

V8 Genexpression - Microarrays

- **Idee:** analysiere Ko-Expression von mehreren Genen um auf funktionelle Ähnlichkeiten zu schließen
- **wichtige Fragen:**
 - (1) wie wird Genexpression reguliert?
 - (2) was wird mit MicroArray-Chips gemessen?
 - (3) wie analysiert man Daten aus MicroArray-Experimenten?
 - (4) was bedeutet Ko-Expression funktionell?
- **Inhalt V8:**
 - (1) Hintergrund zu Transkription und Genregulationsnetzwerken
 - (2) Micro-Arrays
 - (3) Übung: analysiere selbst Daten aus einem MicroArray-Experiment

das Transkriptom

Als **Transkriptom** kennzeichnet man den Level an transkribierter messenger RNA (mRNA) für alle Gene des Genoms.

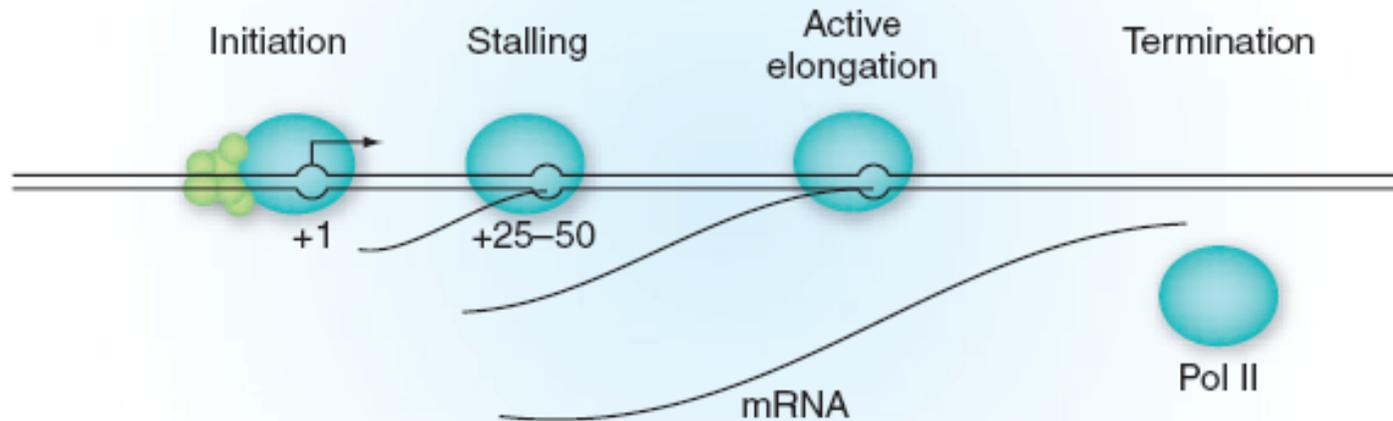
Heutzutage gilt dies sowohl für die Protein-kodierenden Gene als auch für RNA-kodierende Gene, die nicht in Protein translatiert werden.

An die eigentliche Transkription in **pre-mRNA** schließen sich noch viele Prozessierungsschritte zur eigentlichen mRNA an, wie

- die Anheftung eines ca. 250 nt-langen **PolyA-Schwanzes**,
- evtl. Editing (Austausch von Nukleotidbasen), sowie
- Spleißen.

Heute werden wir uns auf den reinen Prozess der DNA-Transkription beschränken.

Transkription durch RNA Polymerase II



Kim Caesar

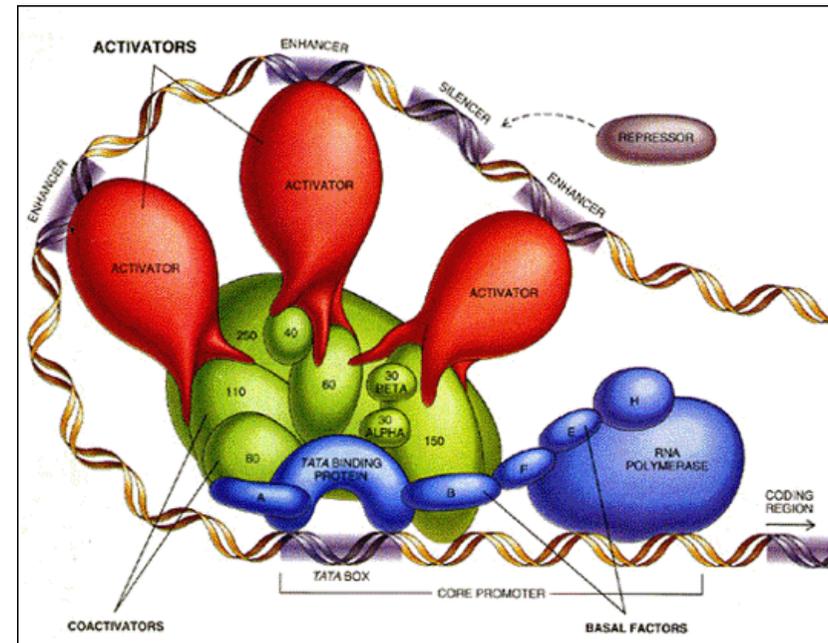
Figure 1 Transcription by RNA polymerase II. Eukaryotic transcription involves a cycle of highly regulated events¹. After clearing the promoter, RNA polymerase II may pause or stall 25–50 base pairs downstream of the transcription start site before transcribing the body of the gene. Pausing is subject to both positive and negative regulation.

Tamkun J. Nat. Gen. 39, 1421 (2007)

Transkriptions – Gen-Regulationsnetzwerke

Die **Maschine**, die ein Gen transkribiert, besteht aus etwa 50 Proteinen, einschließlich der **RNA Polymerase**. Dies ist ein Enzym, das DNA code in RNA code übersetzt.

Eine Gruppe von **Transkriptionsfaktoren** bindet an die DNA gerade oberhalb der Stelle des **Kern-Promoters**, während assoziierte Aktivatoren an Enhancer-Regionen weiter oberhalb der Stelle binden.

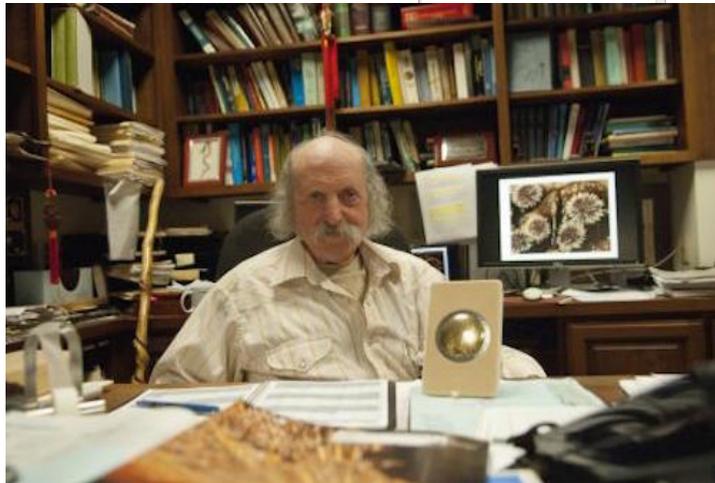
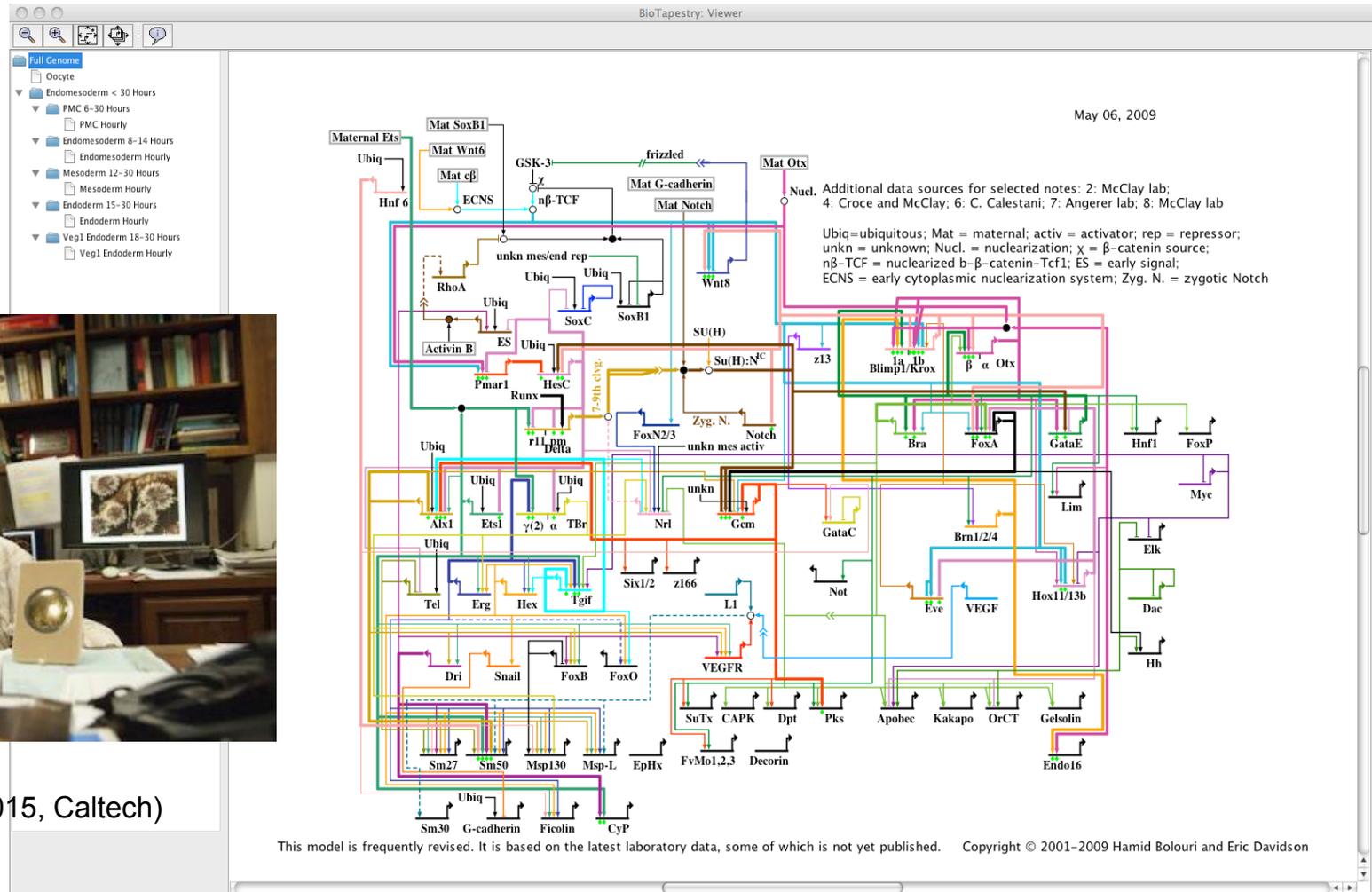


