



# Die 'Daten-Medizin'

**Karsten Borgwardt**

Department Biosysteme, ETH Zürich

Esslingen, 13.4.2018

# Maschinelles Lernen und personalisierte Medizin

## Ziele

- Mit **maschinellem Lernen** sollen **statistische Abhängigkeiten in großen Datensätzen** erkannt werden.

# Maschinelles Lernen und personalisierte Medizin

## Ziele

- Mit **maschinellem Lernen** sollen **statistische Abhängigkeiten** in großen Datensätzen erkannt werden.



# Maschinelles Lernen und personalisierte Medizin

## Ziele

- Mit **maschinellem Lernen** sollen **statistische Abhängigkeiten in großen Datensätzen** erkannt werden.



- **Personalisierte Medizin** versucht Gesundheitsdaten für **verbesserte Diagnosen, Vorhersagen und Therapieentscheidungen** zu nutzen, individuell angepasst auf die Eigenschaften eines jeden Patienten.

# Maschinelles Lernen in der Medizin

## Hauptthemen

# Maschinelles Lernen in der Medizin

## Hauptthemen

- Automatisierung der Diagnostik

### Original Investigation

December 12, 2017

## Diagnostic Assessment of Deep Learning Algorithms for Detection of Lymph Node Metastases in Women With Breast Cancer

Babak Ehteshami Bejnordi, MS<sup>1</sup>; Mitko Veta, PhD<sup>2</sup>; Paul Johannes van Diest, MD, PhD<sup>3</sup>; [et al](#)

[» Author Affiliations](#) | [Article Information](#)

JAMA. 2017;318(22):2199-2210. doi:10.1001/jama.2017.14585

# Maschinelles Lernen in der Medizin

## Hauptthemen

- Automatisierung der Diagnostik
- Entdeckung von Biomarkern

## Machine learning of neural representations of suicide and emotion concepts identifies suicidal youth

Marcel Adam Just , Lisa Pan, Vladimir L. Cherkassky, Dana L. McMakin, Christine Cha, Matthew K. Nock & David Brent

*Nature Human Behaviour* **1**, 911–919 (2017)

doi:10.1038/s41562-017-0234-y

[Download Citation](#)

Received: 06 February 2017

Accepted: 04 October 2017

Published online: 30 October 2017

# Maschinelles Lernen in der Medizin

## Hauptthemen

- Automatisierung der Diagnostik
- Entdeckung von Biomarkern
- Biomedizinisches Datenmanagement



**Roche to buy Flatiron Health for \$1.9 billion to expand cancer care ...**

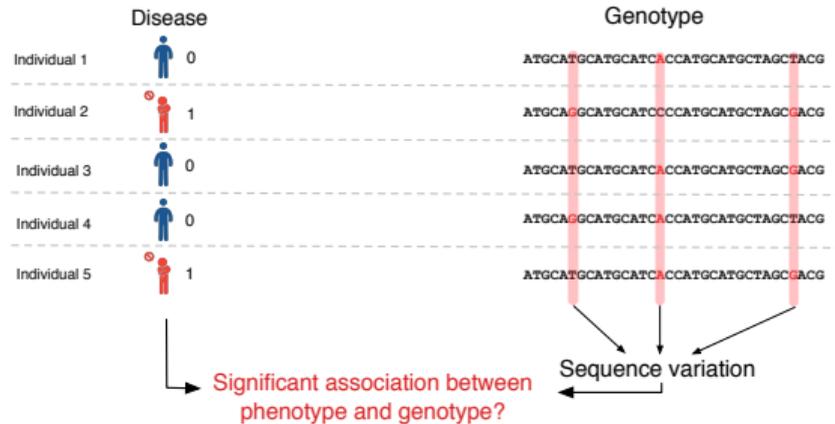
Reuters - 15.02.2018

Roche to buy Flatiron Health for \$1.9 billion to expand cancer care portfolio ... S) said on Thursday it would buy the rest of U.S. cancer data company Flatiron Health for \$1.9 billion to speed development of cancer medicines and support its efforts to ... Privately held Flatiron, backed by Alphabet Inc (GOOGL).

# Maschinelles Lernen in der Medizin

## Hauptthemen

- Automatisierung der Diagnostik
  
- Entdeckung von Biomarkern
  - 1 Methodenentwicklung:  
Kombinatorische Assoziationsuche
  
- Biomedizinisches Datenmanagement



# Maschinelles Lernen in der Medizin

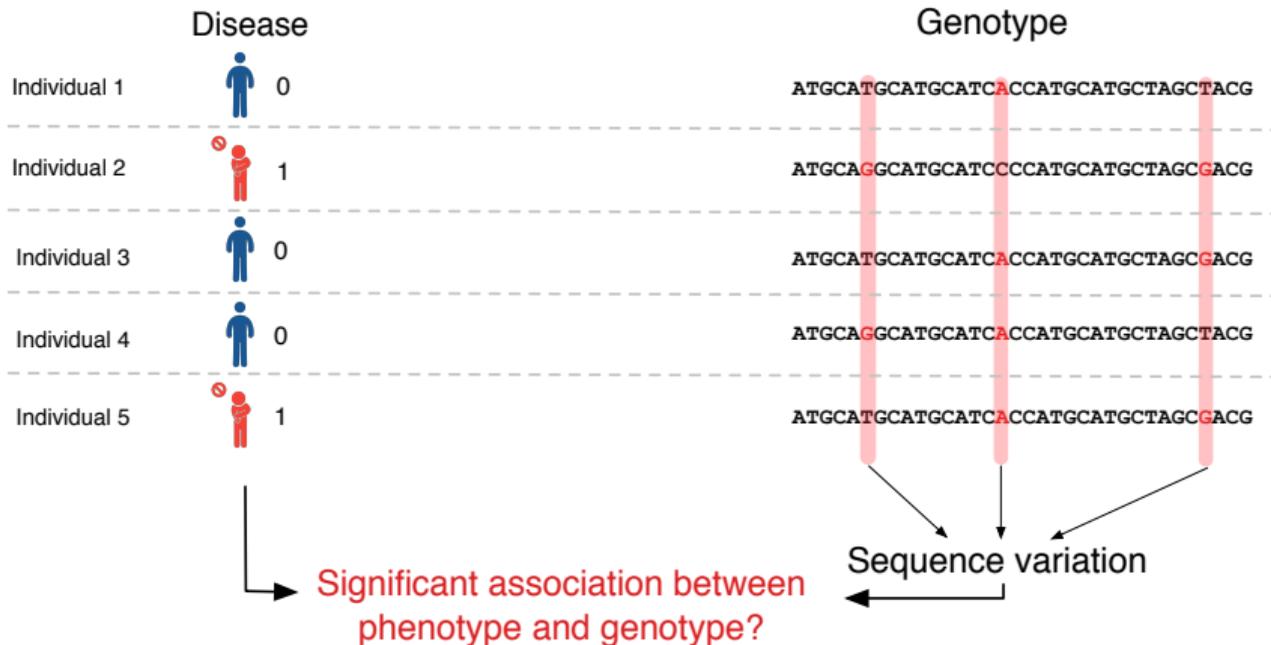
## Hauptthemen

- Automatisierung der Diagnostik
  
- Entdeckung von Biomarkern
  - 1 Methodenentwicklung:  
Kombinatorische Assoziationsuche
  
- Biomedizinisches Datenmanagement
  - 2 Softwareentwicklung:  
easyGWAS



## Kombinatorische Assoziationsuche

# Assoziationsuche: Korrelationen zwischen Phänotyp und Genotyp



Eine **genomweite Assoziationsstudie (GWAS)** untersucht, ob Varianten im Genom (in Form von Einzelnukleotid-Polymorphismen, SNPs) mit Veränderungen des Phänotyps korreliert sind.

