

# Algebraische Statistik – ein junges Forschungsgebiet

Dipl.-Math. Marcus Weber



Disputationsvortrag

15. Februar 2006



1. Statistische Modelle
2. Algebraische Interpretation statistischer Probleme
3. Der Buchberger-Algorithmus
4. Diskussion bisheriger Resultate

---

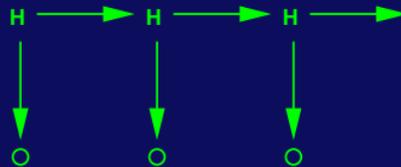
# 1. Statistische Modelle

# Statistisches Modell

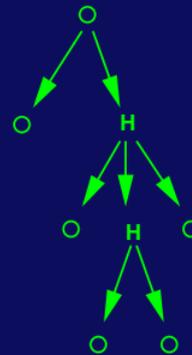
Modellierung bedingter Wahrscheinlichkeiten versteckter (H) und beobachtbarer (O) Ereignisse



**Markov Chain Model**



**Hidden Markov Model**

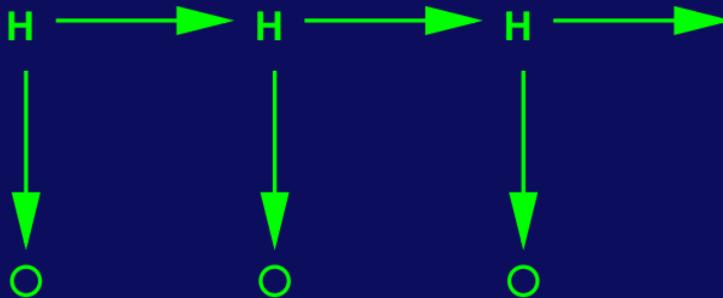


**Tree Model**

# Beobachtung

---

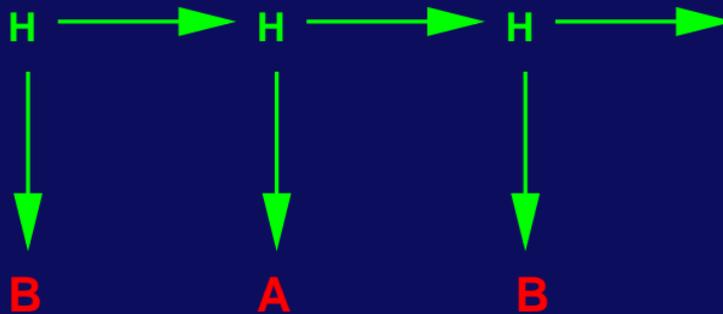
Beispiel: Hidden Markov Model mit  $O \in \{A, B, C\}$  und  $H \in \{a, b\}$ .



# Beobachtung

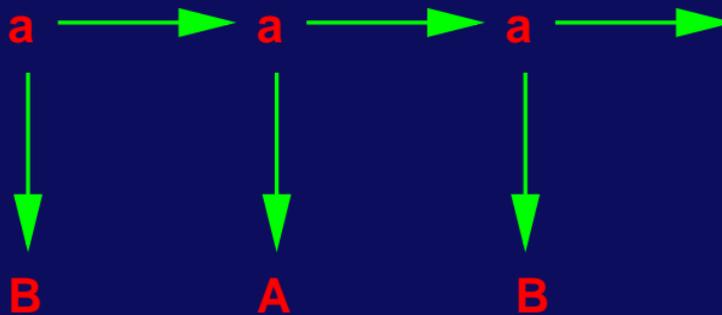
---

Beispiel: Hidden Markov Model mit  $O \in \{A, B, C\}$  und  $H \in \{a, b\}$ .



# Beobachtung

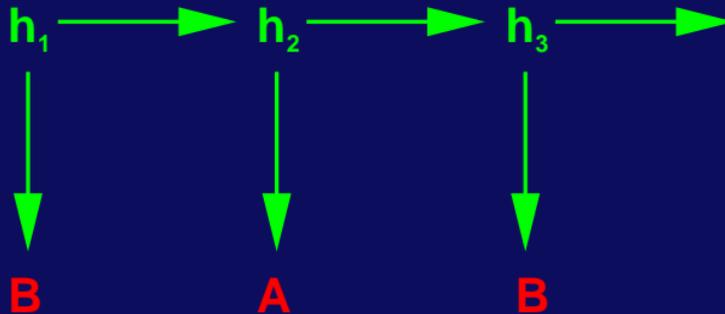
Beispiel: Hidden Markov Model mit  $O \in \{A, B, C\}$  und  $H \in \{a, b\}$ .



$$w_{a \rightarrow B} \cdot w_{a \rightarrow a} \cdot w_{a \rightarrow A} \cdot w_{a \rightarrow a} \cdot w_{a \rightarrow B}$$

