

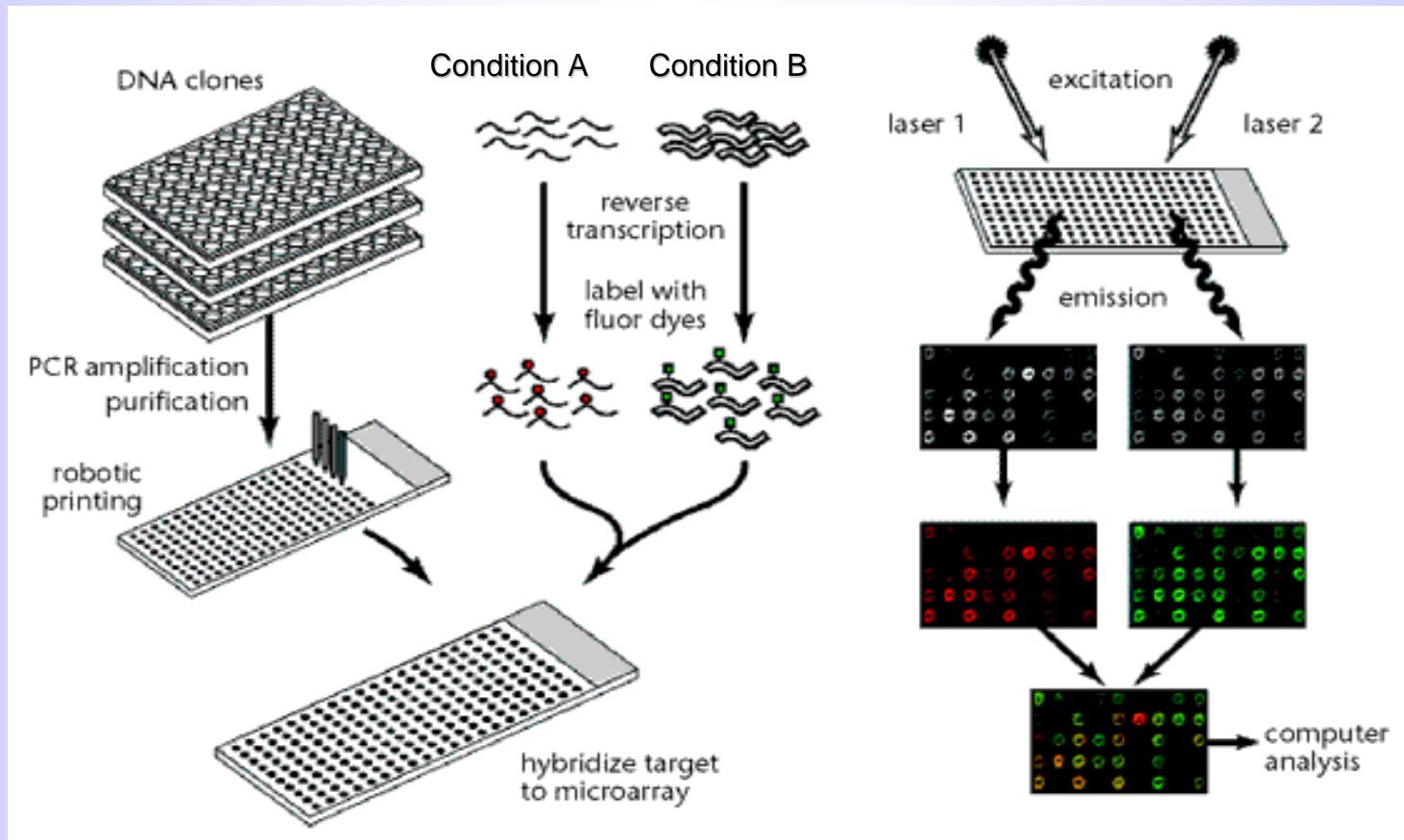
Optimale Versuchsplanung dreifaktorieller cDNA-Microarray-Experimente

1. Einleitung, Modell, Beispiele
2. Optimalität diskreter Designs
3. Ergebnisse
4. Diskussion

Dipl.-Stat. Sven Stanzel
Institut für Medizinische Statistik
RWTH Aachen

GMDS 2006, Leipzig, 11. September 2006

cDNA-Experiment



aus: Duggan et al., Nature Genetics (1999)

Statistisches Modell (1)

2×l×k Design (2 Farben, l Treatments, k Zelllinien)

- **Fixed** effects **gene-specific** linear model for **log-ratios** (Landgrebe et al., 2006):

$$Z_t = \delta_g - \delta_r + \tau_{ij} - \tau_{i'j'} + \varepsilon_t; \quad \varepsilon_t \sim (0; \sigma^2) \quad (1)$$

