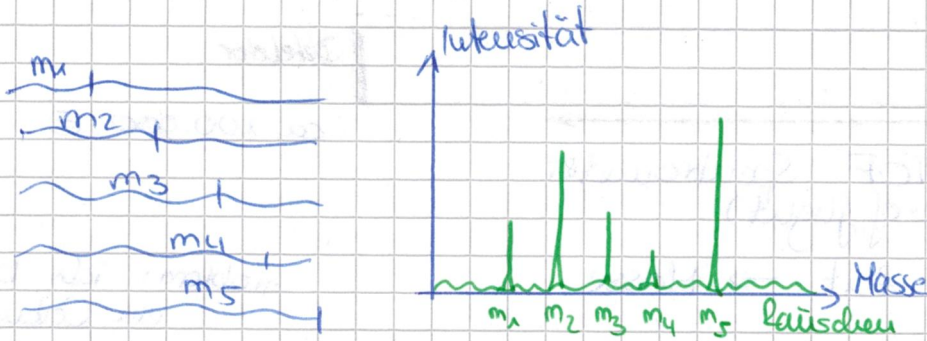
für einzelne Fragmente 2. Spektrometrie



alle Präfixmassen

=> Damit lässt sich dann meist eindeutig das Protein bestimmen.

Nennt sich MS/MS



► Wie kann man das weit abstrahieren, damit man „allgemeine Dinge tun kann“?

Geg: Gewichtetes Alphabet (Σ, μ)

Menge Σ : Zeichen

Funktion $\mu: \Sigma \rightarrow \mathbb{Z}$

String $s \in \Sigma^*$

(das ist keine große Einschränkung)

Def: Speletrum von s : $M = \{ \mu(t) : t \text{ ist echtes Präfix oder Suffix von } s \}$

Mittlermasse (parent mass): $\mu(s)$

Bsp: $\Sigma = \{a, b, c, d\}$ $\mu: \begin{matrix} a: 2 \\ b: 3 \\ c: 7 \\ d: 10 \end{matrix}$

$s = a c a b$ (14)

$M = \{0, 2, 3, 5, 9, 11, 12, 14\}$

Aussage: Alle Präfix- und Suffixmassen sind verschieden. (oft falsch!)

Problem: Geg: $M, \mu(s)$ Gesucht: s

In Praxis:

- Rätschen
- Massen fehlen
- Massen sind falsch
- Massen kommen mehrfach vor
- ...

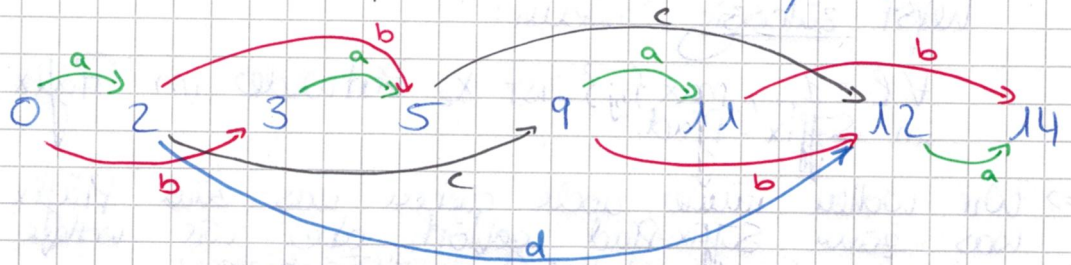
hier: sehr starke Vereinfachung

Def: Spielraum-Graph

Knoten $V := M \cup \{0, M\}$

Kanten: $u \rightarrow v: \exists \sigma \in \Sigma \text{ mit } u + \mu(\sigma) = v$

Bsp:



Idee: Pfad $0 \rightsquigarrow M$ entspricht String

Beob.: $x \in M \Rightarrow M-x \in M$

Benenne Massen in $M \cup \{0, M\}$ als

$$0 = x_0 < x_1 < \dots < x_n < y_n < \dots < y_1 < y_0 = M$$

$\underbrace{\hspace{10em}}_{2n \text{ verschiedene Massen}}$
 $x_i + y_i = M$

Einfacher Algorithmus:

2^n mögl. Zuordnungen x_1, \dots, x_n

→ Präfix oder Suffix

Probiere alle durch, teste ob String zulässig ist

Bsp.: $M = \{0, 2, 3, 5, \dots\}$

$\begin{matrix} p \\ s \\ \vdots \\ s \end{matrix} \Rightarrow \text{Masse} = 1, \text{ geht nicht}$

$\begin{matrix} a \\ a & b \\ a & & ab \\ \underbrace{\hspace{2em}}_{=7} \end{matrix}$

Ende, wir brauchen $u+1$, also d.h. es fehlt $7 = c \Rightarrow a < ab$

Laufzeit: $O(2^n \cdot n)$

→ dann überbrücken!

Idee: Konstruiere 2 Pfade gleichzeitig.

- Präfix-Pfad
- Suffix-Pfad

! disjunkt (wg. Annahme)

! jeder Knoten muss besucht werden

Def: Paar von Pfaden

$x_0 \rightsquigarrow x_i$ Präfix-Pfad
 $x_0 \rightsquigarrow x_j$ Suffix-Pfad

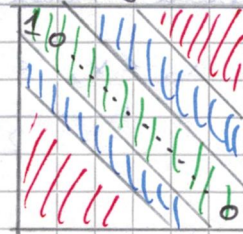
heißt zulässig, wenn:

$\forall l = 1, \dots, \max\{i, j\}$ ist x_l entweder im Präfix-Pfad oder im Suffix-Pfad.

\Rightarrow Wir wollen immer noch eücken was zum Präfix- und was zum Suffix-Pfad gehört, aber wir wollen es besser machen als alles durchzuprobieren!