## Assoziationsuche: Fehlende Erbllichkeit

- Seit 2001: > 59,000 Phänotyp-assoziierte Genomregionen durch GWAS (GWAS Catalog, 6.3.2018)
- Problem: Erklärte Varianz des Phänotyps enttäuschend gering

Vol 461 | 8 October 2009 | doi:10.1038/nature08494

nature

### REVIEWS

## Finding the missing heritability of complex diseases

Teri A. Manolio<sup>1</sup>, Francis S. Collins<sup>2</sup>, Nancy J. Cox<sup>3</sup>, David B. Goldstein<sup>4</sup>, Lucia A. Hindorf<sup>5</sup>, David J. Hunter<sup>6</sup>, Mark I. McCarthy<sup>7</sup>, Erin M. Ramos<sup>8</sup>, Lon R. Cardon<sup>9</sup>, Aravinda Chakravarti<sup>10</sup>, Judy H. Cho<sup>10</sup>, Alan E. Guttmacher<sup>1</sup>, Augustine Kong<sup>11</sup>, Leonid Kruglyak<sup>12</sup>, Elaine Mardis<sup>13</sup>, Charles N. Rotimi<sup>14</sup>, Montgomery Slatkin<sup>15</sup>, David Valle<sup>9</sup>, Alice S. Whittemore<sup>16</sup>, Michael Boehnke<sup>17</sup>, Andrew G. Clark<sup>18</sup>, Evan E. Eichler<sup>19</sup>, Greg Gibson<sup>20</sup>, Jonathan L. Haines<sup>21</sup>, Trudy F. C. Mackay<sup>22</sup>, Steven A. McCarroll<sup>23</sup> & Peter M. Visscher<sup>24</sup>

Genome-wide association studies have identified hundreds of genetic variants associated with complex human diseases and traits, and have provided valuable insights into their genetic architecture. Most variants identified so far confer relatively

## Assoziationsuche: Fehlende Erbllichkeit

- Seit 2001: > 59,000 Phänotyp-assozierte Genomregionen durch GWAS (GWAS Catalog, 6.3.2018)
- Problem: Erklärte Varianz des Phänotyps enttäuschend gering

Vol 461|8 October 2009|doi:10.1038/nature08494

nature

### REVIEWS

#### Finding the missing heritability of complex diseases

Teri A. Manolio<sup>1</sup>, Francis S. Collins<sup>2</sup>, Nancy J. Cox<sup>3</sup>, David B. Goldstein<sup>4</sup>, Lucia A. Hindorf<sup>5</sup>, David J. Hunter<sup>6</sup>, Mark I. McCarthy<sup>7</sup>, Erin M. Ramos<sup>8</sup>, Lon R. Cardon<sup>9</sup>, Aravinda Chakravarti<sup>10</sup>, Judy H. Cho<sup>10</sup>, Alan E. Guttmacher<sup>1</sup>, Augustine Kong<sup>11</sup>, Leonid Kruglyak<sup>12</sup>, Elaine Mardis<sup>13</sup>, Charles N. Rotimi<sup>14</sup>, Montgomery Slatkin<sup>15</sup>, David Valle<sup>9</sup>, Alice S. Whittemore<sup>16</sup>, Michael Boehnke<sup>17</sup>, Andrew G. Clark<sup>18</sup>, Evan E. Eichler<sup>19</sup>, Greg Gibson<sup>20</sup>, Jonathan L. Haines<sup>21</sup>, Trudy F. C. Mackay<sup>22</sup>, Steven A. McCarroll<sup>23</sup> & Peter M. Visscher<sup>24</sup>

Genome-wide association studies have identified hundreds of genetic variants associated with complex human diseases and traits, and have provided valuable insights into their genetic architecture. Most variants identified so far confer relatively

- Mögliche Gründe:
  - Polygenische Architektur komplexer Krankheiten
  - Kleine Effektstärken
  - Unvollständige Integration wichtiger genetischer, epigenetischer und nichtgenetischer Eigenschaften

## Assoziationsuche: Fehlende Erbllichkeit

- Seit 2001: > 59,000 Phänotyp-assozierte Genomregionen durch GWAS (GWAS Catalog, 6.3.2018)
- Problem: Erklärte Varianz des Phänotyps enttäuschend gering

Vol 461|8 October 2009|doi:10.1038/nature08494

nature

### REVIEWS

#### Finding the missing heritability of complex diseases

Teri A. Manolio<sup>1</sup>, Francis S. Collins<sup>2</sup>, Nancy J. Cox<sup>3</sup>, David B. Goldstein<sup>4</sup>, Lucia A. Hindorf<sup>5</sup>, David J. Hunter<sup>6</sup>, Mark I. McCarthy<sup>7</sup>, Erin M. Ramos<sup>8</sup>, Lon R. Cardon<sup>9</sup>, Aravinda Chakravarti<sup>10</sup>, Judy H. Cho<sup>10</sup>, Alan E. Guttmacher<sup>1</sup>, Augustine Kong<sup>11</sup>, Leonid Kruglyak<sup>12</sup>, Elaine Mardis<sup>13</sup>, Charles N. Rotimi<sup>14</sup>, Montgomery Slatkin<sup>15</sup>, David Valle<sup>9</sup>, Alice S. Whittemore<sup>16</sup>, Michael Boehnke<sup>17</sup>, Andrew G. Clark<sup>18</sup>, Evan E. Eichler<sup>19</sup>, Greg Gibson<sup>20</sup>, Jonathan L. Haines<sup>21</sup>, Trudy F. C. Mackay<sup>22</sup>, Steven A. McCarroll<sup>23</sup> & Peter M. Visscher<sup>24</sup>

Genome-wide association studies have identified hundreds of genetic variants associated with complex human diseases and traits, and have provided valuable insights into their genetic architecture. Most variants identified so far confer relatively

- Mögliche Gründe:
  - Polygenische Architektur komplexer Krankheiten → Epistasis
  - Kleine Effektstärken
  - Unvollständige Integration wichtiger genetischer, epigenetischer und nichtgenetischer Eigenschaften

# Assoziationsuche: Fehlende Erbllichkeit

## Epistase als mögliche Ursache

- Die meisten aktuellen Untersuchungen vernachlässigen interaktive Effekte zwischen SNPs
- Bedarf an Methoden für **kombinatorische Assoziationsuche**

Mackay and Moore *Genome Medicine* 2014, 6:42  
<http://genomemedicine.com/content/6/6/42>



### COMMENT

## Why epistasis is important for tackling complex human disease genetics

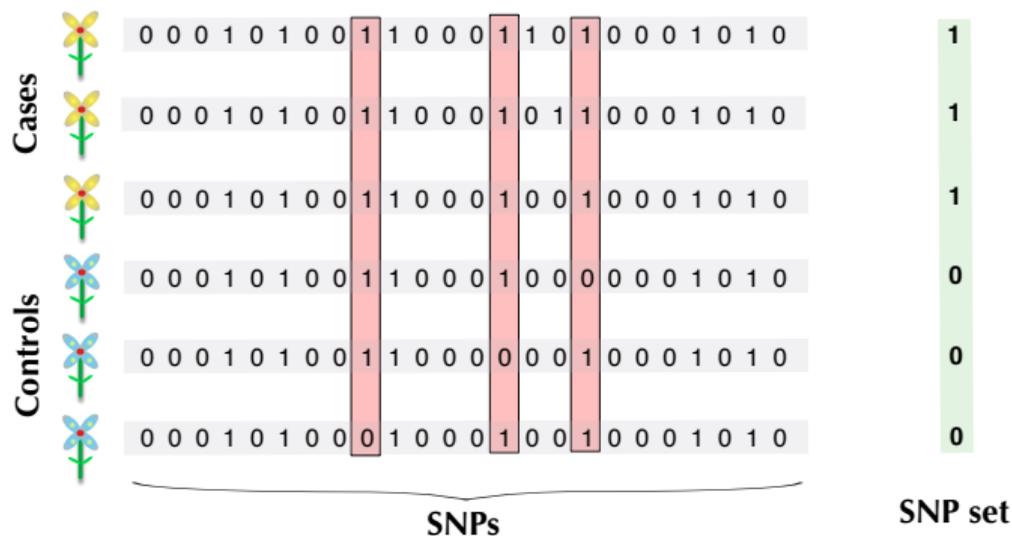
Trudy FC Mackay<sup>1\*</sup> and Jason H Moore<sup>2</sup>

#### Editorial summary

Epistasis has been dismissed by some as having little role in the genetic architecture of complex human disease. The authors argue that this view is the result

and the effects of alleles at these loci are highly sensitive to the environmental circumstances to which the individuals are exposed. Quantitative variation in phenotypes and disease risk must result in part from the perturbation of highly dynamic, interconnected and non-linear net-

## Kombinatorische Assoziationsuche



- Informatische Herausforderung: Kombinatorische Explosion der Anzahl der Kandidatenmengen
- Statistische Herausforderung: Kombinatorische Explosion der Anzahl der Assoziationstests