$t = t(g, r, i, i', j, j'); t = 1, \dots, N$

Z_t : Log Ratio der Fluoreszenzintensitäten auf Array t (für ein spezielles Gen)

δ_g : Farbeffekt des **grünen** Farbstoffs

δ_r : Farbeffekt des **roten** Farbstoffs

τ_{ij} : Kombinationseffekt von Treatment **i** ($i = 1, \dots, l$) und Zelllinie **j** ($j = 1, \dots, k$)

$\tau_{i'j'}$: Kombinationseffekt von Treatment **i'** ($i' = 1, \dots, l$) und Zelllinie **j'** ($j' = 1, \dots, k$)

Statistisches Modell (2)

- Matrixnotation:

$$\mathbf{Z} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\varepsilon} \quad (2)$$

$$\mathbf{Z} = (Z_1, \dots, Z_N)^T \sim N \times 1$$

Vektor der beobachteten Log Ratios

$$\boldsymbol{\theta} = (\delta_g, \delta_r, \tau_{11}, \dots, \tau_{lk})^T \sim s \times 1$$

Parametervektor ($s := lk + 2$)

$$\mathbf{X} \sim N \times s$$

Designmatrix

$$\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)^T \sim (0_N, \sigma^2 I_N)$$

Residualvektor

- Kontraste:

- Kontrastmatrix: $\mathbf{K} \sim s \times q$; $Rg(\mathbf{K}) := r$, $r < s$, $r \leq q$ (abh. von K !!)

- lineare Kontraste: $\mathbf{K}^T \boldsymbol{\theta}$ mit $\mathbf{K}^T \mathbf{1}_s = \mathbf{0}_q$

Beispiel ($l=2$ Treatments: 1, 2 ; $k=3$ Zelllinien: a, b, c):

$$\mathbf{K}^T: \begin{array}{cccccc} g & r & a1 & a2 & b1 & b2 & c1 & c2 \\ 0 & 0 & 1 & -1 & 1 & -1 & 1 & -1 \end{array}$$

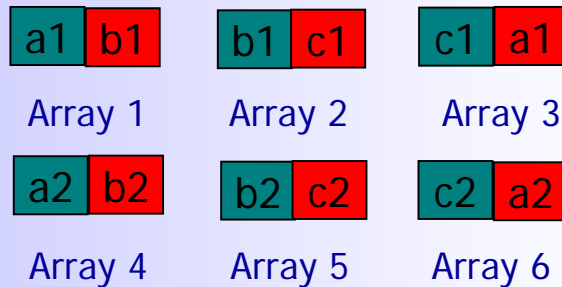
(Treatmenteffekt)

2×2×3 Designs - Beispiele

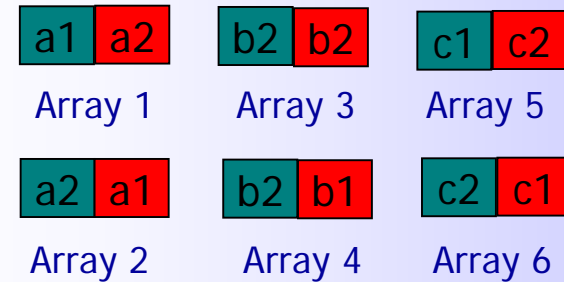
(Landgrebe et al., 2006)

n = 6 Arrays ; **l = 2** Treatments (1, 2); **k = 3** Zelllinien (a, b, c)

Within Treatment Swap (TS)



Within Cell Line Swap (CLS)



Designmatrix (konkretes Design):

Array	g	r	a1	a2	b1	b2	c1	c2
1	1	-1	1	0	-1	0	0	0
2	1	-1	0	0	1	0	-1	0
3	1	-1	-1	0	0	0	1	0
4	1	-1	0	1	0	-1	0	0
5	1	-1	0	0	0	1	0	-1
6	1	-1	0	-1	0	0	0	1

Designmatrix (konkretes Design):

Array	g	r	a1	a2	b1	b2	c1	c2
1	1	-1	1	-1	0	0	0	0
2	1	-1	-1	1	0	0	0	0
3	1	-1	0	0	1	-1	0	0
4	1	-1	0	0	-1	1	0	0
5	1	-1	0	0	0	0	1	-1
6	1	-1	0	0	0	0	-1	1

„E-Effizienz“

(Landgrebe et al., 2006)