Definiere **Zulässigkeitsmatrix** $D[0, \dots, n, 0, \dots, n]$

$D[i, j] := \begin{cases} 1 & \text{wenn es ein zulässiges Pfadpaar } x_0 \rightsquigarrow x_i \text{ Präfixpfad} \\ & x_0 \rightsquigarrow x_j \text{ Suffixpfad} \\ 0 & \text{sonst} \end{cases}$



$D[0, 0] = 1$

$D[j, j] = 0 \quad \forall j = 1, \dots, n$

$i \geq j+2: D[i, j] = \begin{cases} D[i-1, j] & \text{falls } x_{i-1} \rightarrow x_i \text{ Kante ist} \\ 0 & \text{sonst} \end{cases}$

$j \geq i+2$: analog

$i = j+1$ (Nebendiagonale)

$D[j+1, j] = \begin{cases} \max_l D[l, j] & \text{falls } (x_l, x_{j+1}) \text{ ist Kante, } l=0, \dots, j \\ 0 & \text{falls Menge leer} \end{cases}$

$j = i+1$ symmetrisch

da die Diagonale nur 0en hat bis auf die allwissend 1

Laufzeit: $O(n^2)$

($O(n)$ -Zellen mit Laufzeit $O(n)$)
 konstant
 konstant

⊙ \Rightarrow kann man verallgemeinern (ohne falsche Annahme)

Def: $D[i,j] = 1 \iff \exists$ zulässiges Pfadpaar $(P-, S-)$, so dass
 Präfixpfad bei Masse x_i endet,
 Suffixpfad bei Masse x_j endet
 alle Massen bis $x_{\max\{i,j\}}$ werden
 genau k benutzt.

$$D[i,j] = \begin{cases} 1 & \text{falls } i=0, j=0 \\ D[i-1,j] & \text{falls } i \geq j+2, x_{i-1} \rightarrow x_i \\ D[i, j-1] & \text{falls } j \geq i+2, x_{j-1} \rightarrow x_j \\ \max_{\ell=0, \dots, i-1} \{ D[\ell, j] \} & \text{falls } i=j+1 \\ \max_{\ell=0, \dots, j-1} \{ D[i, \ell] \} & \text{falls } j=i+1 \\ 0 & \text{sonst} \end{cases}$$

$\max \beta = 0$

Reale Bedingungen:

- ① zusätzliche Peaks?
- ② fehlende Peaks?
- ③ Präfixe / Suffixe mit gleicher Masse?
- ④ Gewichtung der Peaks?



mit meinem Peptid

Zu ①: Ich will möglichst viele übereinstimmende Peaks haben, bzw. möglichst wenige zusätzliche.
 Dabei kann man entweder nur zählen oder mit Gewichtung von Peaks arbeiten.

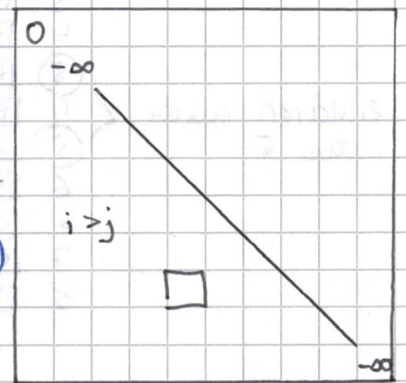
$Q[i,j] :=$ maximale Anzahl (allg. max. Gewicht) von Peaks, die sich durch ~~zusätzl.~~ zulässige Pfadpaare (Präfix $x_0 \rightarrow x_i$, Suffix $x_0 \rightarrow x_j$) erklären lassen.
 $[x_0 = 0] \quad [\max \beta = -\infty]$

$Q[0,0] = 0$

$Q[j,j] = -\infty$

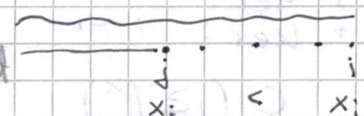
$Q[i,j] = \begin{cases} \max_{\ell=0, \dots, i-1} \{ Q[\ell, j] \} + w(x_\ell, x_i) & \text{Kante von } x_\ell \text{ nach } x_i \text{ hat} \\ \text{od.} & \text{Gewicht} \\ \max_{\ell=0, \dots, j-1} \{ Q[i, \ell] \} + 1 & \text{falls } 0 \leq \ell < i \text{ mit } x_\ell \rightarrow x_i \\ w(x_\ell, x_i) & := \begin{cases} 1 & \text{wenn } x_\ell \rightarrow x_i \\ -\infty & \text{sonst} \end{cases} \end{cases}$

alternativ $:= \begin{cases} f(\text{Gewicht/Intensität der Masse } x_i, x_\ell) & j \\ -\infty & \text{sonst} \end{cases}$



1 zusätzliches Peak erklärt

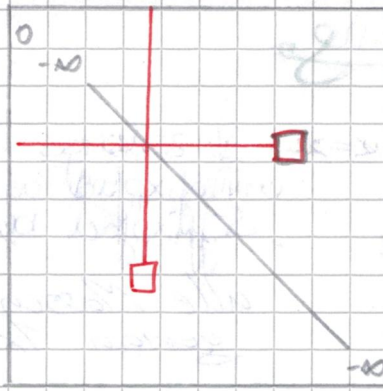
Suffixpfad endet bei i



Zeit: $O(n^3)$

Speicher: $O(n^2)$

Im Allg. kann man hier wenig tun um Speicher oder Laufzeit zu sparen. Aber Peptide sind relativ klein. n ca. 30, deswegen nicht ganz so schlimm. Es ist superschnell!



um das zu berechnen muss ich die ganze Zeile nehmen

um das zu berechnen muss ich die ganze Spalte nehmen

"Eigentlich ist n eine Konstante in dem Algorithmus. Wenn man in der Prüfung sagt $O(1)$. Warum? n ist hier eine Konstante (=30)"