Roger Kornberg
(Stanford Univ)
Noble prize chemistry 2006
„for his studies of the
molecular basis of
eukaryotic transcription“

http://www.berkeley.edu/news/features/1999/12/09_nogales.html

<http://www.osti.gov/>

Gen-Regulationsnetzwerk der Seegurke



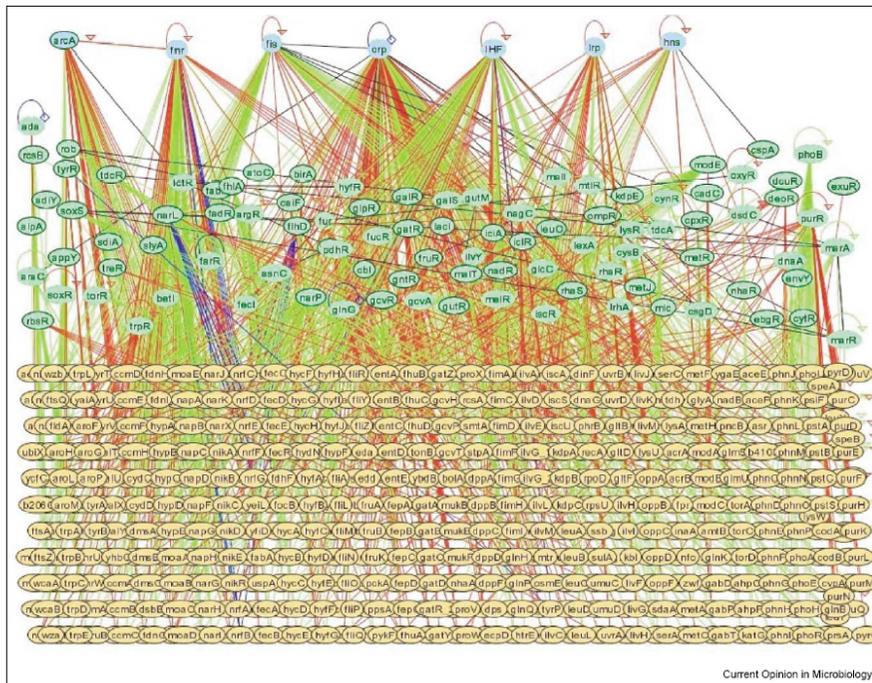
Eric Davidson (1937 – 2015, Caltech)

<http://sugp.caltech.edu/endomes>
<http://www.evolutionnews.org/>

regulatorisches Netzwerk von *E. coli*

RegulonDB: Datenbank mit Information zur transkriptionellen Regulation in *E.coli*; 167 Transkriptionsfaktoren steuern Tausende von Genen.

Durch den hierarchischen Aufbau reichen 7 regulatorische Proteine (CRP, FNR, IHF, FIS, ArcA, NarL and Lrp) aus um die Expression von mehr als der Hälfte aller *E.coli* Gene zu modulieren.



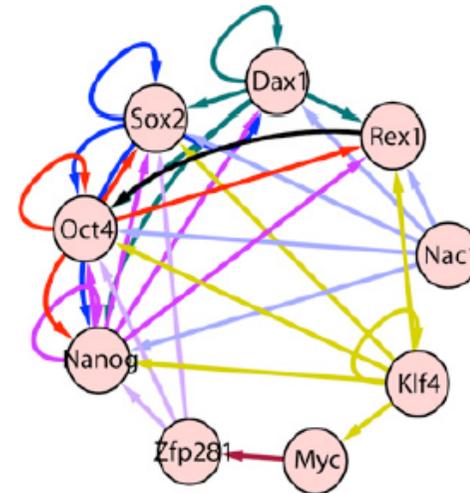
Julio Collado-Vides,
UNAM Mexico-City

Martinez-Antonio, Collado-Vides, Curr Opin Microbiol 6, 482 (2003)

Genregulationsnetzwerk in ESCs um Oct4

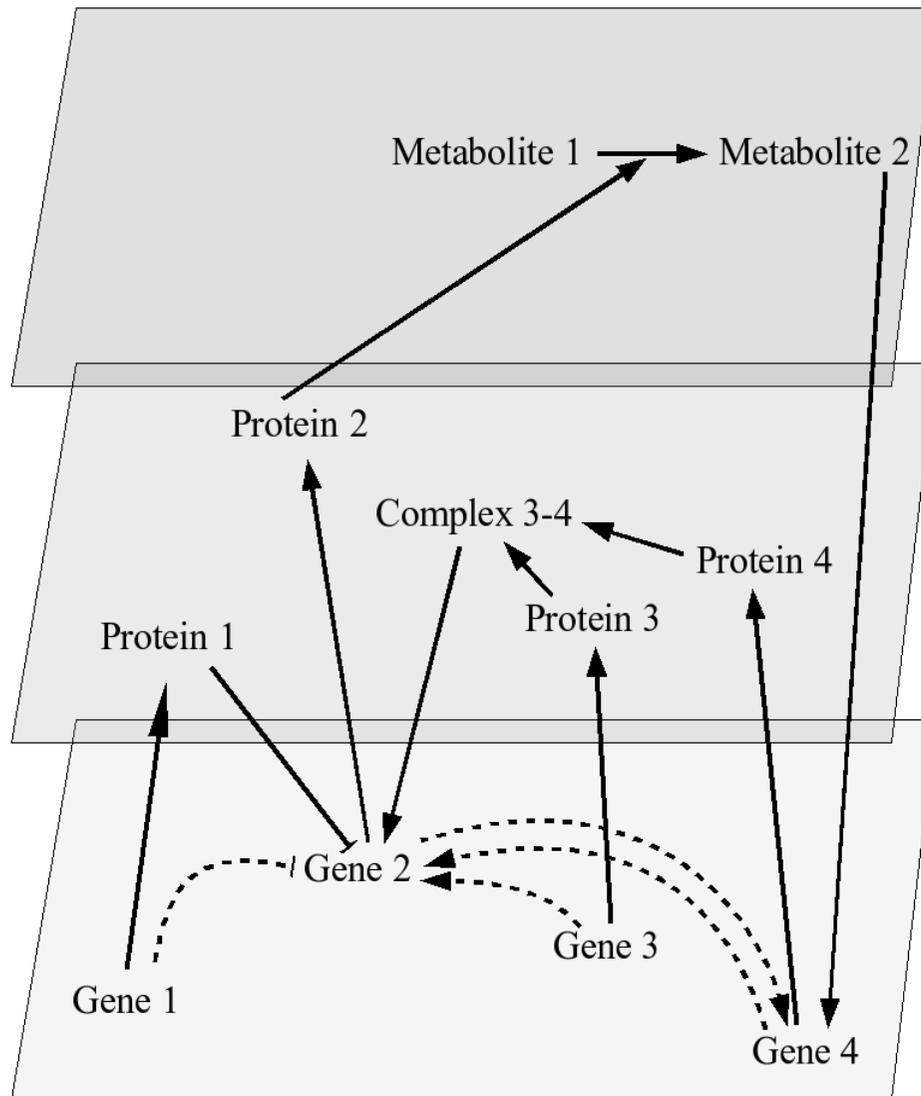
Ein eng verwobenes Netzwerk aus neun Transkriptionsfaktoren hält embryonale Stammzellen (ESC) im pluripotenten Zustand.

Der Masterregulator Oct4 sowie Sox2 und Dax1 haben autoregulatorische Feed-Forward Feedback-Schleifen.

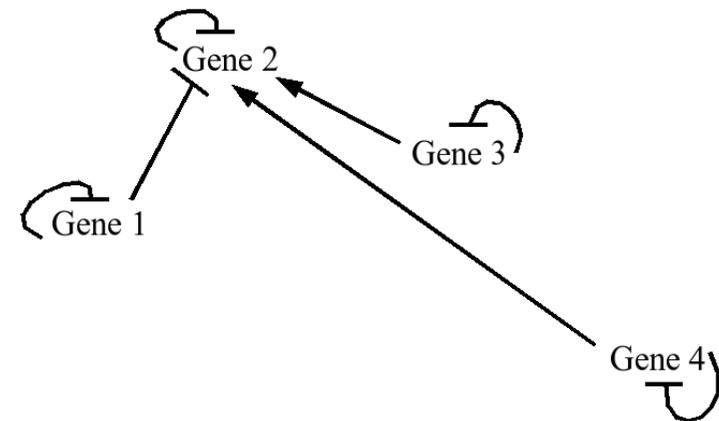


Kim et al. Cell 132, 1049 (2008)

integrierte zelluläre Netzwerke



Statt des komplexen zellulären Netzwerks (links) stellen Genregulationsnetzwerke nur die Projektion auf die Genebene dar (unten).



veränderte Genregulation bei Krankheiten etc.