# Kombinatorische Assoziationsuche

## Das Problem der Alphafehler-Kumulierung

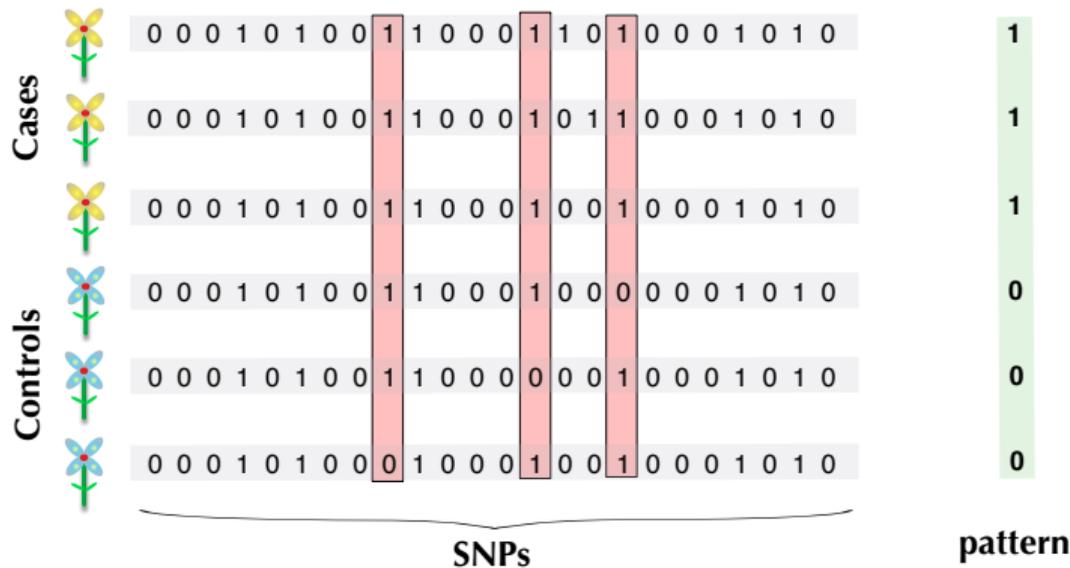
- Was, wenn wir Assoziationen von Gruppen von  $s$  SNPs mit dem Phänotyp betrachten?
- Dies führt **multiplem Hypothesentesten**: Jede der  $k$  SNP-Mengen entspricht einer zu testenden Hypothese ( $k \in O(f^s)$ ), wobei  $f$  die Anzahl der SNPs ist.
- Falls nicht berücksichtigt, könnten selbst bei zufällig generierten Daten  $\alpha$  Prozent aller SNP-Mengen als assoziiert betrachtet werden.
- Es ist zwingend notwendig, multiples Testen zu berücksichtigen, z. B. durch Kontrolle der **family-wise error rate**!
- Falls berücksichtigt, z. B. durch Bonferroni-Korrektur ( $\frac{\alpha}{k}$ ), könnte die **gesamte Teststärke verloren gehen**.

# Kombinatorische Assoziationsuche

## Das Problem der Alphafehler-Kumulierung

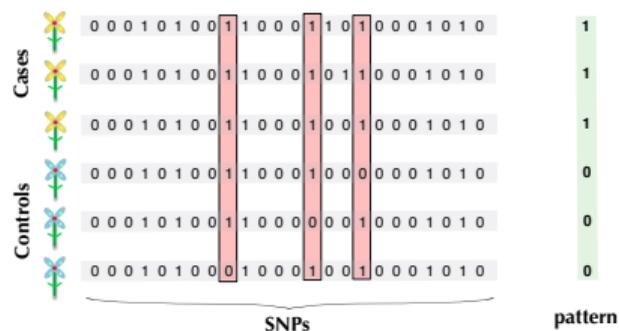
- Was, wenn wir Assoziationen von Gruppen von  $s$  SNPs mit dem Phänotyp betrachten?
- Dies führt **multiplem Hypothesentesten**: Jede der  $k$  SNP-Mengen entspricht einer zu testenden Hypothese ( $k \in O(f^s)$ ), wobei  $f$  die Anzahl der SNPs ist.
- Falls nicht berücksichtigt, könnten selbst bei zufällig generierten Daten  $\alpha$  Prozent aller SNP-Mengen als assoziiert betrachtet werden.
- Es ist zwingend notwendig, multiples Testen zu berücksichtigen, z. B. durch Kontrolle der **family-wise error rate**!
- Falls berücksichtigt, z. B. durch Bonferroni-Korrektur ( $\frac{\alpha}{k}$ ), könnte die **gesamte Teststärke verloren gehen**.
- **Lange als unlösbares Dilemma angesehen**

# Kombinatorische Assoziationsuche als Data-Mining-Problem



- Variablenselektion: Finde die Variablen, die Klassen von Objekten unterscheiden
- Kombinationssuche: Finde **Kombinationen binärer Variablen** höherer Ordnung, um die Klassen zu unterscheiden

# Kombinatorische Assoziationsuche als Data-Mining-Problem



## Kombination

- In unserem Datensatz  $D$  sei der  $i$ -te von  $n$  Patienten dargestellt durch einen binären Vektor  $\mathbf{d}^{(i)} \in \{0, 1\}^f$  und ein Klassenlabel  $y_i \in \{0, 1\}$ .
- Wir wählen eine Teilmenge  $\mathcal{S}$  aller Variablen  $\mathcal{F}$  in einem Datensatz:  $\mathcal{S} \subseteq \mathcal{F}$ .
- Dann enthält ein Objekt  $\mathbf{d}^{(i)}$  die Kombination  $\mathcal{S}$  genau dann, wenn  $\prod_{t \in \mathcal{S}} d^{(i)}(t) = 1$ .



## Suche nach signifikanten Kombinationen

### Tarones Trick

#### ■ Kontingenztabelle

	Kombination vorhanden	Kombination fehlt	
$y=0$	$a$	$n_1 - a$	$n_1$
$y=1$	$x - a$	$n - n_1 - x + a$	$n - n_1$
	$x$	$n - x$	$n$

- Häufig wird der **exakte Test nach Fisher** gewählt, um zu testen, ob eine Kombination in einer von zwei Klassen überrepräsentiert ist.
- Die übliche Methode zur Berechnung der  $p$ -Werte basiert auf der hypergeometrischen Verteilung und setzt voraus, dass die Randhäufigkeiten  $(x, n_1, n)$  fix sind.

## Suche nach signifikanten Kombinationen

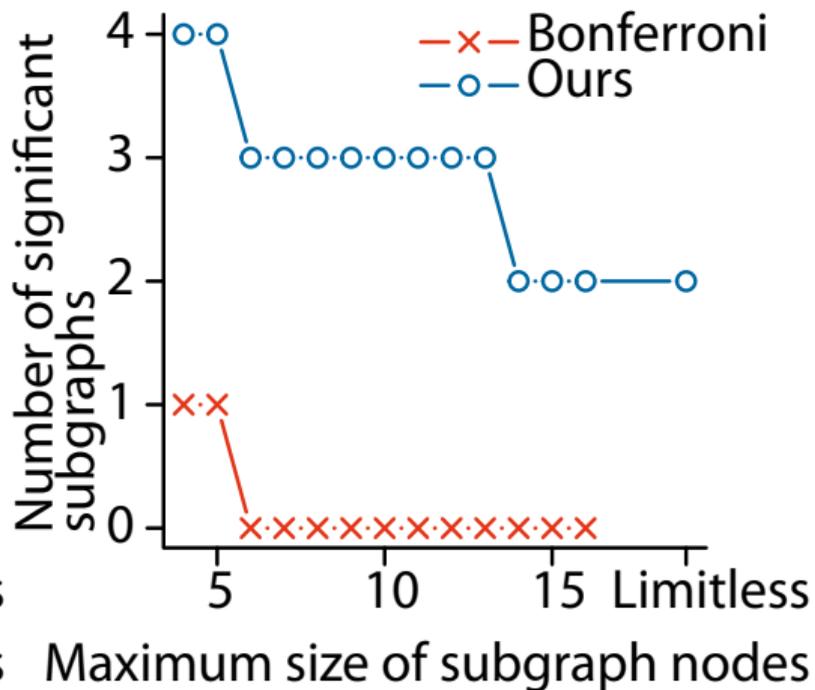
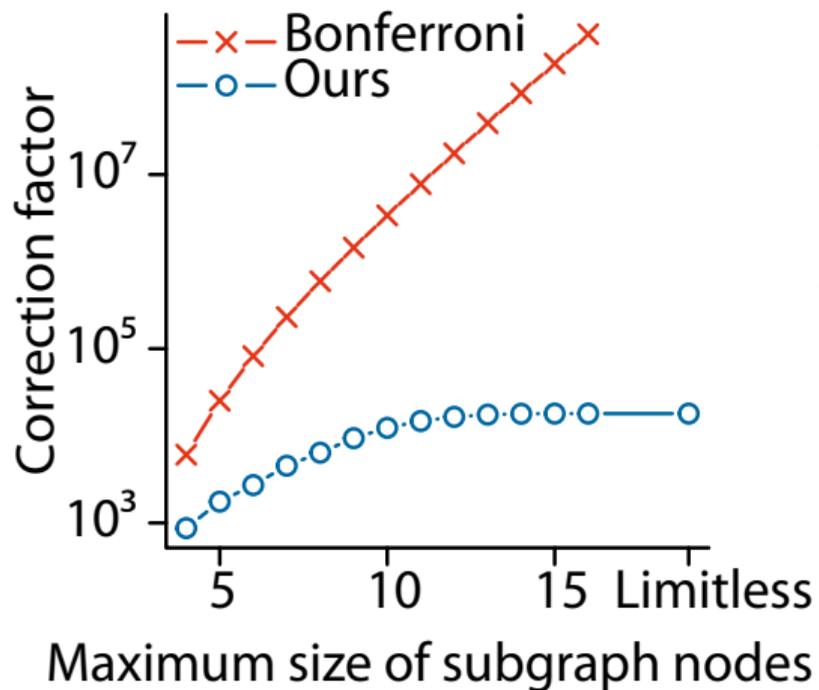
### Tarones Trick