► Zu ②: fehlende Peaks

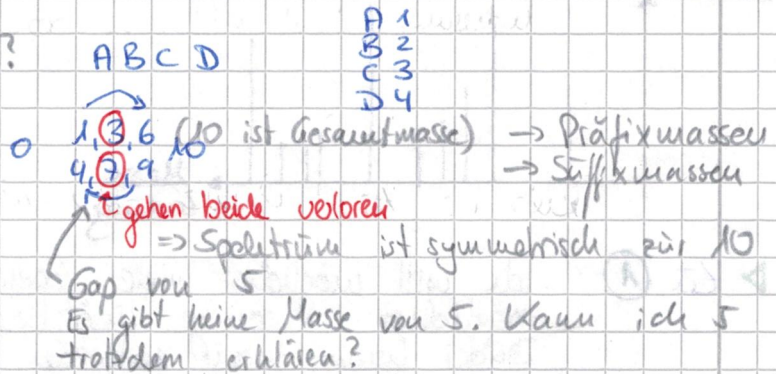
a) von einem Paar fehlt mir P oder S (nicht beide)

$$X \mapsto X_{\text{sym}} = \{X, M-X : X \in X\}$$

symmetrisch

b) beide Massen fehlen

Wie kann das passieren?



⇒ Ich möchte mir nicht einen neuen Algorithmus einführen, deswegen kommt jetzt ein Kunstgriff: Ich vergrößere das Alphabet:

2	AA	
3	AB	= BA
4	AC	= CA
5	AD	= DA
4	BB	
5	BC	= CB
6	BD	= DB
6	CC	
7	CD	= DC
8	DD	

erklären auch die 5

alle Multimengen mit 2 Buchstaben

⇒ möglichst immer die einfachste Lösung nehmen!
Bei Masse 3 → BC

- + 3er Multimengen
- + 4er
- + 5er

erst wenn ich nicht weiterkomme. Explodiert kombinatorisch

⑤ $O(|\Sigma|^k)$

Alphabet auf der-Multimengen erweitern

→ zusätzliche Massenerklärungen.

⇒ Algo aus ① kann verwendet werden. Es muss da nichts mehr zusätzliches gemacht werden.

Es gibt ein kleines Problem mit der Laufzeit ist aber Interpretationssache.

Bisher Alphabet klein → konstante Zeit

Jetzt blähe ich das Alphabet auf. Eigentlich immer noch konstant, dann bleiben die Laufzeiten.

In der Praxis ist es sinnvoll ein Array [oder Hashmap nach Größe sortiert und dann binäre Suche] (oder Hashmap $O(1)$ erwartet)

Zeit: $O(n^3)$ · $\log(|\Sigma|^k)$ (binäre Suche, theoretisch auch perfektes Hashing $O(1)$)
Speicher: $O(n^2)$ + $O(|\Sigma|^k)$

kann irgendwann genauso groß sein wie $O(n^2)$

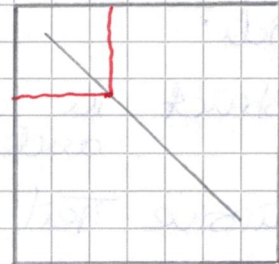
▷ zu ③: Präfixe / Suffixe mit gleicher Masse

Fall $Q[l, j]$ neu überlegen

Peak darf nicht doppelt gezählt werden. Ich will das nicht auch noch bestrafen (ich will es aber auch nicht bestrafen)

Gleiche Formel wie vorher, nur nicht mit +1

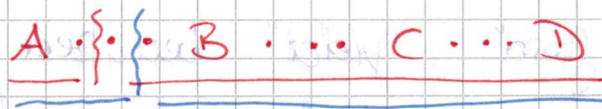
Zulässig ist jetzt: Präfixmasse, Suffixmasse, beides, zusätzliche Masse



$Q[l, j] = \max \left\{ \begin{array}{l} \max \{ Q[l, i] : 0 \leq l < i \text{ mit } x_l \rightarrow x_i \} \\ \max \{ Q[i, j] : 0 \leq i < j \text{ mit } x_i \rightarrow x_j \} \end{array} \right\}$

„Ist elegant und symmetrisch.“

Wir haben jetzt so getan, als könnte man die Strings zwischen den Buchstaben durchschneiden. Treppe sind Mobilfunknetze, die man an verschiedenen Stellen trennen kann.



a-x Paare
b-y -"-
c-z -"-

⇒ Es gibt drei verschiedene Möglichkeiten dazwischen durch zu schneiden.

Einfache Alphabet vergrößern? In Praxis nicht so einfach möglich.

„Kürzungen werden sehr unübersichtlich“. Besser zusätzliche Typen, als Alphabet vergrößern.

Bisher wir kombinatorisch.

(Informations lieber kombinatorisch "zählen Liebe". Statistiker machen das alles mit Wahrscheinlichkeiten.

Statt Maximum summieren über Werten.
=> Uninformulieren in probabilistische Modelle)

-> Max könnte das Ganze mit HMMs (Hidden Markov Modelle) machen.

=> Alle Probleme gelöst!

06.07.17

CO-U

oben steht der "wahre String", dann die "wahren Massen", wie als String. Dann die Klassen ohne ihre Vielfachheiten sind ganz unter der Alphabet.

Trace-back ist implementiert

Das mit dem * ist das, den ich simuliert habe und das sollte wir einmal rauskommen.

! secret.py ist auf dem Laptop, aber nicht im Git. ☺

invalphabet = invertierte Alphabet; Hashtabelle: Masse zu Buchstaben.

Secret in 61 Zeilen möglich (inkl. Kommentare)

Kleines Trick:

- rekonstruiert in innere Flut zerlegen, die rekursiv ist andere Teil der dem Rest macht

Rekursive Teil -> Funktion von Python (überhaupt

mögliche in C++?

falls jemand weiß: 2 longstand-Platz im Atlantico von Sven spendiert. Danach erklären ☺)

-> Hepatoren schreiben ohne Gedanken um Methoden zu machen.