Ausgangspunkt: bestimmte Krankheiten (Krebs ?) entstehen anscheinend durch die veränderte Expression einer Anzahl von Genen, nicht eines einzelnen Gens.

Wie kann man alle Gene identifizieren, die für diese Veränderung des Phänotyps verantwortlich sind?

Am besten müsste man z.B. die Expression aller Gene in den Zellen von gesunden Menschen und von Krebspatienten bestimmen.

Dann möchte man herausfinden, worin die Unterschiede bestehen.

Genau dies ermöglicht die Methode der **Microarrays**.

Microarrays messen die Expression „aller“ Gene zu einem bestimmten Moment im Zellzyklus unter bestimmten Umgebungsbedingungen.

Was mißt man mit Microarrays?

Häufig verwendet werden **Zweifarbenn-MicroAssays**:

Sample A: rot

Sample B: grün

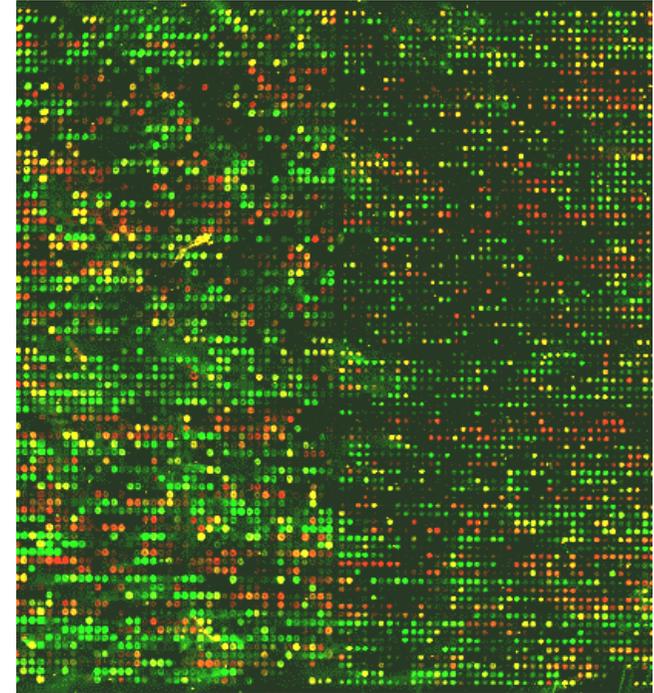
Ziel: bestimme das Verhältnis rot/grün

dunkel: Gen weder in A noch B exprimiert

rot: Gen nur in A exprimiert (bzw. viel stärker)

grün: Gen nur in B exprimiert

gelb: Gen in A und in B exprimiert.



Das Licht wird von zwei Farbstoffen (roter Cy5 und grüner Cy3) erzeugt, die an die cDNA angeheftet wurden (die cDNA wurde „gelabelt“) und die unter Laserlicht fluoreszieren.

Experimentelles Vorgehen

Isolierung einer Zelle im Zustand X

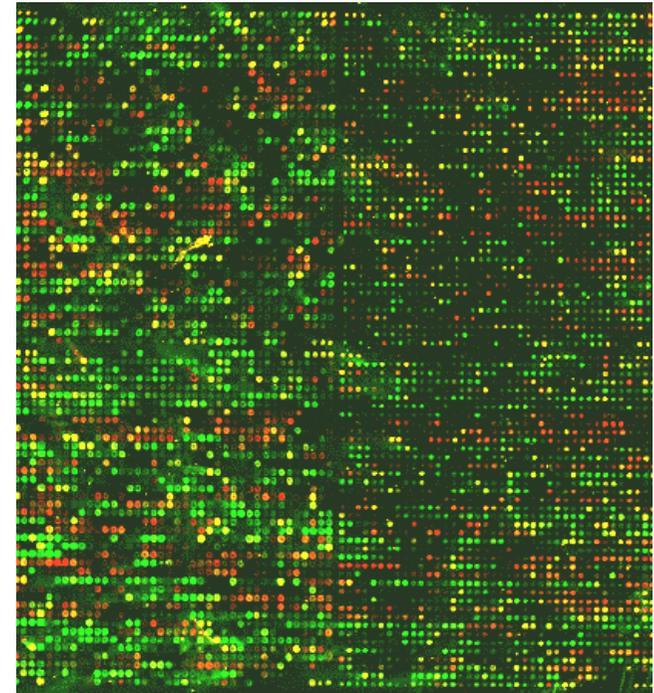
Extraktion aller RNA

Umwandlung in cDNA

Markierung mit Farbstoff (rot oder grün)

Pipette enthält markierte cDNA aller in der Zelle exprimierten Gene.

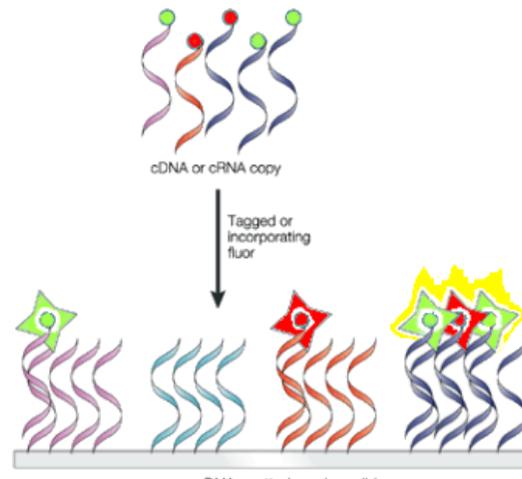
Man bringt nacheinander die cDNA aus zwei verschiedenen Zellpräparationen auf, die unterschiedlich (rot/grün) gelabelt wurden.



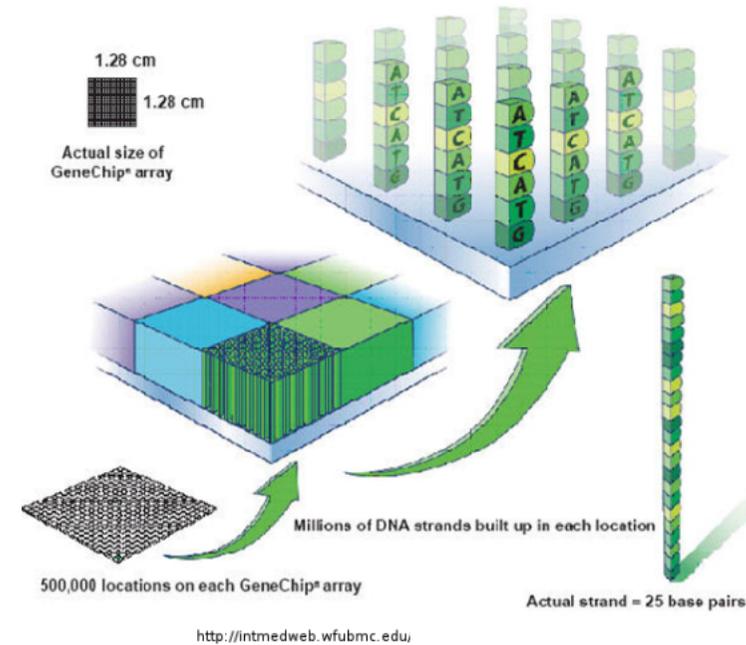
Experimentelles Vorgehen

Aufbringen des zellulären cDNA-Gemischs auf die einzelnen Zellen des Arrays.

Jede Zelle enthält an die Oberfläche funktionalisiert einen cDNA-Klon aus einer cDNA-Bibliothek.



changed from:
A. Butte, Nature Reviews Drug Discovery 1, 951-960, 2002



Jede **Zelle** misst daher die Expression eines **einzelnen Gens**.

Einstellung des Gleichgewichts

Die Gesamtzahl an gebundenen DNA-Strängen zu einer Zeit t sei $n_c(t)$.

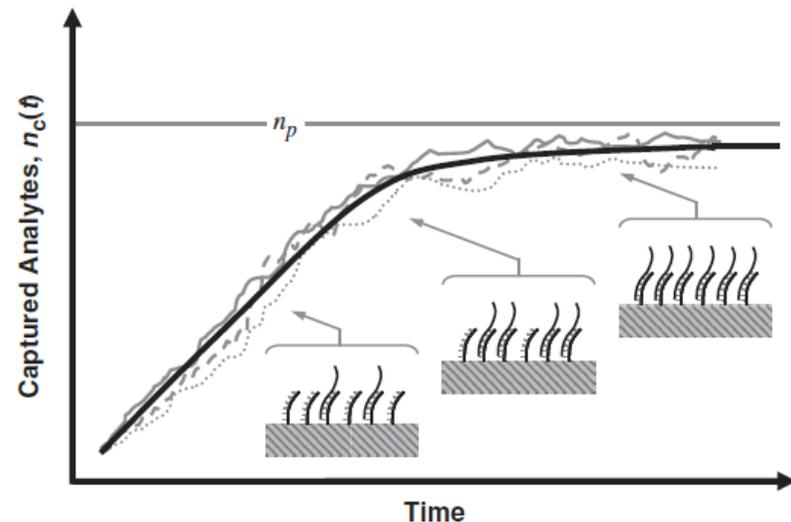
Dann kann man den erwarteten Mittelwert $\langle n_c(t) \rangle$ nach dieser Zeit t durch eine Ratengleichung ausdrücken:

$$\frac{d\langle n_c(t) \rangle}{dt} = k_1^* \left(\frac{n_p - \langle n_c(t) \rangle}{n_p} \right) (n_t - \langle n_c(t) \rangle) - k_{-1} \langle n_c(t) \rangle.$$

k_1^* und k_{-1} : Assoziations- und Dissoziationsraten, mit der die DNA-Stränge der Probe an den Microarray binden,

n_p : Gesamtzahl an freien Bindungsplätzen auf der Microarray-Oberfläche

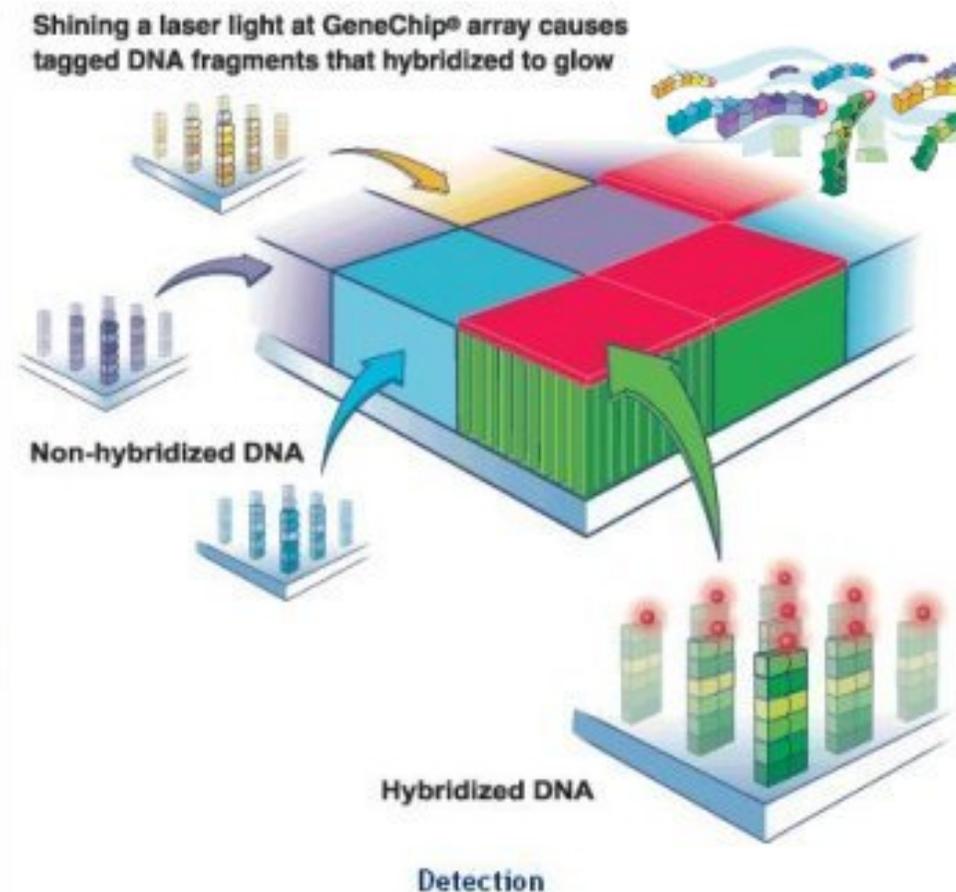
n_t : Gesamtzahl an DNA-Strängen in der Probe



Hassibi et al., Nucl. Ac. Res. 37, e132 (2009)

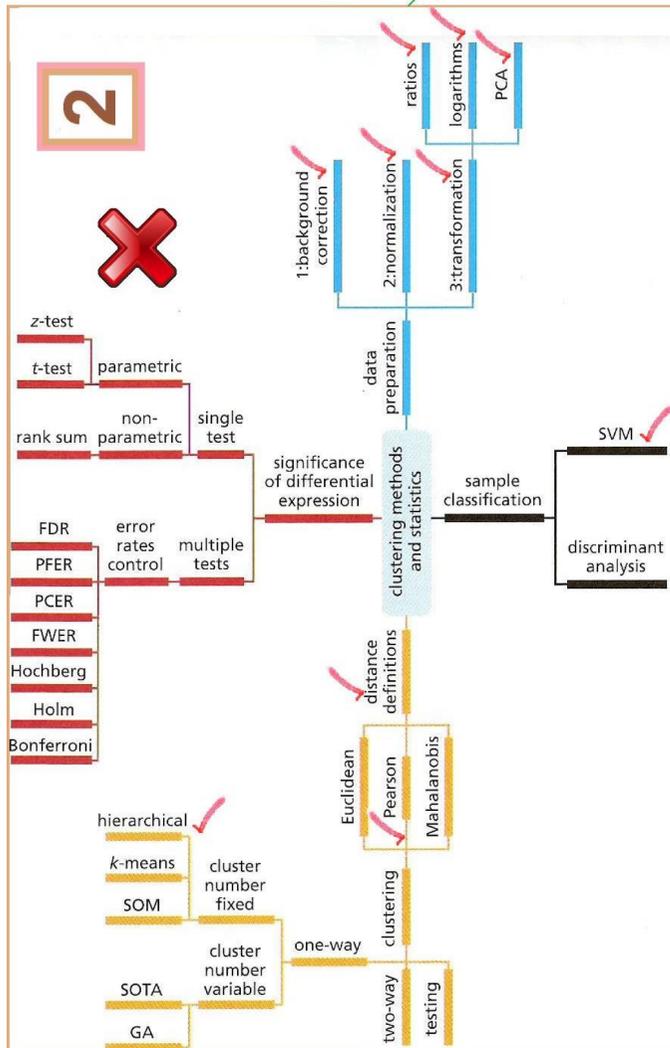
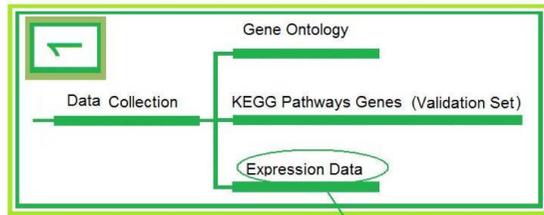
Auslesen der Probe: Laserlicht

Man stimuliert sowohl die Fluoreszenz bei der roten als auch bei der grünen Wellenlänge.



<http://universe-review.ca>

Auswertung von Microarray-Experimenten



Analysis

Korrekte Reihenfolge bei Prozessierung der Daten:

0. Löschen fehlerhafter Daten (Boxplot)
1. Background-correction (Details werden hier nicht behandelt)
2. Normalisierung
3. Transformation (z.B. Log)

Dann

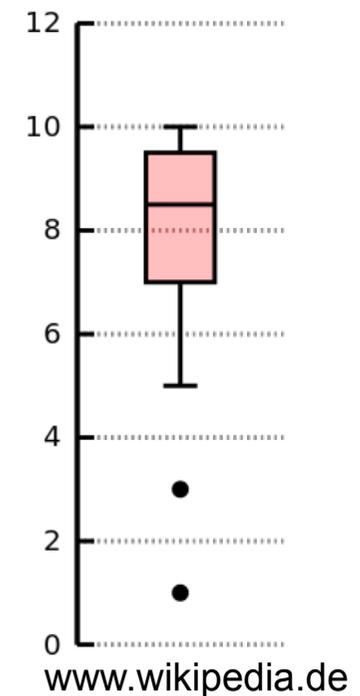
- 4a. Clustering
- bzw.
- 4b. Analyse der Signifikanz

Boxplot

Die Boxplot-Darstellung erlaubt es, schnell einen Überblick über die Werteverteilung in einem Datensatz zu erhalten. Beispiel:

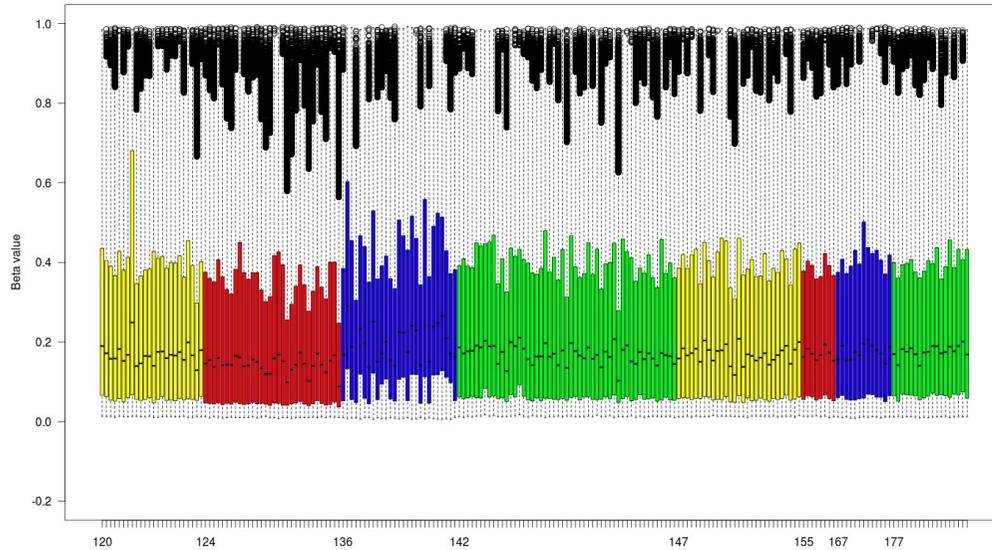
Datenpunkt	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Wert (unsortiert)	9	6	7	7	3	9	10	1	8	7	9	9	8	10	5	10	10	9	10	8
Wert (sortiert)	1	3	5	6	7	7	7	8	8	8	9	9	9	9	9	10	10	10	10	10