#### ■ Kontingenztabelle

	Kombination vorhanden	Kombination fehlt	
y=0	$a$	$n_1 - a$	$n_1$
y=1	$x - a$	$n - n_1 - x + a$	$n - n_1$
	$x$	$n - x$	$n$

- Tarone (1990) erkannte, dass es bei diskreten Teststatistiken, wie z. B. beim exakten Test nach Fisher, einen **minimalen  $p$ -Wert** gibt, den eine Kombination annehmen kann.
- Es gibt viele **nicht testbare Hypothesen**, deren minimaler  $p$ -Wert nicht kleiner ist als  $\frac{\alpha}{k}$ .
- Nur die übrigen  $m(k)$  **testbaren Hypothesen** können überhaupt signifikant sein.
- Man kann für  $m(k)$  **statt  $k$**  Tests korrigieren. Da oft  $m(k) \ll k$ , verbessert das die Teststärke sehr.

## Beispiel: PTC Datensatz (Helma et al., 2001)



## Suche nach signifikanten Kombinationen

### Tarones Methode (1990)

- Sei  $k$  die Anzahl der Tests, für die korrigiert werden soll.
- $m(k)$  ist die Anzahl der testbaren Hypothesen zum Signifikanzniveau  $\frac{\alpha}{k}$ .
- $m(k)$  ist eine Funktion von  $k$ , und wir benötigen  $k \geq m(k)$ , um für alle testbaren Hypothesen zu korrigieren.
- Das Optimierungsproblem lautet dann

min  $k$

so dass  $k \geq m(k)$

## Suche nach signifikanten Kombinationen

### Tarones Methode (1990)

- Sei  $k$  die Anzahl der Tests, für die korrigiert werden soll.
- $m(k)$  ist die Anzahl der testbaren Hypothesen zum Niveau  $\frac{\alpha}{k}$ .

**procedure** Tarone

$k := 1;$

**while**  $k < m(k)$  **do**

$k := k + 1;$

**return**  $k$

## Suche nach signifikanten Kombinationen

### Tarones Methode (1990)

- Sei  $k$  die Anzahl der Tests, für die korrigiert werden soll.
- $m(k)$  ist die Anzahl der testbaren Hypothesen zum Niveau  $\frac{\alpha}{k}$ .

**procedure** Tarone

$k := 1;$

**while**  $k < m(k)$  **do**

$k := k + 1;$

**return**  $k$

- Wie kann  $m(k)$  effizient berechnet werden, ohne alle  $O(f^s)$  möglichen Hypothesen zu durchlaufen?

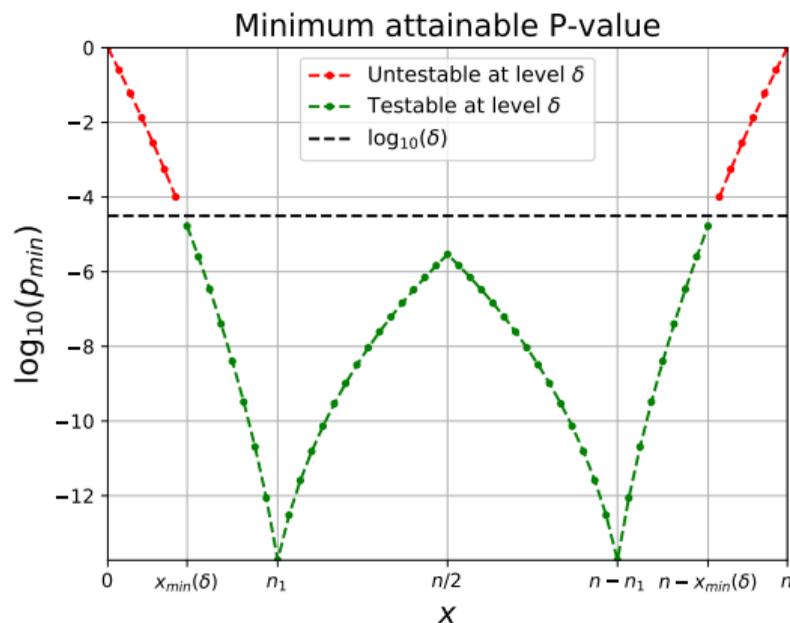
## Suche nach signifikanten Kombinationen

### Data-Mining Problem

- Wie soll man  $m(k)$  effizient berechnen, ohne alle  $O(f^s)$  möglichen Hypothesen zu durchlaufen?
- Lösung: Die Häufigkeit einer Kombination bestimmt deren minimalen  $p$ -Wert.
- Algorithmen zur Suche von häufigen Kombinationen aus dem Data-Mining können benutzt werden, um alle Kombinationen zu finden, deren minimaler  $p$ -Wert einen gewissen Schwellwert unterschreitet (Terada et al., PNAS 2013):
  - frequent itemset mining( $D, \theta$ ) findet alle Kombinationen in einem Datensatz  $D$  mit einer Häufigkeit von mindestens  $\theta$ .

# Suche nach signifikanten Kombinationen

- Häufigkeit versus minimaler  $p$ -Wert



## Suche nach signifikanten Kombinationen

### Tarones Methode mit frequent itemset mining

- Sei  $k$  die Anzahl der Tests für die korrigiert werden soll.
- $m(k)$  ist die Anzahl der testbaren Hypothesen mit Signifikanzniveau  $\frac{\alpha}{k}$ .

**procedure** Tarone( $D, \alpha$ )

$k := 1$ ;

**while**  $k < m(k)$  **do**

$k := k + 1$ ;

$m(k) :=$  frequent itemset mining( $D, \phi(\frac{\alpha}{k})$ );

**return**  $k$

## Suche nach signifikanten Kombinationen

### Tarones Methode mit frequent itemset mining

- Sei  $k$  die Anzahl der Tests für die korrigiert werden soll.
- $m(k)$  ist die Anzahl der testbaren Hypothesen mit Signifikanzniveau  $\frac{\alpha}{k}$ .

**procedure** Tarone( $D, \alpha$ )

$k := 1$ ;

**while**  $k < m(k)$  **do**

$k := k + 1$ ;

$m(k) := \text{frequent itemset mining}(D, \phi(\frac{\alpha}{k}))$ ;

**return**  $k$

- Hinweis:  $\phi(\frac{\alpha}{k})$  ist die minimale Häufigkeit einer Kombination, die mit Niveau  $\frac{\alpha}{k}$  testbar ist.

## Suche nach signifikanten Kombinationen

### Tarones Methode mit frequent itemset mining

- Sei  $k$  die Anzahl der Tests für die korrigiert werden soll.
- $m(k)$  ist die Anzahl der testbaren Hypothesen mit Signifikanzniveau  $\frac{\alpha}{k}$ .

**procedure** Tarone( $D, \alpha$ )

$k := 1;$

**while**  $k < m(k)$  **do**

$k := k + 1;$

$m(k) := \text{frequent itemset mining}(D, \phi(\frac{\alpha}{k}));$

**return**  $k$

- Hinweis:  $\phi(\frac{\alpha}{k})$  ist die minimale Häufigkeit einer Kombination, die mit Niveau  $\frac{\alpha}{k}$  testbar ist.
- Für kleine  $k$  ist  $\phi(\frac{\alpha}{k})$  klein. Frequent itemset mining wird extrem rechenaufwändig!