-> Es kann mehrere Vorgänge geben

-> (...)
"yield from", "yield" benutzen

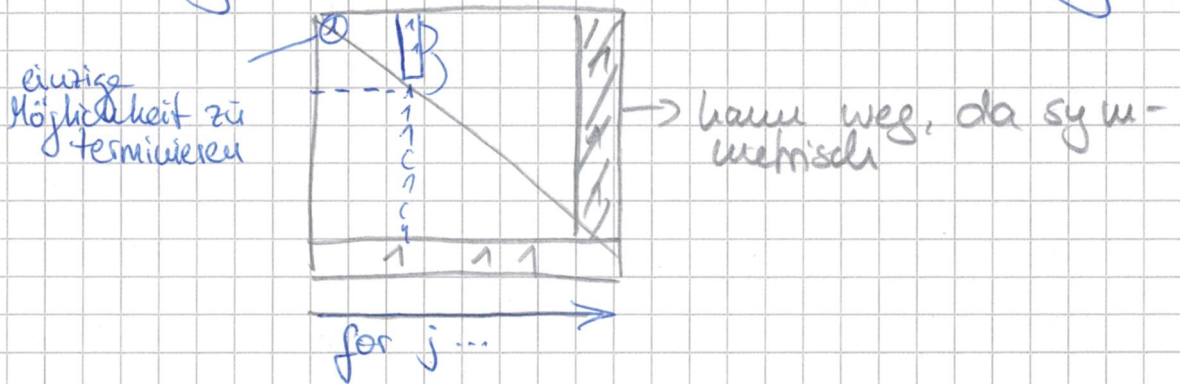
Traceback: wenn oben links, dann ist Prefix und Suffix aufgebaut und wenn nicht dann muss ich gucken, wie ich es aufbaue?

2 Möglichkeiten rekursiv aufzubauen:

- globale Liste, anfangs leer, bei möglicher Lösung erweitern kann groß werden, da es viele Lsg. gibt

ich bin zu langsam und habe nicht alles verstanden

- elegante: Kontrollfluss so zu gestalten, dass
Ausgabe verschärft ist mit Generierung



Die aktuelle Funktion wird mitten im Ablauf unterbrochen, damit die andere Funktion sie direkt weiterverarbeitet.

for s in reconstruct (...)

↓
dann "yield s"

(...)

↓
wieder ein "yield r"

(...)

Jens hat es verstanden!

⇒ ich definiere keine normale Fkt, sondern einen Generator.

Vgl. "[x**2 for x in count()]" vs. "g=(x**2 for ...)"
(bzw. generator-funktion)

➔ in Python Generators-funktion mit yield-Statements

⇒ Secret noch ins Repo

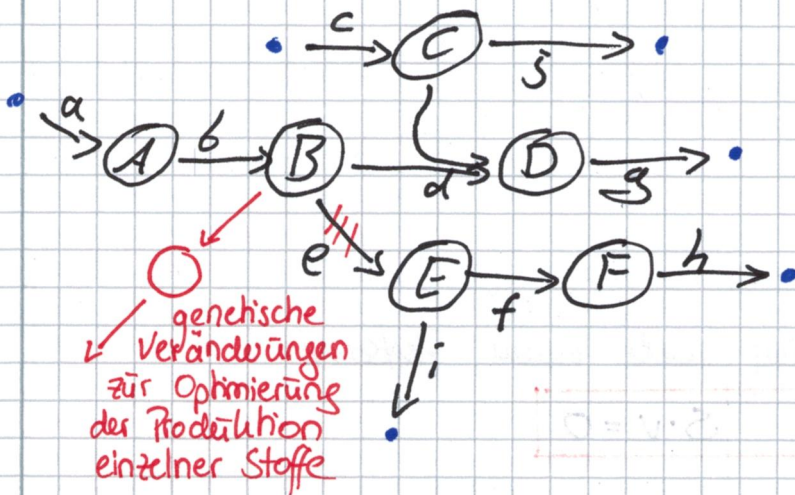
▷ Noch erweitern:

- Argumente für loss, additional
- dann für @-Matrix
- Alphabetvergrößerung

Man berechnet in der Klasse i.d.R. die TOP50 oder 20 Formeln (bei Popeln vllt. Top 10).
Im Prinzip auch suboptimale (sg.) berechnen.

Metabolomics

Analyse metabolischer Netzwerke



- externe Metabolite (nicht betrachtet)
- A-F: Stoffe, compounds, Metabolite
- a-j: Reaktionen

- häufig tatsächlich grafisch dargestellt
-> schlecht.

→ stöchiometrische Matrix

► Allgemein: chemisches Reaktionsnetzwerk "ist" eine stöchiometrische Matrix

$$S \in \mathbb{Z}^{\# \text{Stoffe} \times \# \text{Reaktionen}}$$

S_{ij} : Anzahl Moleküle von Stoff i , die in Reaktion j produziert ($S_{ij} > 0$) bzw. verbraucht ($S_{ij} < 0$) werden.

S =

	a	b	c	d	e	f	g	h	i	j
A	1	-1								
B		1		-1	-1					
C			1	-1						-1
D				1			-1			
E					1	-1			-1	
F						1	-1			

• Nullen sind nicht eingetragen!