Kennwert	Beschreibung	Lage im Boxplot
Minimum	Kleinsten Datenwert des Datensatzes	Ende eines Whiskers oder entferntester Ausreißer
Unteres Quartil	Die kleinsten 25% der Datenwerte sind kleiner oder gleich diesem Wert	Beginn der Box
Median	Die kleinsten 50% der Datenwerte sind kleiner oder gleich diesem Kennwert	Strich innerhalb dieser Box
Oberes Quartil	Die kleinsten 75% der Datenwerte sind kleiner oder gleich diesem Kennwert	Ende der Box
Maximum	Größter Datenwert des Datensatzes	Ende eines Whiskers oder entferntester Ausreißer

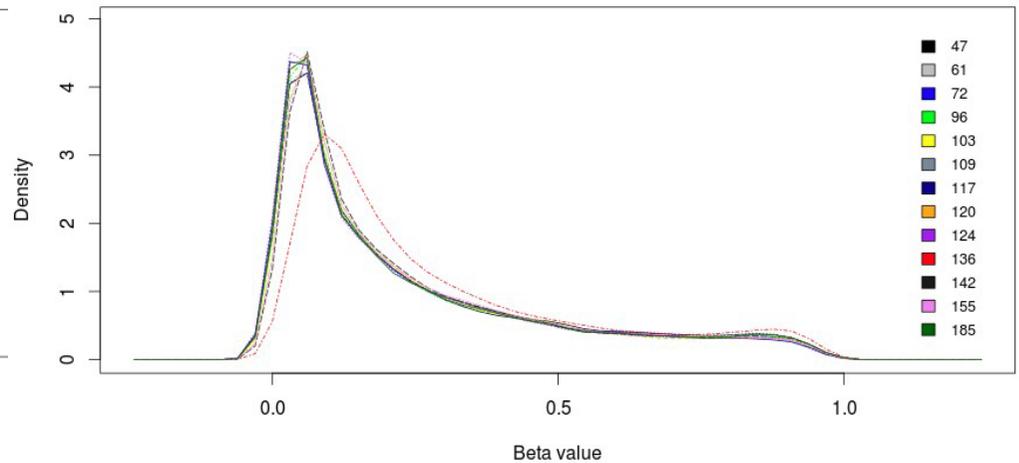
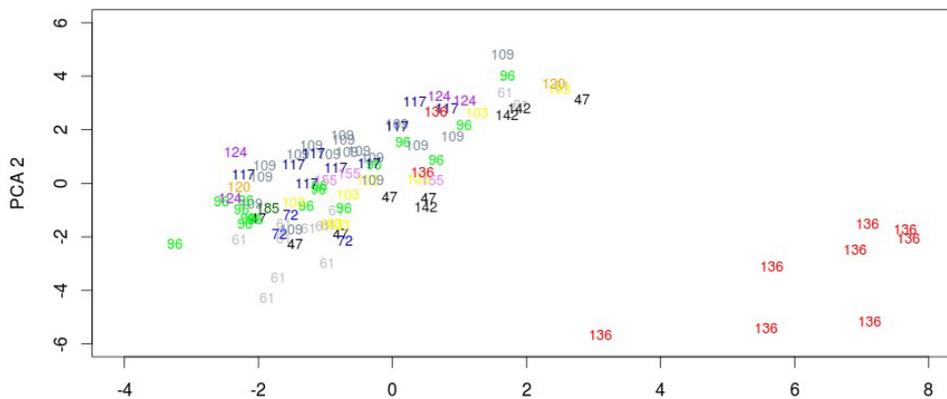


0. Ausreißer-Datenpunkte?

Datensatz 136 in diesen DNA-Methylierungsdaten (Boxplot-Darstellung) verhält sich anders als die anderen Datensätze.

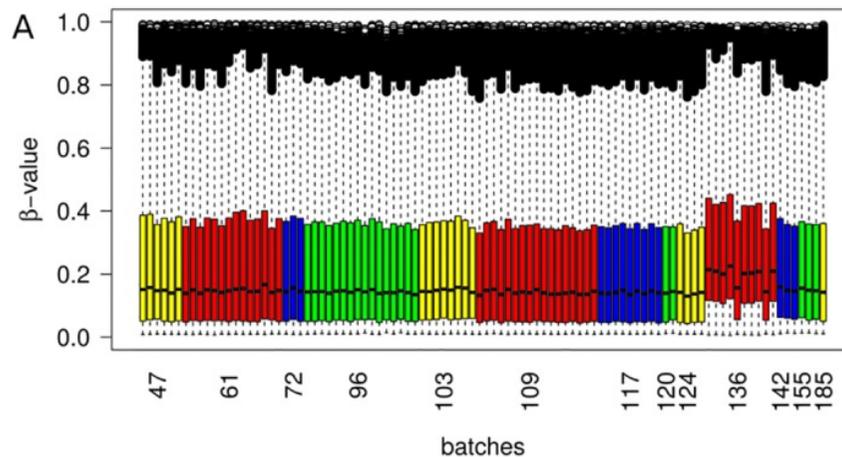


Dies sieht man auch im PCA-Plot (unten links) bzw. im Plot der Werteverteilung (unten rechts).



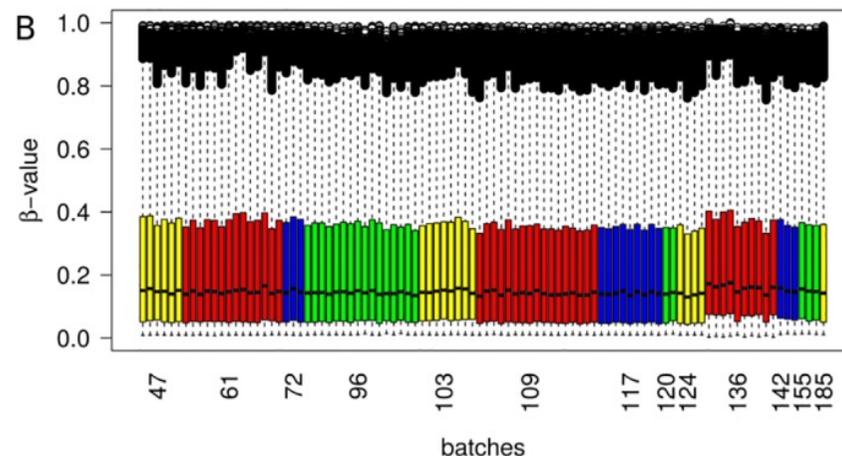
PCA: principle component analysis;
Projektion der Daten auf PC1 und PC2

0. Korrektur von Ausreißer-Datenpunkten



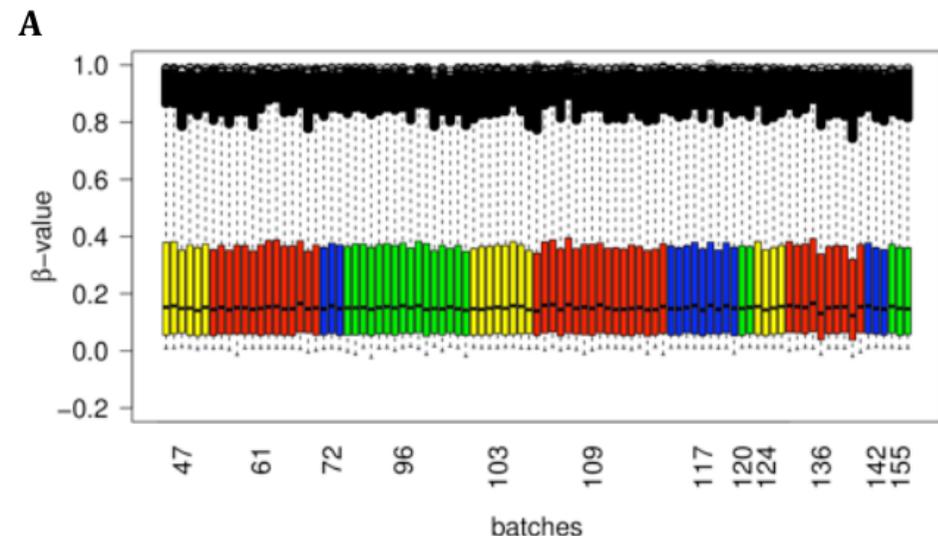
(Bild links oben): Anteil von methylierten CpG-Basen in verschiedenen Samples. Sample 136 ist Ausreißer.

(unten) Korrektur mit unserem Tool BEclear: Nur stark abweichende Werte werden korrigiert: diese Werte werden aus den Werten benachbarter Datenpunkte vorhergesagt. Effekt: natürliche Variation bleibt erhalten.



(Bild rechts) Batch-Effekt-Korrektur desselben Datensatzes mit Tool ComBat: Natürliche Variation der Werte wird stark „geglättet“; alle Werte werden geändert.

Akulenko, Merl, Helms (2016)
PloS ONE 11: e0159921



2. Normalisierung von Arrays

Wie alle anderen biologischen Experimente zeigen auch Microarrays **zufällige** und **systematische Abweichungen**.

Zufällige Schwankungen treten auf

- in der absoluten Menge an mRNA, die eingesetzt wird,
- in der Hybridisierungs-Technik und
- in Waschsritten.

Systematische Unterschiede gibt es z.B. bei den physikalischen Fluoreszenzeigenschaften der beiden Farbstoffmoleküle bzw. durch inhomogen hergestellte Microarray-Chips.

Um diese systematischen Abweichungen der Genexpressionslevel zwischen zwei Proben zu unterdrücken, verwendet man **Normalisierungsmethoden**.