# Von der Suche nach signifikanten Kombinationen zur kombinatorischen Assoziationsuche

Fragen, die 2014 noch unbeantwortet waren

- 1 Wie kann das optimale  $k$  **effizient** gefunden werden? (SDM 2015)

# Von der Suche nach signifikanten Kombinationen zur kombinatorischen Assoziationsuche

## Fragen, die 2014 noch unbeantwortet waren

- 1 Wie kann das optimale  $k$  **effizient** gefunden werden? (SDM 2015)
  - Wir präsentierten eine effiziente Suchstrategie mit Kriterium zum vorzeitigen Abbruch (falls  $m(k) > k$ ).

# Von der Suche nach signifikanten Kombinationen zur kombinatorischen Assoziationsuche

## Fragen, die 2014 noch unbeantwortet waren

- 1 Wie kann das optimale  $k$  **effizient** gefunden werden? (SDM 2015)
  - Wir präsentierten eine effiziente Suchstrategie mit Kriterium zum vorzeitigen Abbruch (falls  $m(k) > k$ ).
- 2 Kombinationen stehen in Untermengen-/Obermengenbeziehungen. Wie kann diese **Abhängigkeit der Tests** berücksichtigt werden? (KDD 2015)

# Von der Suche nach signifikanten Kombinationen zur kombinatorischen Assoziationsuche

## Fragen, die 2014 noch unbeantwortet waren

- 1 Wie kann das optimale  $k$  **effizient** gefunden werden? (SDM 2015)
  - Wir präsentierten eine effiziente Suchstrategie mit Kriterium zum vorzeitigen Abbruch (falls  $m(k) > k$ ).
- 2 Kombinationen stehen in Untermengen-/Obermengenbeziehungen. Wie kann diese **Abhängigkeit der Tests** berücksichtigt werden? (KDD 2015)
  - Wir führen Westfall-Young Permutationen durch, um die Abhängigkeit zu berücksichtigen.
  - Durch dynamische Aktualisierung des Häufigkeitsschwellwerts benötigen wir auch bei 10,000 Permutationen lediglich 1 einzige Anwendung von frequent itemset mining.

## Von der Suche nach signifikanten Kombinationen zur kombinatorischen Assoziationsuche

### Fragen, die 2014 noch unbeantwortet waren

- 3 Können wir Effizienz und Teststärke beibehalten, wenn wir **kategorische Kovariaten** wie Alter und Geschlecht berücksichtigen? (NIPS 2016)

# Von der Suche nach signifikanten Kombinationen zur kombinatorischen Assoziationsuche

## Fragen, die 2014 noch unbeantwortet waren

- 3 Können wir Effizienz und Teststärke beibehalten, wenn wir **kategorische Kovariaten** wie Alter und Geschlecht berücksichtigen? (NIPS 2016)
  - Wir erweiterten Tarones Trick auf den Cochran-Mantel-Haenszel-Test für stratifizierte Kontingenztabelle.

## Von der Suche nach signifikanten Kombinationen zur kombinatorischen Assoziationsuche

### Fragen, die 2014 noch unbeantwortet waren

- 3 Können wir Effizienz und Teststärke beibehalten, wenn wir **kategorische Kovariaten** wie Alter und Geschlecht berücksichtigen? (NIPS 2016)
  - Wir erweiterten Tarones Trick auf den Cochran-Mantel-Haenszel-Test für stratifizierte Kontingenztabelle.
- 4 Können wir **neue kombinatorische Assoziationsverfahren** entwickeln, die auf Tarones Trick basieren? (ISMB 2015, OUP Bioinformatics 2017, ISMB 2018)

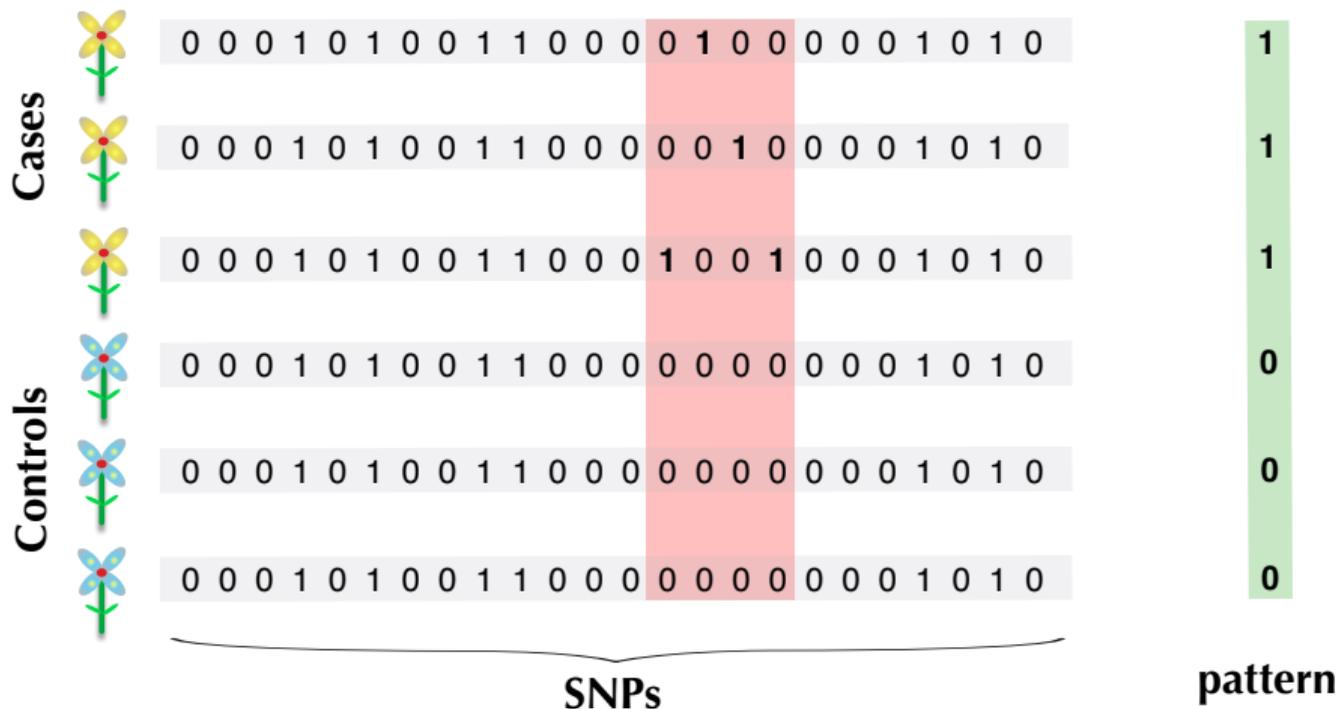
## Von der Suche nach signifikanten Kombinationen zur kombinatorischen Assoziationsuche

### Fragen, die 2014 noch unbeantwortet waren

- 3 Können wir Effizienz und Teststärke beibehalten, wenn wir **kategorische Kovariaten** wie Alter und Geschlecht berücksichtigen? (NIPS 2016)
  - Wir erweiterten Tarones Trick auf den Cochran-Mantel-Haenszel-Test für stratifizierte Kontingenztabellen.
- 4 Können wir **neue kombinatorische Assoziationsverfahren entwickeln, die auf Tarones Trick basieren?** (ISMB 2015, OUP Bioinformatics 2017, ISMB 2018)

# Kombinatorische Assoziationskartierung für die Erkennung von genetischer Heterogenität

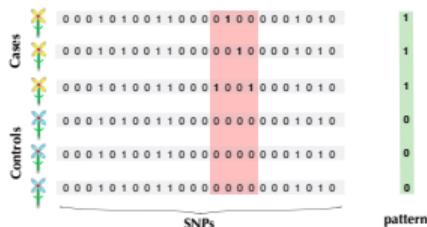
## Erkennung von genetischer Heterogenität



# Erkennung von genetischer Heterogenität

## Genetische Heterogenität

- Genetische Heterogenität bezeichnet das Phänomen, dass mehrere verschiedene Gene oder Sequenzvarianten zu demselben Phänotyp führen können.



- Spezielles Problem, das hier betrachtet wird: Finde ein Intervall im Genom, so dass die Eigenschaft, dass
  - eine seltene Variante,
  - ein rezessiver Genotyp, oder
  - ein Minderheitsallel
 in diesem Intervall liegt, mit dem Phänotyp assoziiert ist.