• hier der Einheitswert wäre genaugenommen, aber es können auch andere Werte gewählt werden
=> im Gegensatz zum Netzwerkbild oben kann man Verhältnisse darstellen

sehr dünn besetzt

► Flux-Vektor:

$v \in \mathbb{R}^{\# \text{Reaktionen}}$
beschreibt Aktivität jeder Reaktion (Rate, mit der Reaktion Stoffe umgesetzt)

$$E = \begin{pmatrix} a & 1 & 1 & 1 \\ b & 1 & 1 & 1 \\ c & 0 & 0 & 1 \\ d & 1 & 1 & 0 \\ e & 1 & 1 & 1 \\ f & 1 & 1 & 1 \\ g & 1 & 1 & 1 \\ h & 1 & 1 & 1 \\ i & 1 & 1 & 1 \\ j & 1 & 1 & 1 \end{pmatrix}$$

↓ ↓ ↓ ↓
jedes. wervon ist
in Z ~~ist~~
elementare Flux-Moden

Siehe 4 Flüsse der
letzten Seite

Idee: Beschreibe Z als $\{v: v = E \cdot \lambda, \lambda \geq 0\}$

nicht negative Linearkombinationen von Spalten (Flüsse aus Z) von E , E minimal.

Spalten von E : "Elementarmoden"

EFMs

↑ elementare Flux-Moden

(in den Spalten)

Frage: $S \rightarrow E$?

Matrix der Basisvektoren des Nullraums (Kerns) von S
Lösung: Vektorräume

Bekannt: $\{v: Sv = 0\} = \{E \cdot \lambda, \lambda \in \mathbb{R}^d\}$

Ges.: S
Ges.: E

homogenes lineares Gleichungssystem

Hier: $\{v: Sv = 0, v \geq 0\} = \{E \cdot \lambda, \lambda \geq 0\}$

erhöht enorm die Komplexität

Ges.: S
Ges.: E

In der Praxis kann man Probleme ohne E lösen.

Typisches Problem (ohne E):

max. $w^T v$
Gewichtsvektor

s.d. $Sv = 0$
 $v \geq E$
 $1^T v = 1$
 $v \geq E$

w : Zielfunktion
 E : Vektor von Mindestflüssen

$$(1 \dots 1) \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} = 1$$

$$= v_1 + \dots + v_n = 1$$

S =

	a	b	c	d	e	f	g	h	i	j
A	1	-1	0	0	0	0	0	0	0	0
B	0	1	0	-1	-1	0	0	0	0	0
C	0	0	1	-1	0	0	0	0	0	-1
D	0	0	0	1	0	0	-1	0	0	0
E	0	0	0	0	1	-1	0	0	-1	0
F	0	0	0	0	0	1	0	-1	0	0

$Sv = 0$

Nulldraum von S,

Basis?

freie Variablen

$f v: v = 3 \cdot \lambda, \lambda \in \mathbb{R}^d$

Beispiel um Rechnung zu verdeutlichen:

$$\begin{matrix} \text{I} \\ \text{II} \\ \text{III} \end{matrix} \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{matrix} \text{II} - 4 \cdot \text{I} \\ \text{III} - 7 \cdot \text{I} \end{matrix} \begin{pmatrix} 1 & 2 & 3 \\ 0 & -3 & -6 \\ 0 & -6 & -12 \end{pmatrix}$$

$$\text{III} - 2 \cdot \text{II} \begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

x_3 beliebig: λ_1

x_3 : beliebig

$x_2 = -2x_3$

$x_1 = -3x_2 - 2x_3$

$= -3(-2x_3) - 2x_3 = 6x_3 - 2x_3 = 4x_3$

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 4x_3 \\ -2x_3 \\ x_3 \end{pmatrix} = \lambda_1 \cdot \begin{pmatrix} 4 \\ -2 \\ 1 \end{pmatrix}$$

alle Vektoren dieser Form lösen dieses System

$v = \begin{bmatrix} 4 \\ -2 \\ 1 \end{bmatrix} \cdot [\lambda_1]$ Koeffizient

Basisvektor

mit ein Basisvektor da Dimension 1

=> Unsere Matrix S ist schon in Stufenform

v_g, v_h, v_i, v_j beliebig = $\lambda_1, \lambda_2, \lambda_3, \lambda_4$

$$v_f = v_h$$

$$v_e = v_f + v_i \\ = v_h + v_i$$

$$v_d = v_g$$

$$v_c = v_d + v_j \\ = v_g + v_j$$

$$v_b = v_d + v_e \\ = v_g + v_h + v_i$$

$$v_a = v_g + v_h + v_i$$

$$\begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix} = A$$

$$\{0 \leq A \leq 1\}$$

$$\{0 \leq A \leq 1\} \Leftrightarrow$$

Systeme mit den Matrix A

hauptideagonal ist

4 verschiedene beliebige Variablen, deswegen hat meine Basis die "Größe" 4. 4 Basisvektoren!

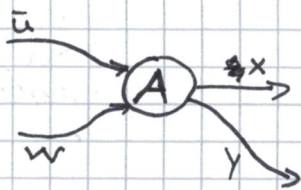
$$v = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \lambda_4 \end{bmatrix} = \begin{bmatrix} v_g \\ v_h \\ v_i \\ v_j \end{bmatrix}$$

ab hier die Identität

zufällig alle ≥ 0

i.d.R. viel mehr Elementarmoden als Basisvektoren!