# Beobachtung

Beispiel: Hidden Markov Model mit  $O \in \{A, B, C\}$  und  $H \in \{a, b\}$ .



$$w_{BAB} = \sum_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

# Algebraisches Statistisches Modell

$$w_{BAB} = \sum_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$w_{H \rightarrow H}$	a	b	$w_{H \rightarrow O}$	A	B	C
a	$\theta_1$	$\theta_2$	a	$\theta_5$	$\theta_6$	$\theta_7$
b	$\theta_3$	$\theta_4$	b	$\theta_8$	$\theta_9$	$\theta_{10}$

$$\mathbf{p} : \Delta_{2 \times 2, 2 \times 3} \subset \mathbb{R}^{10} \rightarrow \Delta_{27} \subset \mathbb{R}^{27}$$

# Algebraisches Statistisches Modell

$$w_{BAB} = \sum_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$w_{H \rightarrow H}$	a	b	$w_{H \rightarrow O}$	A	B	C
a	$\theta_1$	$\theta_2$	a	$\theta_5$	$\theta_6$	$\theta_7$
b	$\theta_3$	$\theta_4$	b	$\theta_8$	$\theta_9$	$\theta_{10}$

$$\mathbf{p} : \mathbb{C}^{10} \rightarrow \mathbb{C}^{27}$$

$\mathbf{p}_i(\Theta)$  ist ein Polynom in  $\Theta = [\theta_1, \dots, \theta_{10}]$ .

---

## 2. Algebraische Interpretation statistischer Probleme

# Wahrscheinlichste Pfade

---

$$w_{BAB} = \sum_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

# Wahrscheinlichste Pfade

---

$$w_{BAB} = \sum_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$$w_{BAB}^* = \max_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

# Wahrscheinlichste Pfade

---

$$w_{BAB} = \sum_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$$w_{BAB}^* = \max_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$$l_{BAB}^* = \max_{h_1, h_2, h_3 \in \{a, b\}} l_{h_1 \rightarrow B} + l_{h_1 \rightarrow h_2} + l_{h_2 \rightarrow A} + l_{h_2 \rightarrow h_3} + l_{h_3 \rightarrow B}$$

# Wahrscheinlichste Pfade

$$w_{BAB} = \sum_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$$w_{BAB}^* = \max_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$$l_{BAB}^* = \max_{h_1, h_2, h_3 \in \{a, b\}} l_{h_1 \rightarrow B} + l_{h_1 \rightarrow h_2} + l_{h_2 \rightarrow A} + l_{h_2 \rightarrow h_3} + l_{h_3 \rightarrow B}$$

$$l_{BAB}^* = \bigoplus_{h_1, h_2, h_3 \in \{a, b\}} l_{h_1 \rightarrow B} \odot l_{h_1 \rightarrow h_2} \odot l_{h_2 \rightarrow A} \odot l_{h_2 \rightarrow h_3} \odot l_{h_3 \rightarrow B}$$

# Wahrscheinlichste Pfade

$$w_{BAB} = \sum_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$$w_{BAB}^* = \max_{h_1, h_2, h_3 \in \{a, b\}} w_{h_1 \rightarrow B} \cdot w_{h_1 \rightarrow h_2} \cdot w_{h_2 \rightarrow A} \cdot w_{h_2 \rightarrow h_3} \cdot w_{h_3 \rightarrow B}$$

$$l_{BAB}^* = \max_{h_1, h_2, h_3 \in \{a, b\}} l_{h_1 \rightarrow B} + l_{h_1 \rightarrow h_2} + l_{h_2 \rightarrow A} + l_{h_2 \rightarrow h_3} + l_{h_3 \rightarrow B}$$

$$l_{BAB}^* = \bigoplus_{h_1, h_2, h_3 \in \{a, b\}} l_{h_1 \rightarrow B} \odot l_{h_1 \rightarrow h_2} \odot l_{h_2 \rightarrow A} \odot l_{h_2 \rightarrow h_3} \odot l_{h_3 \rightarrow B}$$

Geschicktes Ausrechnen (Distributiv-Gesetz, Horner-Schema) dieses "Polynoms":  
 Viterbi-Algorithmus (Hidden Markov Models), Needleman-Wunsch-Algorithmus (Sequenz-Alignment-Graph)