Normalisierung

1. Schritt einer Normalisierung: wähle einen Satz von Genen, für die ein Expressionsverhältnis von 1 erwartet wird (z.B. House keeping-Gene).
2. Berechne eine Normierungsfaktor aus der beobachteten Variabilität der Expression für diese Gene.
3. Wende diesen Normierungsfaktor auf die Expressionswerte der anderen Gene an.

Wichtig: durch diese Normierung werden die Daten verändert.

2.a Normalisierung

Man nimmt grundsätzlich an, dass die absolute Menge an RNA in beiden Messungen (bzw. Proben) dieselbe ist und dass die selbe Anzahl an RNA-Molekülen in beiden Messungen mit dem Microarray hybridisieren.

Dann sollten auch die Hybridisierungsintensitäten der beiden Gen-Mengen gleich sein.

Berechne Normierungsfaktor

$$N_{total} = \frac{\sum_{k=1}^{N_{gene-set}} R_k}{\sum_{k=1}^{N_{gene-set}} G_k}$$

Reskaliere die Intensitäten, so dass
und

$$G'_k = G_k \times N_{total}$$
$$R'_k = R_k$$

D.h. nur die Grünwerte werden skaliert, die Rotwerte bleiben unverändert.

2.b alternativ: Quantile Normalisierung

Gegeben: 3 Messungen von 4 Variablen A – D.

Ziel: alle Messungen sollen eine identische Werte-Verteilung bekommen

A	5	4	3
B	2	1	4
C	3	4	6
D	4	2	8

Originaldaten

A	2	1	3
B	3	2	4
C	4	4	6
D	5	4	8

Ordne jede Spalte nach Größe

A	5.67	4.67	3
B	2	2	3
C	3	4.67	4.67
D	4.67	3	5.67



A	iv	iii	i
B	i	i	ii
C	ii	iii	iii
D	iii	ii	iv

Bestimme in jeder Spalte den Rang jedes Wertes

A	2	Rang i
B	3	Rang ii
C	4.67	Rang iii
D	5.67	Rang iv

Bilde Mittelwert jeder Reihe

Ersetze die Originalwerte durch die Mittelwerte entsprechend dem Rang des Datenfeldes.
Nun enthalten alle Spalte dieselben Werte (bis auf doppelte Datenpunkte) und können leicht miteinander verglichen werden.

Expressionsverhältnis

Der relative Expressions-Wert eines Gens kann als Menge an rotem oder grünen Licht gemessen werden, die nach Anregung ausgestrahlt wird.

Man drückt diese Information meist als **Expressionsverhältnis** T_k aus:

$$T_k = \frac{R_k}{G_k}$$

Für jedes Gen k auf dem Array ist hier R_k der Wert für die Spot-Intensität für die Test-Probe und G_k ist die Spot-Intensität für die Referenz-Probe.

Man kann entweder absolute Intensitätswerte verwenden, oder solche, die um den mittleren Hintergrund (Median) korrigiert wurden.

In letzterem Fall lautet das Expressionsverhältnis für einen Spot:

$$T_{median} = \frac{R_{median}^{spot} - R_{median}^{background}}{G_{median}^{spot} - G_{median}^{background}}$$

Bereich der Expressionsverhältnisse

Das Expressionsverhältnis stellt auf intuitive Art die Änderung von Expressions-Werten dar. Gene, für die sich nichts ändert, erhalten den Wert 1.

Allerdings ist die Darstellung von Hoch- und Runterregulation nicht balanciert.

Wenn ein Gen um den Faktor 4 hochreguliert ist, ergibt sich ein Verhältnis von 4.

$$R/G = 4G/G = 4$$

Wenn ein Gen jedoch um den Faktor 4 runterreguliert ist, ist das Verhältnis 0.25.

$$R/G = R/4R = 1/4.$$

D.h. Hochregulation wird aufgebläht und nimmt Werte zwischen 1 und unendlich an, während die Runterregulation komprimiert wird und lediglich Werte zwischen 0 und 1 annimmt.

3. Logarithmische Transformation

Eine bessere Methode zur Transformation ist, den Logarithmus zur Basis 2 zu verwenden.

d.h. $\log_2(\text{Expressionsverhältnis})$

Dies hat den großen Vorteil, dass Hochregulation und Runterregulation gleich behandelt werden und auf ein kontinuierliches Intervall abgebildet werden.

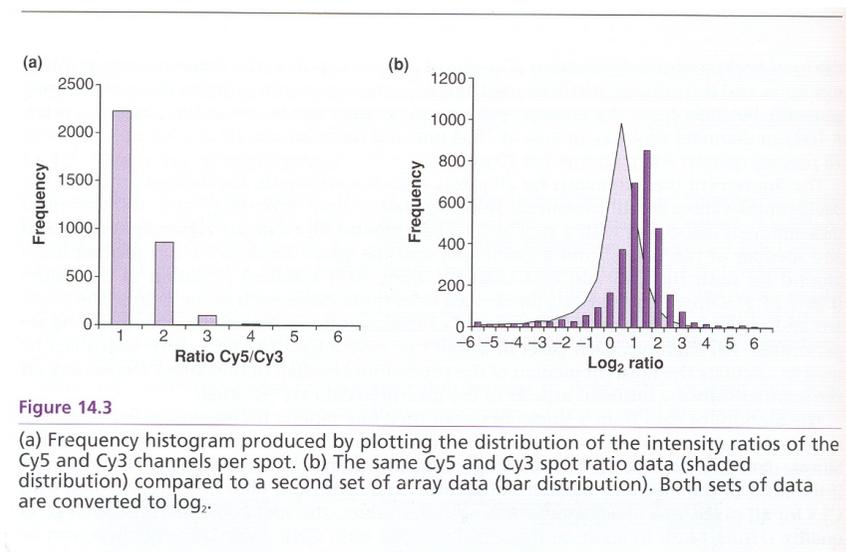
Für ein Expressionsverhältnis von 1 ist $\log_2(1) = 0$, das keine Änderung bedeutet.

Für ein Expressionsverhältnis von 4 ist $\log_2(4) = 2$,

für ein Expressionsverhältnis von $1/4$ ist $\log_2(1/4) = -2$.

Für die **logarithmierten Daten** ähneln die Expressionsraten dann oft einer **Normalverteilung** (Glockenkurve).

M. Madan Babu, An Introduction to Microarray Data Analysis



Orengo-Buch

Daten-Interpretation von Expressionsdaten

Annahme:

Funktionell zusammenhängende Gene sind oft ko-exprimiert.

Z.B. sind in den 3 Situationen

$X \rightarrow Y$	(Transkriptionsfaktor X aktiviert Gen Y)
$Y \rightarrow X$	(Transkriptionsfaktor Y aktiviert Gen X)
$Z \rightarrow X, Y$	(Transkriptionsfaktor Z aktiviert Gene X und Y)

die Gene X und Y ko-exprimiert.

Durch Analyse der Ko-Expression (beide Gene an bzw. beide Gene aus) kann man also funktionelle Zusammenhänge im zellulären Netzwerk entschlüsseln.

Allerdings nicht die kausalen Zusammenhänge, welches Gen das andere reguliert.

4.a Hierarchisches Clustering zur Analyse von Ko-Expression

Man unterscheidet beim Clustering zwischen anhäufenden Verfahren (**agglomerative clustering**) und teilenden Verfahren (**divisive clustering**).

Bei den anhäufenden Verfahren, die in der Praxis häufiger eingesetzt werden, werden schrittweise einzelne Objekte zu Clustern und diese zu größeren Gruppen zusammengefasst, während bei den teilenden Verfahren größere Gruppen schrittweise immer feiner unterteilt werden.

Beim Anhäufen der Cluster wird zunächst jedes Objekt als ein eigener Cluster mit einem Element aufgefasst.

Nun werden in jedem Schritt die jeweils einander nächsten Cluster zu einem Cluster zusammengefasst.

Das Verfahren kann beendet werden, wenn alle Cluster eine bestimmte Distanz zueinander überschreiten oder wenn eine genügend kleine Zahl von Clustern ermittelt worden ist.

Hierarchisches Clustering

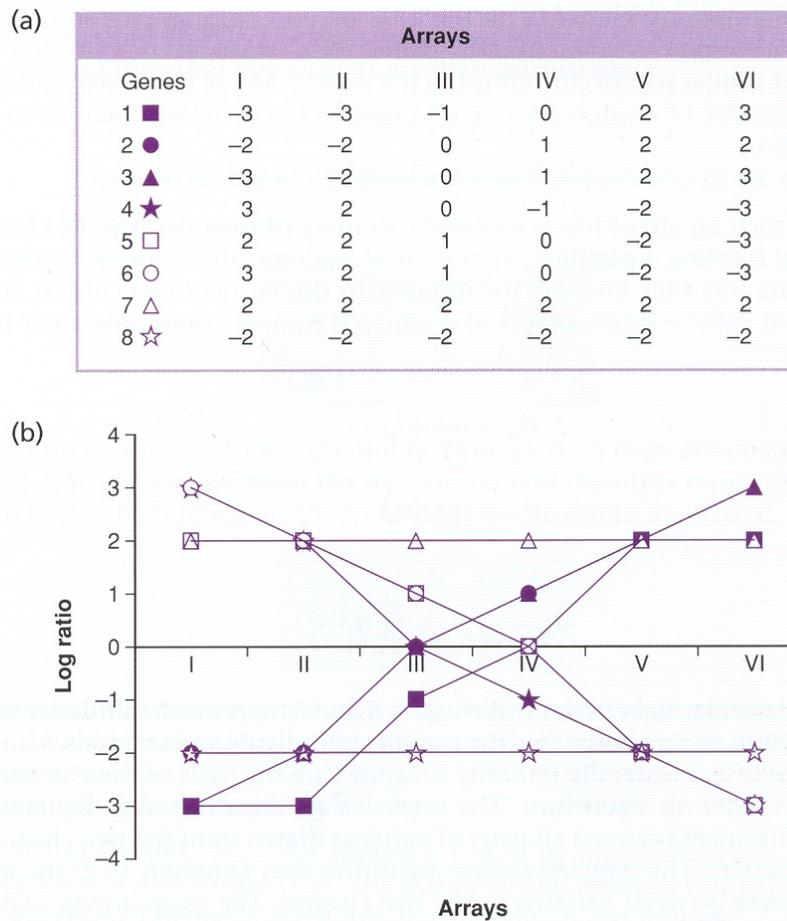


Figure 15.1

(a) Gene names are in column 1 (gene 1 to 8) and each gene is coded per row, with respect to (b). Arrays (I to VI) are in columns (2–7) with each feature being a gene expression \log_2 ratio value (see Chapter 14). (b) A graphical representation of the data for eight genes measured across six arrays corresponding the gene matrix in (a).