ANDERES BEISPIEL



4 Elementarmoden: $\begin{matrix} \bar{u} & x \\ \bar{u} & y \\ w & x \\ w & y \end{matrix}$

Weg: $\bar{u} + w + 2y$ ist keine Elementarmode, da es eine Kombination von $\bar{u}y$ und wy ist.
 \Rightarrow Elementarmoden dürfen keine Kombination aus anderen sein

$$S = A \begin{array}{c|ccc} \bar{u} & w & x & y \\ \hline 1 & 1 & -1 & -1 \end{array}$$

3 freie Variablen

$$Sv = 0$$

$$v_{\bar{u}} = (-1 \ 1 \ 1) \begin{pmatrix} v_w \\ v_x \\ v_y \end{pmatrix}$$

$$v = \begin{pmatrix} \bar{u} \\ w \\ x \\ y \end{pmatrix} = \begin{pmatrix} -1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{pmatrix}$$

Allgemeiner:

$$A = \begin{bmatrix} S \\ -S \\ I \end{bmatrix} \rightarrow \text{bleutität}$$

$$\{v: Av \geq 0\}$$

$$\Leftrightarrow \{P\lambda: \lambda \geq 0\}$$

"P steht für polyedrischer Kegel oder so"

$$\begin{matrix} A \rightarrow P \\ P \rightarrow A \end{matrix}$$

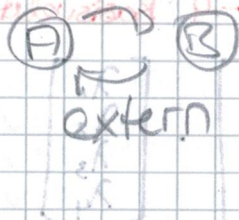
=> Problem aus der Geometrie

Double Description Method

Da man nicht reversible Zyklen zulässt, bekommt

man durch Elementarmethoden die Unabhängigkeit

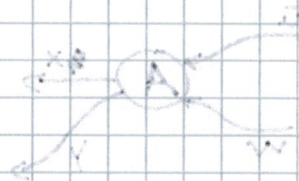
blähen, aber keinen Mehrwert haben.



0	1	1	1
0	1	1	1
1	0	0	1
0	0	0	1
0	1	1	0
0	0	1	0
0	0	0	1
0	1	0	0
1	0	0	0

Andersherum

$$\begin{matrix} x & \bar{x} \\ y & \bar{y} \\ x & \bar{x} \\ y & \bar{y} \end{matrix} = \text{submodularität } \mu$$



submodularität ist für $x+y+\bar{x}+\bar{y}$ (oder) ...

1	1	1	1
0	0	1	1
0	1	0	1
0	1	1	0

$$\begin{pmatrix} x & y \\ x & y \\ x & y \end{pmatrix} (1,1,1) = \mu$$

$$\begin{matrix} x & y & \bar{x} & \bar{y} \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{matrix} \mu = 2$$

Geg.: Matrix S , dann Fluxkegel

$$F := \{v : Sv = 0, v \geq 0\}$$

$$= \{R \cdot \lambda, \lambda \geq 0\}$$

$S \mapsto R$?
 R minimal

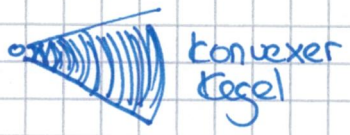
Warum Kegel?

Def.: Punktmenge heißt

• Kegel, wenn $x \in C$
auch $\alpha x \forall \alpha \geq 0$ ent-
halten ist



• konvex, wenn mit $x, y \in K$
auch $\lambda x + (1-\lambda) \cdot y \forall \lambda \in [0, 1]$
enthalten ist (Verbindungsstrecke)



• Konvexer Kegel heißt polyedrisch, wenn
 $P = \{x : Ax \geq 0\}$ für eine Matrix A .

$\Rightarrow F$ ist polyedrischer konvexer Kegel $(A = \begin{bmatrix} S \\ -S \\ I \end{bmatrix}) \quad \left. \begin{matrix} Sx \geq 0 \\ -Sx \geq 0 \\ x \geq 0 \end{matrix} \right\} Sx = 0$

relativ spezielle Form, die
wir auch ausnutzen, aber man
kann auch alles was wir sagen
allgemein fassen

• Für einen konvexen polyedrischen Kegel P heißt (A, R)
Doppelbeschreibungs-paar, wenn

$$P = \{x : Ax \geq 0\} = \{R\lambda : \lambda \geq 0\}$$

Schnitt von Halbräumen
(Zeile $\hat{=}$ Hyperebene,
Normalvektor)

positive Kombination von
("Extrem)strahlen" des Kegels
(Spalte $\hat{=}$ Strahl in P)

Algorithmus: Doppelbeschreibungsmethode (engl. double description method)

$$A \mapsto R \mapsto A$$



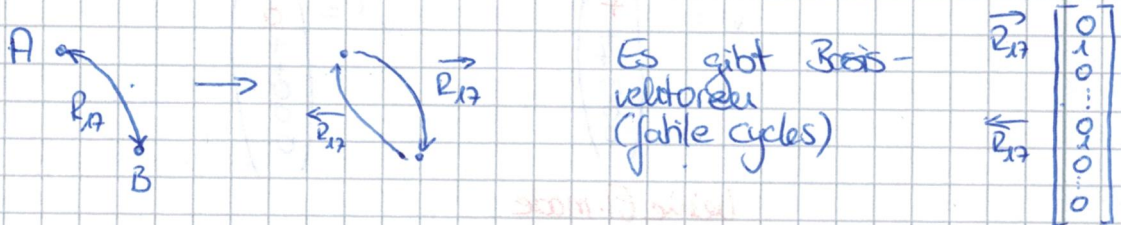
Ist die allgemeinst mögliche Lösung für das Problem.
Sollte man mal gehört haben, aber da unser Problem
spezieller ist, kann man das ausnutzen.

Nullraum-Methode (Wagner, ca. 2005)

1. Bestimme Basis von $\{v: Sv=0\}$ (Nullraum von S)

Basis B hat die Form $B = \left[\begin{array}{c} \pm \\ \mp \\ \pm \\ \dots \\ \mp \end{array} \right] \left\{ \begin{array}{l} \# \text{Reaktionen für ein } K \\ \text{Dim. Nullraum} \end{array} \right.$

(Falls es reversible Reaktionen gab und wir daraus zwei irreversible gemacht haben:
Betrachte diese separat; danach nicht mehr.
 $B = \left[\begin{array}{c} \pm \\ \mp \end{array} \right]$