- **CLS (2×)** hat **maximale „E-Effizienz“** (unter 9 betrachteten 12-Array-Designs) für das Schätzen von **Treatmenteffekt** sowie **Interaktionseffekt** von Treatment und Zelllinie.
- **TS (2×)** hat **maximale „E-Effizienz“** (unter 9 betrachteten 12-Array-Designs) für das Schätzen des **Zelllinieneffektes**.
- **Schwachstellen:**
 - (globale) Optimalität nicht gesichert, da keine wirkliche E-Effizienz
 - stattdessen: Vergleich ausgewählter Designs anhand des E-Kriteriums
 - Effizienzberechnungen abhängig von Anzahl Arrays (→ wenig flexibel !)
 - wenige, ausgewählte Designs in Effizienzvergleich einbezogen (**lokale** Optima)
 - Abhängigkeit vom Optimalitätskriterium (z.B. D-, A-Kriterium) nicht untersucht
- **Ziele:**
 - 1) Herleitung **global** optimaler Versuchspläne (+ **analytischer** Beweis!)
 - 2) **flexible** Anzahl von N, I, k
 - 3) Untersuchung der **Robustheit** gegenüber **Kriterium**

Äquivalenztheorem (Φ_p -Optimalität)

- $m = 2 \cdot \binom{lk}{2}$: Anzahl theoretisch möglicher verschiedener Designpunkte x_1, \dots, x_m

- $\chi = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$: Designraum (diskrete Menge von Designpunkten)

- Ω : Menge aller möglichen **diskreten** Designs

$$\xi = \left\{ \begin{matrix} x_1, \dots, x_m \\ p_1, \dots, p_m \end{matrix} \right\}; 0 \leq p_i \leq 1 \forall i = 1, \dots, m; \sum_{i=1}^m p_i = 1$$

- $\mathbf{M} := \sum_{i=1}^m p_i \mathbf{x}_i \mathbf{x}_i^T$ ($\sim s \times s$) **Momentenmatrix**

- ξ^* : Kandidatendesign; *schätzbar* bzgl. $\mathbf{K}^T \theta$

Theorem 1 (Äquivalenztheorem für Matrizenmittelwerte, *generalisiert*):

- (a) ξ^* ist ϕ_p -**optimal** für $\mathbf{K}^T \theta$ dann und nur dann, wenn eine *g-Inverse* \mathbf{G} von \mathbf{M} existiert, die die **Normalitätsungleichung**

$$\mathbf{x}^T \mathbf{G} \mathbf{K} (\mathbf{K}^T \mathbf{G} \mathbf{K})^+ (\mathbf{K}^T \mathbf{G} \mathbf{K})^{1-p} (\mathbf{K}^T \mathbf{G} \mathbf{K})^+ \mathbf{K}^T \mathbf{G}^T \mathbf{x} \leq \text{Tr}[(\mathbf{K}^T \mathbf{G} \mathbf{K})^+ (\mathbf{K}^T \mathbf{G} \mathbf{K})^{1-p}] \quad (\mathbf{N})$$

für **alle** Designpunkte $x \in \chi$ erfüllt (für alle $p \in (-\infty; 1]$).

- (b) Für die **Trägerpunkte** von ξ^* gilt **Gleichheit** in **(N)**.

Within Cell Line Swap Design (CLS)

Trägerpunkte vergleichen jeweils **zwei Treatments** bei Verwendung **derselben Zelllinie**.

Theorem 2 (Optimalität des Within Cell Line Swap Designs):

Das Design **CLS** ist

- (a) D-optimal ($p=0$)
- (b) A-optimal ($p=-1$)
- (c) E-optimal ($p=-\infty$)

für das Schätzen der folgenden (linearen) Kontraste:

- (1) **Treatmenteffekt** (für alle $k, l \in N$)
- (2) **Interaktionseffekt** von Treatment und Zelllinie (**für** $k \geq l ; k, l \in N$)
- (3) **Kombination** von **Treatmenteffekt** und **Interaktionseffekt** (für alle $k, l \in N$)

Beweisskizze (Treatmenteffekt; ϕ_p)

Design CLS:
$$\mathbf{M} = c_1 \cdot \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_k \otimes (\mathbf{A}_1^T \mathbf{A}_1) \end{bmatrix}, \quad \mathbf{G} = \frac{1}{c_1} \cdot \begin{bmatrix} \mathbf{A}^- & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_k \otimes \left(\frac{1}{4l^2} \mathbf{A}_1^T \mathbf{A}_1 \right) \end{bmatrix}, \quad \mathbf{A}_1 = \begin{bmatrix} \tilde{\mathbf{A}}_1 \\ -\tilde{\mathbf{A}}_1 \end{bmatrix}.$$