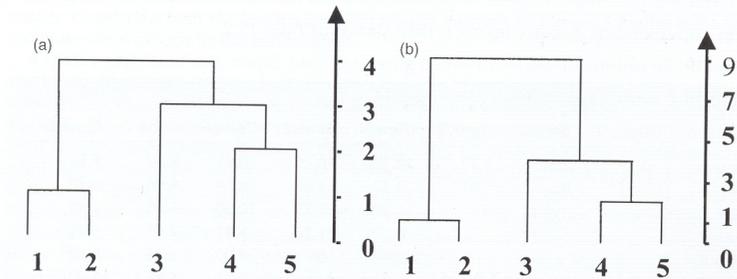


Figure 15.2

Dendrograms from single-linkage (a) and complete-linkage clustering (b)

k-means Clustern

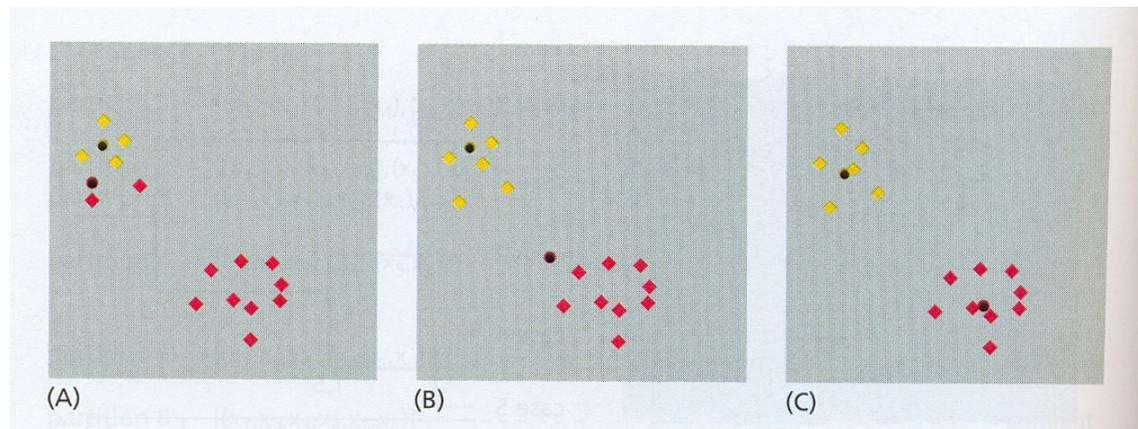
Ein Durchlauf der k -means Clustering Methode erzeugt eine Auftrennung der Datenpunkte in k Cluster. Gewöhnlich wird der Wert von k vorgegeben.

Zu Beginn wählt der Algorithmus k Datenpunkte als Centroide der k Cluster. Anschließend wird jeder weitere Datenpunkt dem nächsten Cluster zugeordnet.

Nachdem alle Datenpunkte eingeteilt wurden, wird für jedes Cluster das Centroid als Schwerpunkt der in ihm enthaltenen Punkte neu berechnet.

Diese Prozedur (Auswahl der Centroide - Datenpunkte zuordnen) wird so lange wiederholt bis die Mitgliedschaft aller Cluster stabil bleibt.

Dann stoppt der Algorithmus



4.b Abschätzung der Signifikanz

Cancer Research



Lipid Metabolism Signatures in NASH-Associated HCC— Letter

Sonja M. Kessler, Stephan Laggai, Ahmad Barghash, et al.

Cancer Res Published OnlineFirst April 28, 2014.

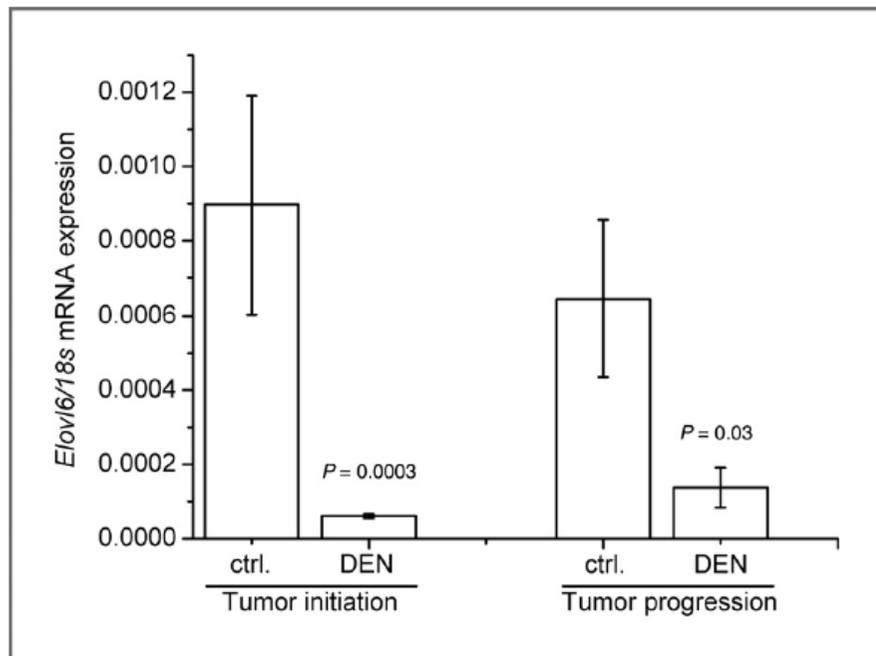


Figure 2. Wild-type mice were treated with the carcinogen DEN at the age of 2 weeks. Livers were analyzed after 24 weeks to assess the tumor initiation state. Analyses in the tumor progression stadium were done after 36 weeks. *Elov16* mRNA expression as determined by real-time reverse transcriptase PCR with $n = 8-18$ per group. Data were normalized to *18S*. Statistical differences compared with untreated animals of the same age (ctrl.) were calculated by Mann-Whitney U test.