2. Iteration: Schritt j: Annahme: Zeilen bis j-1 enthalten nur Einträge ≥ 0
Zeilen ab j evtl. auch neg. Einträge

Zu Beginn: Schritt $j = d+1$

Teile Spalten von B am Hand Vorzeichen in Zeile j in 3 Klassen

	≥ 0
j	$+ 0 - 0 - -$
	$? \ ? \ ?$

3 Möglichkeiten
+
-
oder 0

$K^+ := \{k: B_{jk} > 0\}$ übernehmen

$K^0 := \{k: B_{jk} = 0\}$ übernehmen

$K^- := \{k: B_{jk} < 0\}$ weglassen

⊕ neue zulässige Strahlen generieren, die die aus K^- enthalten

Generierung neuer Strahlen: Alle Paare aus $K^- \times K^+$, so dass in Zeile j Nullen erzeugt werden:

$$\Gamma_{k^+,k^+}^* := \Gamma_{j,k^+} \cdot \Gamma_{j,k^-} - \Gamma_{j,k^-} \cdot \Gamma_{j,k^+}$$

Zahl > 0 $\begin{bmatrix} \geq 0 \\ - \\ \end{bmatrix}$ \cdot $\underbrace{- \text{Zahl} < 0}_{\text{pos. Linear-kombination}}$ $\cdot \begin{bmatrix} \geq 0 \\ + \\ \end{bmatrix} = \begin{bmatrix} \geq 0 \\ 0 \\ \end{bmatrix} j$

3. Abschluss: Reduziere das finale B auf Elementarmatrix:

1. $v \in F$

2. Minimalitätseigenschaft:

$Z(v) := \{j : v_j = 0\}$ Nullmenge von v

Es gibt kein v' in $F (v' \neq 0)$ mit $Z(v') \supsetneq Z(v)$

$$v = \begin{pmatrix} 0 \\ + \\ + \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$v' = \begin{pmatrix} 0 \\ 0 \\ + \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

keine El.-matrix

CO-Ü

Metabolomics - Beispieldatei im GIT yeastmodel.txt

1. Datei-format

// Kommentarzeile
 - Steigung / Überschrift
 ENZ $\hat{=}$ Enzym
 REV $\hat{=}$ Reversibel
 MET $\hat{=}$ Metaboliten
 INT $\hat{=}$ intern
 EXT $\hat{=}$ extern
 CAT $\hat{=}$?

als Excel-Datei gespeichert (nur 1 Spalte) \Rightarrow jetzt wieder als

Name: Edukte = Produkte

ecdimodel.txt 2. Dateiformat

\rightarrow nicht so schön

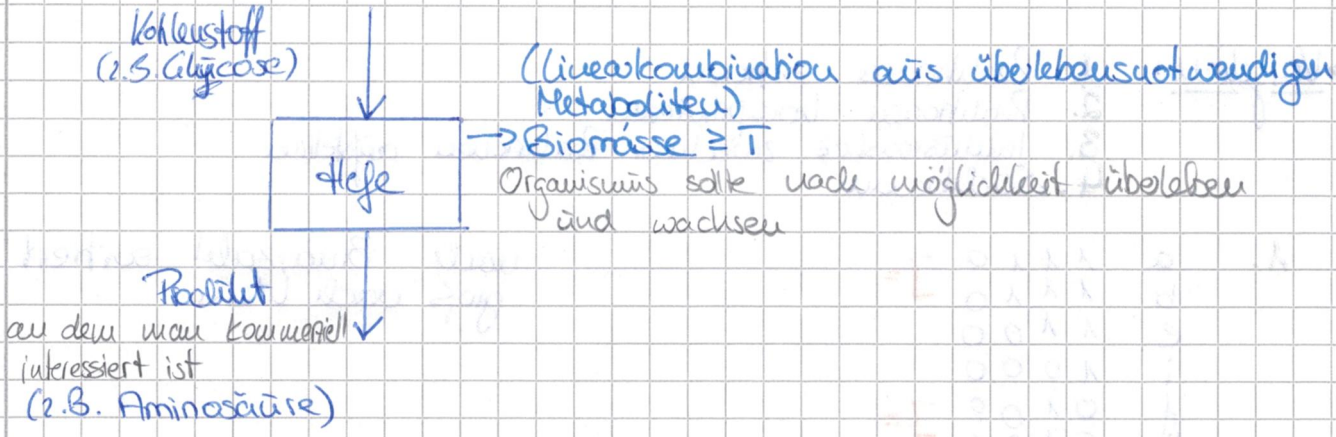
\rightarrow man sieht nicht was reversibel ist

Metabolic Engineering

z.B. Hefe, E.coli

- Löschen von Reaktionen (Heute)
- Hinzufügen von Reaktionen (Reaktion $\hat{=}$ Enzym(e) $\hat{=}$ Gene(e))

► Ausgangssituation: Stöchiometrische Matrix S , Elementarmoden - Matrix E



Wie optimiert man das?
 Nachdem man das einmal gemacht hat, braucht man keine Elementarmoden.

Löschen: Welche Reaktionen muss ich löschen?

► Idee: Lösche ~~wenige~~ ~~alle (?)~~ Reaktionen, so dass Hefe nur Produkte (~~hoch~~ ~~max~~) und Biomasse ($\geq T$) produziert
 $\hat{=} (1-E)max$

"Bleibt" nach kombinatorischem Auswahlproblem.

Elementarmoden E

Reaktionen

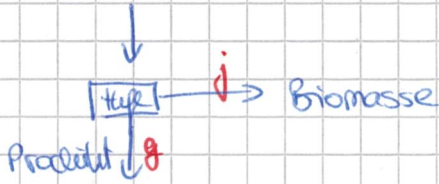
Reaktionen	1	2	3	4
a	1	1	1	0
b	1	1	1	0
c	0	0	1	1
d	0	0	2	0
e	1	1	0	0
f	0	1	0	0
g	0	1	0	0
h	1	0	0	0
i	0	0	0	1
yield y Biomasse?	0	0	1	0
	0	0	0	1

ich brauche alle Reaktionen die da drin sind

e, f, h, i sind zu viel wenn ich e lösche, dann sind dielet El.moden 1 und 2 weg

Elementarmoden 1-4

Wir haben Interesse an Produkt g und Biomasse j



es muss für e rausgeworfen werden (f, h und j unnötig)

=> s. Matrix auf der letzten Seite.