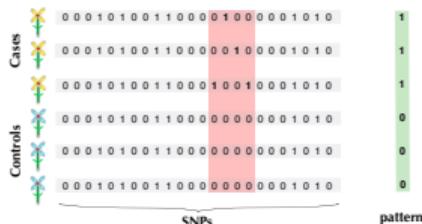
## Erkennung von genetischer Heterogenität

### Fast Automatic Interval Search (Llinares-Lopez et al., ISMB 2015)

- **Aktueller Stand der Forschung:** Beschränkung der Suche auf Intervalle, die Genen oder Exons entsprechen (Lee et al., AJHG 2014).
- Unser Ziel ist es, **nach Intervallen zu suchen, in denen genetische Heterogenität vorkommen könnte**, und dabei
  - beliebige Start- und Endpunkte der Intervalle zu erlauben,
  - die Alphafehler-Kumulierung des inhärenten multiplen Testproblems zu korrigieren, und
  - Teststärke und Recheneffizienz beizubehalten.

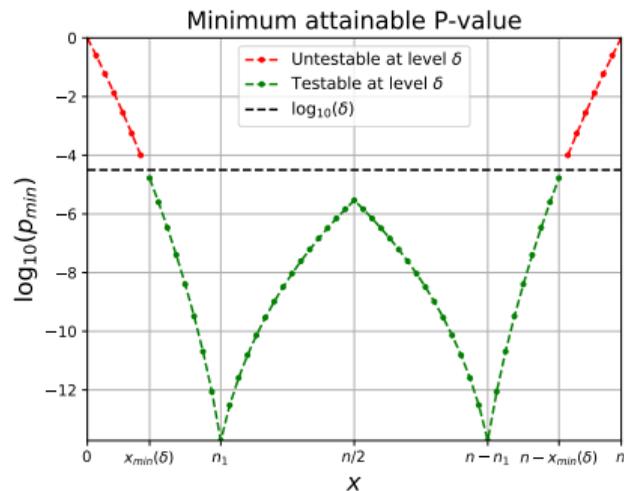
## FAIS: Suche nach Intervallen mit genetischer Heterogenität

Erkennung von genetischer Heterogenität als Fragestellung der Kombinationsuche



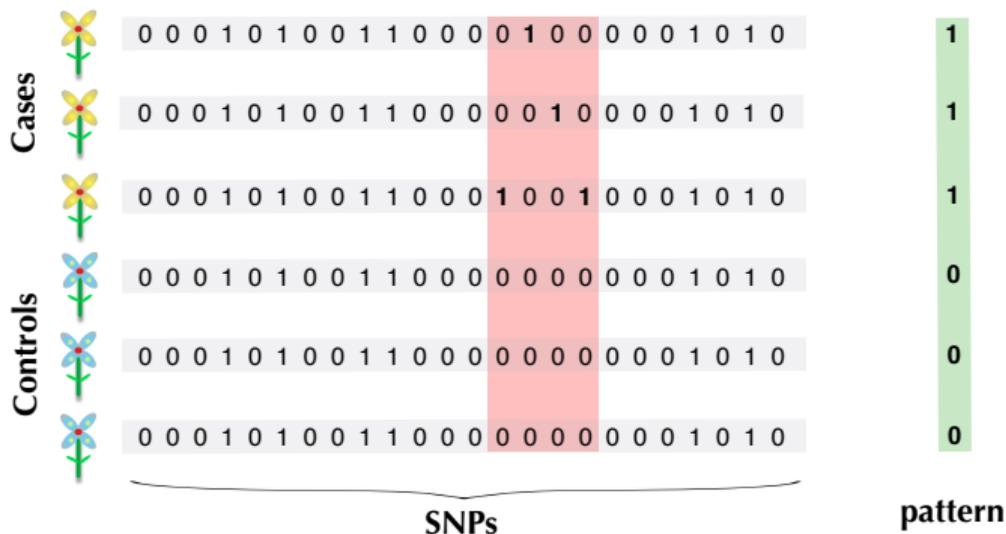
- Wir modellieren die Suche als **Problem der Kombinationsuche**: Gegeben ein Intervall, dann besitzt ein Individuum eine Kombination, wenn es mindestens eine seltene Variante in diesem Intervall hat.
- Ein Intervall wird durch seinen maximalen Wert dargestellt. Je länger das Intervall, desto höher ist die Wahrscheinlichkeit, dass dieses Maximum 1 ist.
- Assoziation wird mit dem exakten Test nach Fisher gemessen, und die Family Wise Error Rate mit Tarones Methode beschränkt.

## FAIS: Suche nach Intervallen mit genetischer Heterogenität



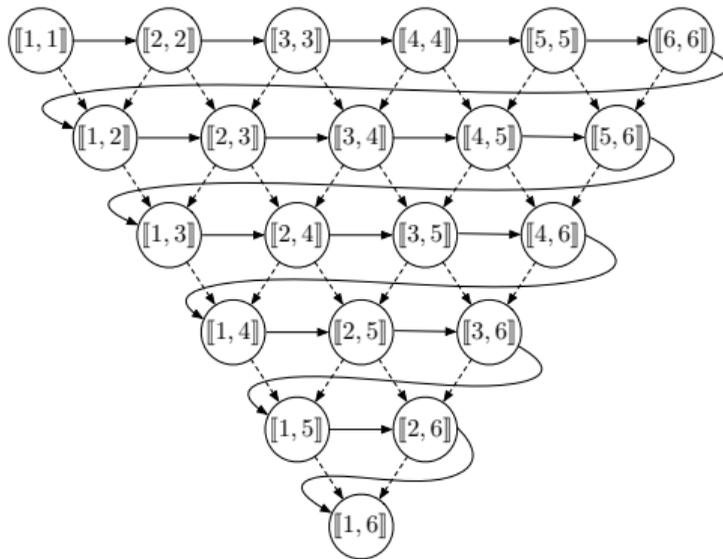
- Falls zu viele Individuen eine bestimmte Kombination besitzen, dann ist das entsprechende Intervall nicht testbar.

## FAIS: Suche nach Intervallen mit genetischer Heterogenität



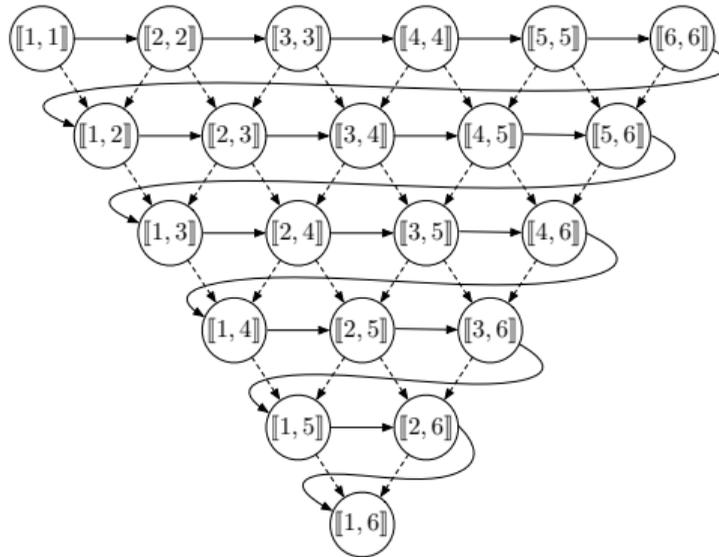
- **Pruning-Kriterium:** Falls eine Kombination zu häufig ist, um testbar zu sein, dann ist kein Superintervall testbar.

# FAIS: Suche nach Intervallen mit genetischer Heterogenität



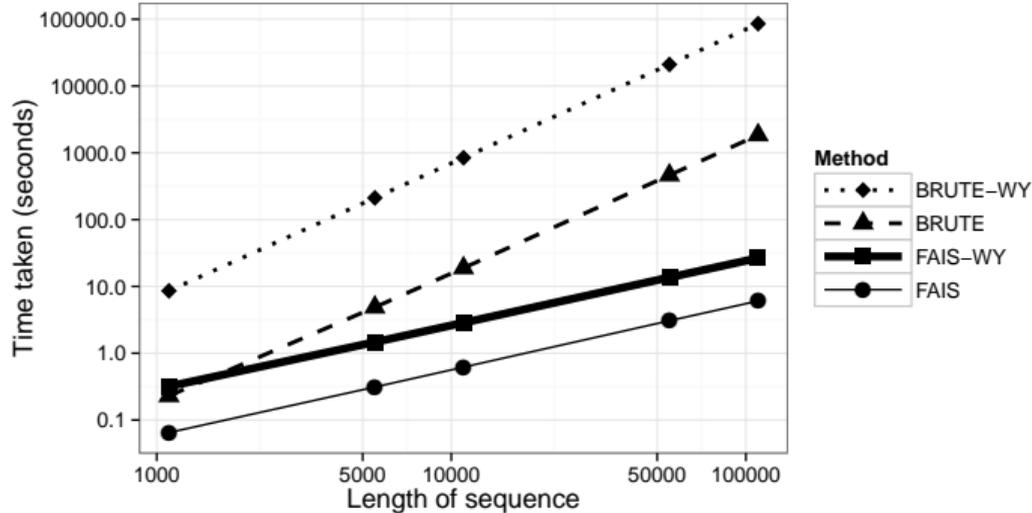
- **Suchstrategie:** Wir suchen Intervalle steigender Länge  $l$  und verwerfen nicht-testbare Superintervalle.