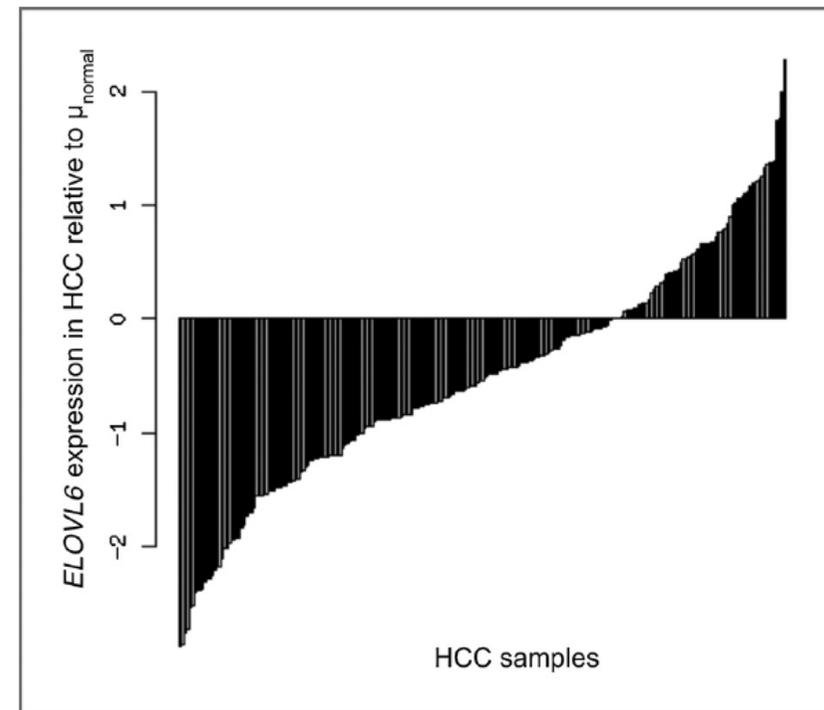
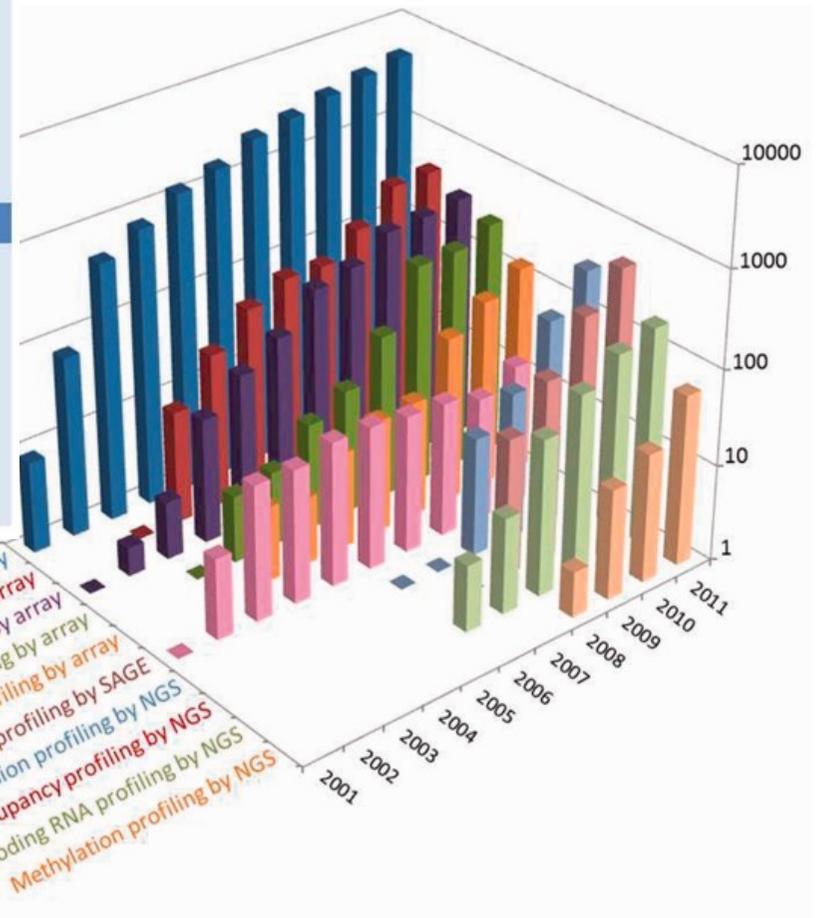
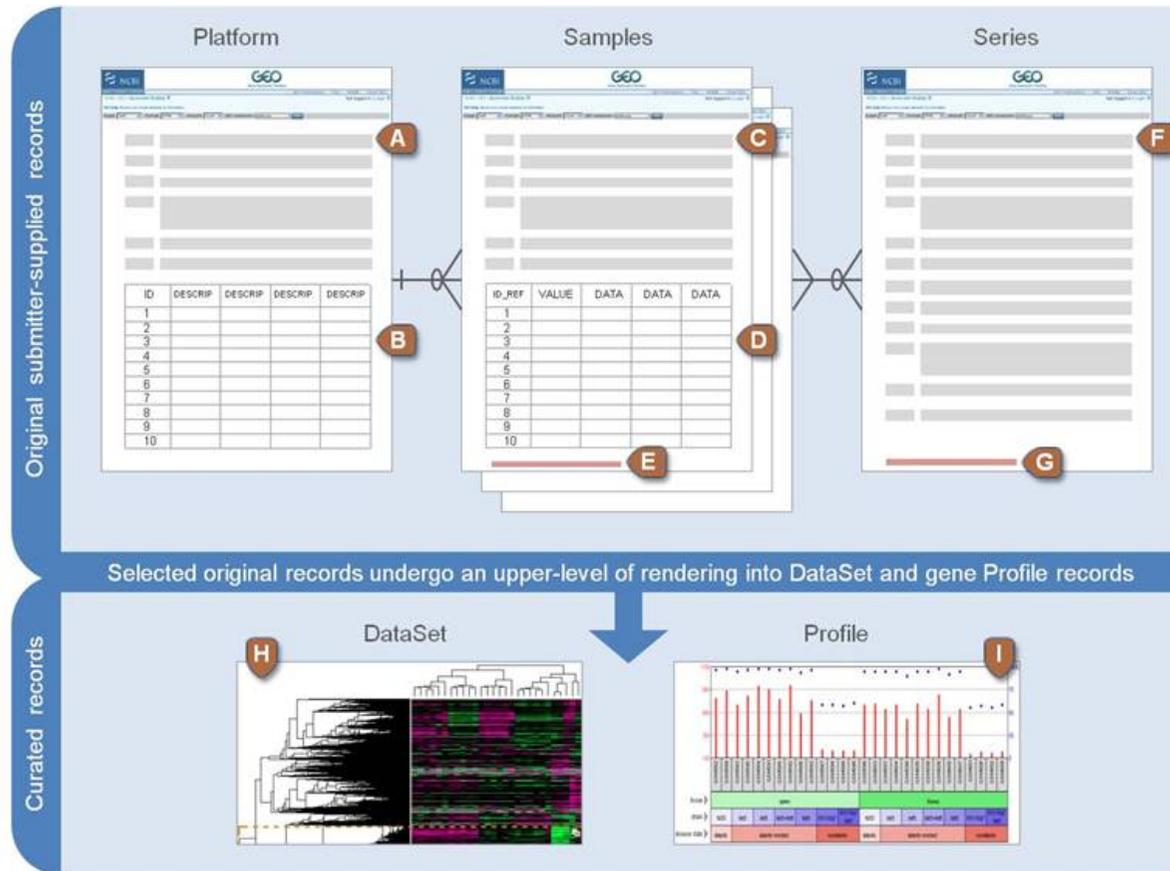


Figure 1. mRNA levels of *ELOVL6* in 247 human HCC samples relative to the mean of 239 nontumor liver tissue (μ_{normal}). Samples of dataset GSE14520 [\log_2 (expression) values from GEO after Robust Multi-array Average normalization] were mapped to hgu133a.db using bioconductor. Significance values: $P = 3.8E-11$, Kolmogorov-Smirnov test; $P = 6.7E-11$, t test; $5.1E-11$, Mann-Whitney U test.

GEO: Gene Expression Omnibus



<http://www.ncbi.nlm.nih.gov/geo/info/overview.html>

Nucleic Acids Res. 41, D991-D995 (2013)

Bewertung von Signifikanz: Mann Whitney Text

Dies ist ein nicht-parametrischer Test. Die abhängige Variable muss nicht normalverteilt sein.

Beispiel: durchschnittliche Noten der Schüler in 2 Schulklassen.

	Schulnoten											Median
Schulklasse A	4.2	6	4.5	4.9	3.9	5	3.6	4.7	5.5	4.3	4.6	4.6
Schulklasse B	4.8	5.8	5.9	4	5.4	3.5	3.8	3.7	5.3	4.4	4.1	4.4

Median : Schüler in Klasse A bessere Noten (Schweiz: 1 bis 6 (am besten)).

Ist der Unterschied statistisch signifikant?

Bilde eine gemeinsame Rangreihe:

Schulnoten	3.5	3.6	3.7	3.8	3.9	4	4.1	4.2	4.3	4.4	4.5	4.6	4.7	4.8	4.9	5	5.3	5.4	5.5	5.8	5.9	6
Schulklasse	B	A	B	B	A	B	B	A	A	B	A	A	A	B	A	A	B	B	A	B	B	A
Gemeinsamer Rangplatz	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22

Bei 2 Stichproben mit identischer zentraler Tendenz würden sich die Rangplätze der beiden Stichproben gleichmässig verteilen und z.B. folgende Muster ergeben:

ABABABABABAB oder AABBBBBAA

www.methodenberatung.uzh.ch/datenanalyse/unterschiede/zentral/mann.html

Bewertung von Signifikanz: Mann Whitney Text

Die Teststatistik U überprüft nun die Gleichmässigkeit der Verteilung der Rangplätze in der gemeinsamen Rangreihe.

Für die erste Stichprobe (Schulklasse A)

lautet die Teststatistik

$$U_1 = n_1 \cdot n_2 + \frac{n_1 \cdot (n_1 + 1)}{2} - T_1$$

mit n_k = Stichprobengrösse der Stichprobe k

T_1 = Rangsumme der Stichprobe 1

Entsprechend gilt für die zweite Stichprobe

$$U_2 = n_1 \cdot n_2 + \frac{n_2 \cdot (n_2 + 1)}{2} - T_2$$

Zwischen beiden Werten besteht folgender Zusammenhang $U_1 + U_2 = n_1 n_2$

Die Rangsumme T_1 für Schulklasse A ist die Summe aller Rangplätze von Werten für Schulklasse A: $2+5+8+9+11+12+13+15+16+19+22 = 132$

Dies ergibt $U_1 = 55$

Für Schulklasse B gilt $T_2 = 121$, $U_2 = 66$

Schulnoten	3.5	3.6	3.7	3.8	3.9	4	4.1	4.2
Schulklasse	B	A	B	B	A	B	B	A
Gemeinsamer Rangplatz	1	2	3	4	5	6	7	8

Bewertung von Signifikanz: Mann Whitney Text

Als Prüfgrösse wird immer der kleinere der beiden Werte verwendet, hier also 55. U gibt die Summe der Rangplatzüberschreitungen an.

Die Frage ist daher, wie oft ein solches Ungleichgewicht der Rangplätze zufällig auftreten kann.

Dazu vergleicht man den kleineren U-Wert mit dem kritischen Wert auf der theoretischen U-Verteilung.

Im konkreten Beispiel ergibt dies eine Signifikanz (p-Wert) von 0.718.

Daher liegt kein statistisch signifikanter Unterschied der zentralen Tendenz zwischen den Klassen vor.

Genauso geht man vor, wenn man den Unterschied der Expression eines bestimmten Gens zwischen zwei Mengen von Proben bewerten möchte.

Zusammenfassung

Die Methode der Microarrays erlaubt es, die Expression aller möglichen kodierenden DNA-Abschnitte eines Genoms experimentell zu testen.

Die **Zwei-Farben-Methode** ist weit verbreitet um differentielle Expression zu untersuchen.

Aufgrund der natürlichen biologischen Schwankungen müssen die Rohdaten **prozessiert** und *normalisiert* werden.

Durch **Clustering** von Experimenten unter verschiedenen Bedingungen erhält man Gruppen von **ko-exprimierten Genen**.

Diese haben vermutlich **funktionell** miteinander zu tun.

Die Signifikanz der unterschiedliche Expression in zwei Gruppen von Proben bewertet man mit statistischen Testverfahren.