Aufstellen von LP oder selber durchbrauchen? Was geht (schneller)?

- Vorgehen:
1. Reaktionen sortieren
 2. Reaktionen komprimieren
 3. Inklusivitätsbez. zwischen Reaktionen aufstellen
 4. Suchbaum

1.

a	1	1	1	0
b	1	1	1	0
e	1	1	0	0
i	1	0	0	0
f	0	1	0	0
h	0	1	0	0
c	0	0	1	1
d	0	0	1	0
g	0	0	1	0
i	0	0	0	1

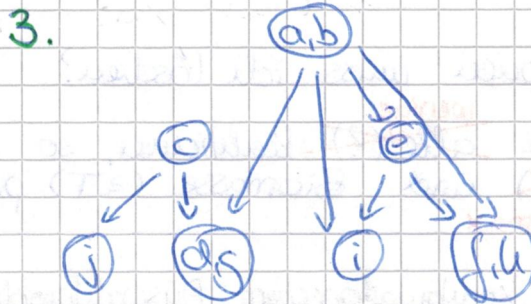
J=

wach Binärzahl sortiert von groß nach klein

Transitive Inklusionen werden mit reingezogenen (deswegen kein Hasse-Diagramm)

2.

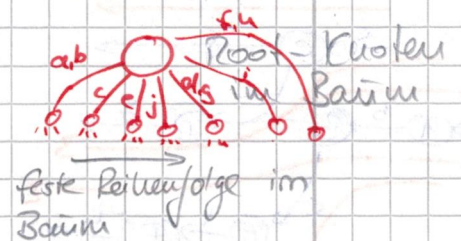
a,b	1	1	1	0
e	1	1	0	0
i	1	0	0	0
f,h	0	1	0	0
c	0	0	1	1
d	0	0	1	0
g	0	0	1	0
i	0	0	0	1



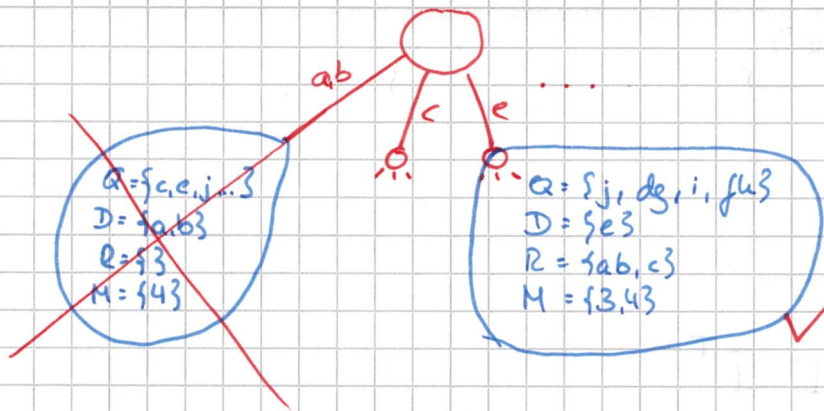
4. Erschöpfende Suche

Q: Reaktionen, über die noch entschieden werden muss
 D: gelöschte Reaktionen → M verbleibende Elementarmoden
 R: nicht gelöschte Reaktionen

Initialisierung Q = alle Reaktionen
 D = {} M = alle Moden
 R = {}



- Entscheidung, welche Reaktion zuerst gelöscht wird (z.B. e: nicht a,b, nicht c löschen)



• Prüfe verbleibende Elementarmoden:

min product yield	0	0
max product yield	1	1
biomass	1	1
reactions deleted	1	1
	ab	e

- Lösung?
- unzulässig?
- weiter verzweigen?

↓
eigentlich will ich es vermeiden, dass es Elementarmoden gibt, die kein Produkt machen, aber das kann man nicht immer vermeiden

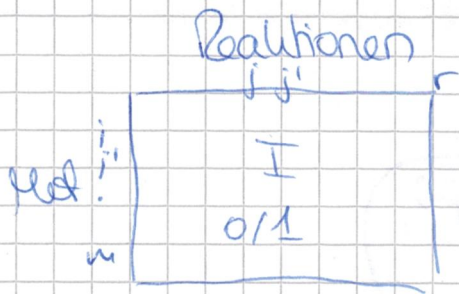
Sehr stumpfsinniger Algorithmus. Alles aufzählen. Früh abbrechen, wenn es keinen Sinn macht.
"Eigentlich macht ein LP-Solver auch nichts anderes"

Es werden alle Lösungen angezeigt. Dann kann Biotechnologie pragmatisch auswählen, z.B. eine 4er Lösung auswählen, weil man eines der Gene nicht austauschen kann und stattdessen lieber die 5er Lösung wählen.

Berechnungen für's Löschen braucht wesentlich weniger lang als Berechnung der Elementarmoden!

Funktioniert wirklich gut!

Sichbaumalgorithmen sind wichtiges Werkzeug (auch wenn es intellektuell nicht so ausspricht ist).



$$I^T \cdot I$$

$(r \times r)$

$$I \cdot I^T$$

$(m \times r)(r \times m) = (m \times m)$

$$(I I^T)_{i, i'} = \sum_j I_{ij} \cdot I_{i'j}$$

$$= \sum 1 [i, i' \text{ in } R_j]$$

= # Reaktionen mit i und i' gemeinsam

Frage aus Fragekatalog (1 Seite 4. Pkt. von unten)

2. Frage: Kreis = Sättiger

15 besteht der De Bruijn Graph voll, also mind.

16 Gerade ist da, also 17. Da Abstand: 25