## FAIS: Suche nach Intervallen mit genetischer Heterogenität



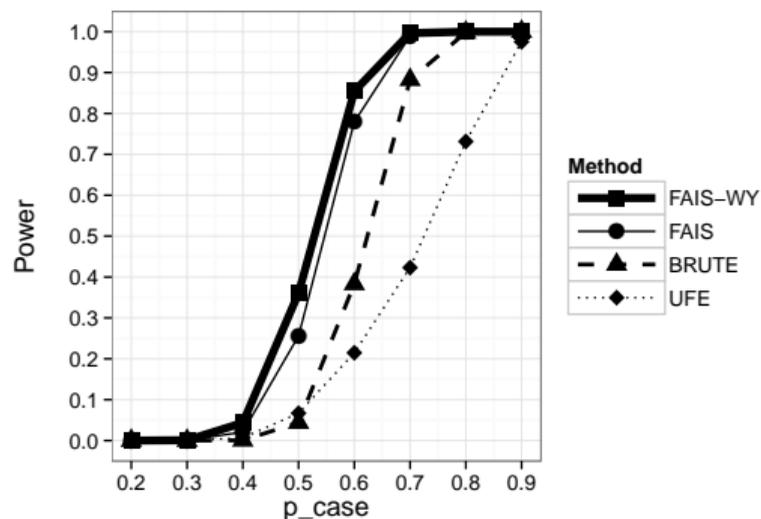
- **Suchstrategie:** Konkret überspringen wir ein Intervall der Länge  $l$ , falls mindestens eines seiner beiden Subintervalle der Länge  $l - 1$  zu häufig ist, um testbar zu sein.

# FAIS: Suche nach Intervallen mit genetischer Heterogenität



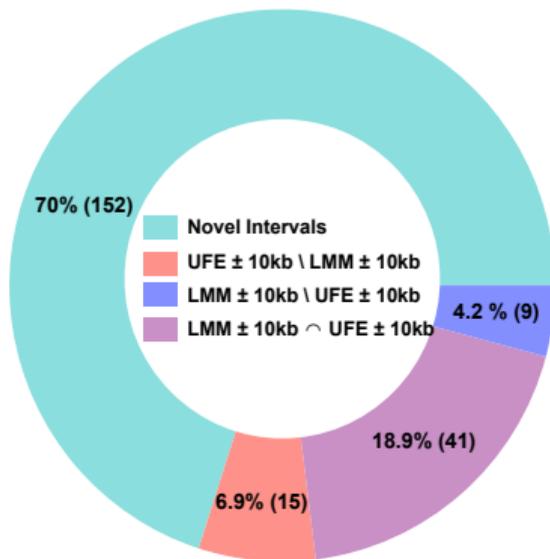
- Unsere Methode FAIS (Fast Automatic Interval Search) erzielt in Simulationen eine bessere Laufzeit als die 'brute-force' Intervallsuche.

# FAIS: Suche nach Intervallen mit genetischer Heterogenität



- Unsere Methode FAIS (Fast Automatic Interval Search) erzielt in Simulationen eine bessere Teststärke als die 'brute-force' Intervallsuche und univariate Methoden.

## FAIS: Suche nach Intervallen mit genetischer Heterogenität



- FAIS entdeckt 217 signifikante Intervalle bei 21 binären Phänotypen von *Arabidopsis thaliana*, mit 214,051 SNPs und bis zu 194 Linien (Atwell et al., Nature 2010).
- 70% wären mit univariaten Methoden (UFE and LMM) übersehen worden.

# Kombinatorische Assoziationsuche: Zusammenfassung und Ausblick

## Zusammenfassung

- Die **Suche nach signifikanten Kombinationen** galt lange als unlösbar.

# Kombinatorische Assoziationsuche: Zusammenfassung und Ausblick

## Zusammenfassung

- Die **Suche nach signifikanten Kombinationen** galt lange als unlösbar.
- Wir haben die kombinatorische Assoziationsuche auf mehreren Ebenen verbessert, wodurch genomweite kombinatorische Assoziationsuche (z. B. für die **Erkennung von genetischer Heterogenität**) ermöglicht wird.

# Kombinatorische Assoziationsuche: Zusammenfassung und Ausblick

## Zusammenfassung

- Die **Suche nach signifikanten Kombinationen** galt lange als unlösbar.
- Wir haben die kombinatorische Assoziationsuche auf mehreren Ebenen verbessert, wodurch genomweite kombinatorische Assoziationsuche (z. B. für die **Erkennung von genetischer Heterogenität**) ermöglicht wird.

`www.significant-patterns.org`

# Kombinatorische Assoziationsuche: Zusammenfassung und Ausblick

## Zusammenfassung

- Die **Suche nach signifikanten Kombinationen** galt lange als unlösbar.
- Wir haben die kombinatorische Assoziationsuche auf mehreren Ebenen verbessert, wodurch genomweite kombinatorische Assoziationsuche (z. B. für die **Erkennung von genetischer Heterogenität**) ermöglicht wird.

[www.significant-patterns.org](http://www.significant-patterns.org)

## Kommende Herausforderungen

- Wie kann genetische Heterogenität in biologischen Netzwerken entdeckt werden?
- Wie kann die False Discovery Rate kontrolliert werden?
- Wie können nichtbinäre Features und nichtbinäre Phänotypen behandelt werden?

## Biomedizinische Softwareentwicklung

## easyGWAS

- Wir haben [easygwas.org](http://easygwas.org) (Grimm et al., 2017), eine Cloud-Plattform für genomweite Assoziationsstudien (1362 Benutzer am 10. April 2018), entwickelt:



## Biomarkerentdeckung für Sepsis

## Personalized Swiss Sepsis Study

- Konsortium aus 22 Forschungsgruppen und 5 Universitätskrankenhäusern in der Schweiz
- Ziel: Vorhersage von Sepsis und der damit zusammenhängenden Sterblichkeit
- Vorgehen: Integration von klinischen und molekularen Daten zur gemeinsamen Biomarkerentdeckung



**Adrian Egli**  
PI SPHN  
Clinical Microbiology, University Hospital Basel



**Karsten Borgwardt**  
PI PHRT  
MLCB, D-BSSE, ETH Zürich



- Dauer:  
3 Jahre  
(2018-2021)
- Gesamtbudget:  
5.3 Millionen CHF

## Vorhersage von Sepsis

### Hintergrund: Was ist Sepsis und warum ist sie relevant?

- Sepsis ist ein lebensbedrohliches Organversagen, verursacht durch eine fehlregulierte Reaktion des Körpers auf eine Infektion (Singer et al., 2016).
- Nach der Blutentnahme kann es noch 24 bis 48h dauern, eine Bakterienspezies im Blut zu identifizieren (Osthoff et al., 2017).
- Nach Ausbruch der Erkrankung steigt das Sterberisiko mit jeder Stunde, um die die Antibiotikabehandlung verzögert wird (Ferrer et al., 2014).

# Vorhersage von Sepsis

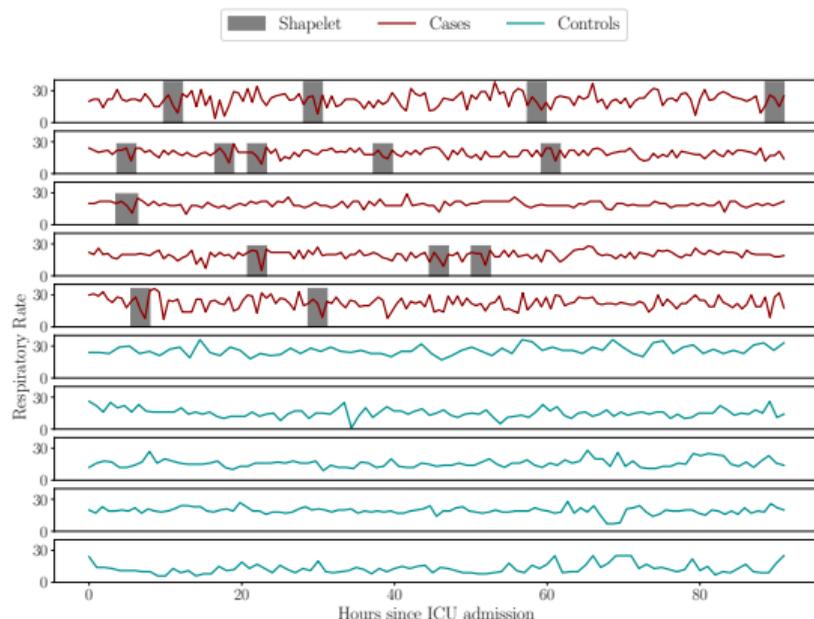
## Hintergrund: Was ist Sepsis und warum ist sie relevant?