Kontrastmatrix:
$$\mathbf{K} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{1}_k^T \otimes \tilde{\mathbf{A}}_1 \end{bmatrix}, \quad \tilde{\mathbf{A}}_1 \sim \binom{l}{2} \times l.$$

$$\mathbf{K}^T \mathbf{G}^T = \frac{1}{c_1} \cdot \begin{bmatrix} \mathbf{0} & \mathbf{1}_k^T \otimes \frac{1}{2l} \tilde{\mathbf{A}}_1 \end{bmatrix}, \quad \mathbf{K}^T \mathbf{G} \mathbf{K} = c_2 \cdot \tilde{\mathbf{A}}_1 \tilde{\mathbf{A}}_1^T, \quad (\mathbf{K}^T \mathbf{G} \mathbf{K})^+ = c_3 \cdot \tilde{\mathbf{A}}_1 \tilde{\mathbf{A}}_1^T$$

Designpunkte von CLS:
$$\mathbf{x} = [1, -1, x_{11}, \dots, x_{lk}]^T \quad \text{mit}$$

$$x_{uv} = \begin{cases} +1, & \text{Treatment } u \leftrightarrow \text{Zelllinie } v \text{ (g)} \\ -1, & \text{Treatment } u \leftrightarrow \text{Zelllinie } v \text{ (r)} \\ 0, & \text{sonst} \end{cases}; \quad u = 1, \dots, l; \quad v = 1, \dots, k \quad \text{und} \quad \bar{x}_{\cdot w} = \frac{1}{l} \sum_{u=1}^l x_{uw}.$$

Dann gilt:
$$\mathbf{x}^T \mathbf{G} \mathbf{K} (\mathbf{K}^T \mathbf{G} \mathbf{K})^+ (\mathbf{K}^T \mathbf{G} \mathbf{K})^{1-p} (\mathbf{K}^T \mathbf{G} \mathbf{K})^+ \mathbf{K}^T \mathbf{G}^T \mathbf{x} \leq \text{Tr} [(\mathbf{K}^T \mathbf{G} \mathbf{K})^+ (\mathbf{K}^T \mathbf{G} \mathbf{K})^{1-p}]$$

$$\Leftrightarrow \sum_{u=1}^l \sum_{v=1}^k \left[x_{uv} \cdot \sum_{w=1}^k (x_{uw} - \bar{x}_{\cdot w}) \right] \leq 2. \quad (\Delta)$$

Gültigkeit von (Δ) für alle \mathbf{x} (Gleichheit für Trägerpunkte von CLS) kann über *Fallunterscheidung* gezeigt werden.

Within Treatment Swap Design (TS)

Trägerpunkte vergleichen jeweils **zwei Zelllinien** bei Verwendung **desselben Treatment**.

Theorem 3 (Optimalität des Within Treatment Swap Designs)

Das Design **TS** ist für $I \geq k$

(a) D-optimal ($p=0$)

(b) A-optimal ($p=-1$)

(c) E-optimal ($p=-\infty$)

für das Schätzen des linearen Kontrastes für den **Interaktionseffekt** von Treatment und Zelllinie.

Beweis: Analog zum Beweis von Theorem 2. ■

Diskussion

- Für zahlreiche *dreifaktorielle* Designsituationen konnten **global optimale** cDNA-Microarray-Designs gefunden und deren Optimalität **analytisch** bewiesen werden.
- **Robustheit**: Wahl des optimalen Designs **unabhängig** vom **Designkriterium !!**
- N , l , k nicht auf bestimmte Werte beschränkt („hohe Flexibilität“).
- In manchen Situationen ist die Lösung abhängig von der *Zahl der Zelllinien* bzw. der *Zahl der Treatments*.
- **Praxis**:
 - Empfehlungen für Wahl des effizientesten cDNA-Microarray-Designs
 - Möglichkeit der Einsparung finanzieller Ressourcen

Ausblick:

- Beweis der Optimalität von **TS** bzgl. Zelllinienneneffekt
- Betrachtung **anderer Kontraste** (z.B. Helmert-Kontraste)
- Betrachtung **vier- und mehrfaktorieller** Designsituationen
- Optimales Design **modellabhängig** ??
- Berücksichtigung von **Korrelationsstrukturen** zwischen Genen

Literatur

Duggan, D., Bittner, M., Chen, Y., Meltzer, P., Trent, J. (1999):

Expression profiling using cDNA Microarrays. Nature Genetics 21 (1 Suppl): 10-14.

Landgrebe, J., Bretz, F., Brunner, E. (2006): *Efficient design and analysis of two colour factorial microarray experiments.* Computational Statistics and Data Analysis 2006; 50: 499-517.

Pukelsheim, F. (1993): *Optimal Design of Experiments.* Wiley, New York.

Searle, S. R. (1971): *Linear Models.* Wiley, New York.

Searle, S. R. (1982): *Matrix Algebra useful for Statistics.* Wiley, New York.