- Sepsis ist ein lebensbedrohliches Organversagen, verursacht durch eine fehlregulierte Reaktion des Körpers auf eine Infektion (Singer et al., 2016).
  - Nach der Blutentnahme kann es noch 24 bis 48h dauern, eine Bakterienspezies im Blut zu identifizieren (Osthoff et al., 2017).
  - Nach Ausbruch der Erkrankung steigt das Sterberisiko mit jeder Stunde, um die die Antibiotikabehandlung verzögert wird (Ferrer et al., 2014).
- Die ersten Stunden sind von entscheidender Wichtigkeit.
- Derzeit ist der Organschaden bereits fortgeschritten, wenn Sepsis diagnostiziert wird.
- **Schnellere Erkennung und Behandlung von Sepsis** und bessere Identifikation von Hochrisikogruppen könnte von höchster klinischer Bedeutung sein.

## Vorhersage von Sepsis

- Datensatz: MIMIC (<https://mimic.physionet.org>)
- Labels:
  - Fall: Sepsis-3-Kriterien (Singer et al., 2016) während Aufenthalts auf der Intensivstation (mindestens 4 Stunden nach Aufnahmen) erfüllt, wobei der Begriff 'suspicion of infection' wie in (Seymour et al., 2016) definiert wird.
  - Kontrolle: kein Infektionsverdacht (zumindest nicht während Aufenthalt auf Intensivstation und den 2 vorangegangenen Wochen)
  - Anstieg des SOFA-Scores, gemessen als der maximale SOFA-Score im 3d-Fenster um den Infektionsverdacht (-2 bis +1 Tage) im Vergleich zum Grundwert (baseline SOFA im 3d-Fenster davor).
- Features: Zeitreihen mit Herzfrequenz, systolischem Blutdruck und Atemfrequenz während intensivmedizinischer Behandlung
- Stichprobengröße: 355 intensivmedizinische Behandlungen im Fall, 21,079 Kontrollen (Zufallsstichprobe 355)

# Vorhersage von Sepsis



- Wir entdecken Kombinationen in Zeitreihen mit der Atemfrequenz, die statistisch signifikant mit Sepsis assoziiert sind (Bock et al., ISMB 2018).

## Data-Mining in den Lebenswissenschaften

# Data-Mining in Genetik, Medizin und den Lebenswissenschaften

## Ausblick

- Automatisierung, Biomarkerentdeckung und biomedizinisches Datenmanagement werden Hauptforschungsthemen bleiben.

# Data-Mining in Genetik, Medizin und den Lebenswissenschaften

## Ausblick

- **Automatisierung, Biomarkerentdeckung und biomedizinisches Datenmanagement** werden Hauptforschungsthemen bleiben.
- **Datenwachstum in drei Dimensionen** wird extreme neue Herausforderungen im Data-Mining in Genetik und Medizin generieren:
  - Datensätze von Individuen einer ganzen Population
  - Lebenslange Aufzeichnung des Gesundheitszustandes
  - Höchstaufgelöste Information über den Gesundheitszustand

# Data-Mining in Genetik, Medizin und den Lebenswissenschaften

## Ausblick

- Automatisierung, Biomarkerentdeckung und biomedizinisches Datenmanagement werden Hauptforschungsthemen bleiben.
- Datenwachstum in drei Dimensionen wird extreme neue Herausforderungen im Data-Mining in Genetik und Medizin generieren:
  - Datensätze von Individuen einer ganzen Population
  - Lebenslange Aufzeichnung des Gesundheitszustandes
  - Höchstaufgelöste Information über den Gesundheitszustand
- Wie können diese Daten ausgewertet und genutzt werden?

# Data-Mining in Genetik, Medizin und den Lebenswissenschaften

## Ausblick

- Automatisierung, Biomarkerentdeckung und biomedizinisches Datenmanagement werden Hauptforschungsthemen bleiben.
- Datenwachstum in drei Dimensionen wird extreme neue Herausforderungen im Data-Mining in Genetik und Medizin generieren:
  - Datensätze von Individuen einer ganzen Population
  - Lebenslange Aufzeichnung des Gesundheitszustandes
  - Höchstaufgelöste Information über den Gesundheitszustand
- Wie können diese Daten ausgewertet und genutzt werden?
- Viele Zweige der Lebenswissenschaften stehen vor sehr ähnlichen oder analogen Problemen.

# Data-Mining in Genetik, Medizin und den Lebenswissenschaften

## Ausblick

- **Automatisierung, Biomarkerentdeckung und biomedizinisches Datenmanagement** werden Hauptforschungsthemen bleiben.
- **Datenwachstum in drei Dimensionen** wird extreme neue Herausforderungen im Data-Mining in Genetik und Medizin generieren:
  - Datensätze von Individuen einer ganzen Population
  - Lebenslange Aufzeichnung des Gesundheitszustandes
  - Höchstaufgelöste Information über den Gesundheitszustand
- **Wie können diese Daten ausgewertet und genutzt werden?**
- Viele Zweige der Lebenswissenschaften stehen vor sehr ähnlichen oder analogen Problemen.

**Viele Möglichkeiten für Data-Mining in den Lebenswissenschaften**

## Vielen Dank



- Marie-Curie-Initial Training Network for 'Machine Learning for Personalized Medicine' (mlpm.eu, 2013-2016)
- Starting Grant (Temporäre ERC-Ersatzmaßnahmen des SNSF)
- Alfried-Krupp-Förderpreis für junge Hochschullehrer
- SPHN-PHRT Driver Project 'Personalized Swiss Sepsis Study'

<http://www.bsse.ethz.ch/mlcb>

## Referenzen I

-  C. E. Bonferroni, *Teoria statistica delle classi e calcolo delle probabilita* (1936). Published: (Pubbl. d. R. Ist. Super. di Sci. Econom. e Commerciali di Firenze. 8) Firenze: Libr. Internaz. Seeber. 62 S. (1936).
-  R. Ferrer, *et al.*, *Critical Care Medicine* **42**, 1749 (2014).
-  D. G. Grimm, *et al.*, *The Plant Cell* **29**, 5 (2017).
-  T. F. Mackay, J. H. Moore, *Genome Medicine* **6**, 42 (2014).
-  T. A. Manolio, *et al.*, *Nature* **461**, 747 (2009).
-  S. Lee, *et al.*, *The American Journal of Human Genetics* **95**, 5 (2014).
-  F. Llinares-López, *et al.*, *KDD* (2015), pp. 725–734.
-  F. Llinares-López, *et al.*, *Bioinformatics* **31**, i240 (2015).
-  F. Llinares-López, *et al.*, *Bioinformatics (Oxford, England)* (2017).
-  M. Osthoff, *et al.*, *Clinical Microbiology and Infection: The Official Publication of the European Society of Clinical Microbiology and Infectious Diseases* **23**, 78 (2017).

## Referenzen II

-  L. Papaxanthos, *et al.*, *NIPS*, D. D. Lee, *et al.*, eds. (2016), pp. 2271–2279.
-  C. W. Seymour, *et al.*, *JAMA* **315**, 762 (2016).
-  M. Singer, *et al.*, *JAMA* **315**, 801 (2016).
-  M. Sugiyama, *et al.*, *SIAM Data Mining*, S. Venkatasubramanian, J. Ye, eds. (SIAM, 2015), pp. 37–45.
-  R. E. Tarone, *Biometrics* **46**, 515 (1990).
-  A. Terada, *et al.*, *Proceedings of the National Academy of Sciences* **110**, 12996 (2013).
-  P. H. Westfall, *et al.*, *Biometrics* **49**, 941 (1993).

Icon source: Icons made by Freepik from [www.flaticon.com](http://www.flaticon.com), licensed under CC BY 3.0.