# Modell-Invarianten



Zusammenhänge zwischen  $w_{BAB}, w_{BAC}, \dots$ , die unabhängig von der Wahl der Parameter  $\Theta$  sind, sondern nur das Statistische Modell betreffen? Z.B.:

$$w_{AAA} - w_{BAC} = w_{CCC} - w_{CAB}$$

$$\mathbf{p} : \mathbb{C}^{10} \rightarrow \mathbb{C}^{27}$$

$$w_{AAA} = \mathbf{p}_1(\Theta)$$

$$w_{AAB} = \mathbf{p}_2(\Theta)$$

$$\vdots$$

$$w_{CCC} = \mathbf{p}_{27}(\Theta)$$

# Modell-Invarianten

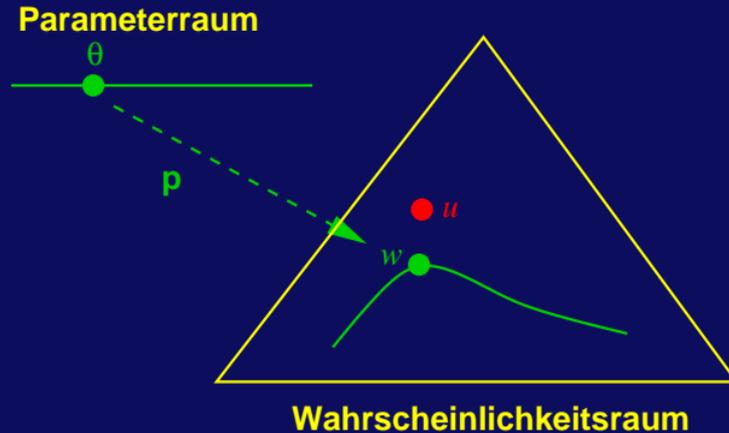
---

$$\mathbf{p} : \mathbb{C}^{10} \rightarrow \mathbb{C}^{27}$$

$$\begin{aligned} 0 &= w_{AAA} - \mathbf{p}_1(\Theta) \\ 0 &= w_{AAB} - \mathbf{p}_2(\Theta) \\ &\vdots \\ 0 &= w_{CCC} - \mathbf{p}_{27}(\Theta) \end{aligned}$$

27 algebraische Gleichungen mit  $37=27+10$  "Unbekannten". Varietät  $\mathcal{V}$ .  
Verfahren: Elimination der 10 Parameter ("Eliminationstheorie").

# Parameterschätzung, Log-Likelihood



# Parameterschätzung, Log-Likelihood

Maximiere bezüglich  $\Theta$  die Log-Likelihood-Funktion:

$$L_u(\Theta) = \sum_{i=1}^{27} u_i \log(\mathbf{p}_i(\Theta))$$

Algebraischer Ansatz: Berechne *alle* kritischen Punkte von  $L_u$ , auch alle komplexwertigen.

$$0 = \frac{\partial L_u(\Theta)}{\partial \theta_j} = \sum_{i=1}^{27} \frac{u_i}{\mathbf{p}_i(\Theta)} \frac{\partial \mathbf{p}_i(\Theta)}{\partial \theta_j}, \quad j = 1, \dots, 10.$$

## Parameterschätzung, Log-Likelihood

---

Maximiere bezüglich  $\Theta$  die Log-Likelihood-Funktion:

$$L_u(\Theta) = \sum_{i=1}^{27} u_i \log(\mathbf{p}_i(\Theta))$$

Algebraischer Ansatz: Berechne *alle* kritischen Punkte von  $L_u$ , auch alle komplexwertigen.

$$0 = \frac{\partial L_u(\Theta)}{\partial \theta_j} = \sum_{i=1}^{27} u_i z_i \frac{\partial \mathbf{p}_i(\Theta)}{\partial \theta_j}, \quad j = 1, \dots, 10,$$
$$0 = z_i \mathbf{p}_i(\Theta) - 1, \quad i = 1, \dots, 27.$$

---

### 3. Der Buchberger-Algorithmus

# Varietäten

---

$$\begin{aligned}0 &= p_1(x_1, x_2, \dots, x_m) \\0 &= p_2(x_1, x_2, \dots, x_m) \\&\vdots \\0 &= p_n(x_1, x_2, \dots, x_m)\end{aligned}$$

$$\mathcal{V}(p_1, \dots, p_n) \subseteq \mathbf{C}^m$$

## Idee der Gauß-Elimination:

Umformung eines *linearen* Gleichungssystems in ein äquivalentes *lineares* System “einfacherer” Gestalt.

$$0 = 2x_1 - 5x_2 + 3x_3 - 4$$

$$0 = x_1 - \frac{5}{2}x_2 + \frac{3}{2}x_3 - 2$$

$$0 = 4x_1 + x_2 + x_3 + 1 \Leftrightarrow$$

$$0 = x_2 - \frac{8}{11}x_3 + \frac{9}{11}$$

$$0 = x_1 - x_2 - x_3 - 6$$

$$0 = x_3 + \frac{23}{8}$$

# Idee des Buchberger-Algorithmus

---

## Idee des Buchberger-Algorithmus (Buchberger, 1965):

Umformung eines *algebraischen* Gleichungssystems in ein äquivalentes *algebraisches* System “einfacherer” Gestalt.

# Idee des Buchberger-Algorithmus

---

Idee des Buchberger-Algorithmus (Buchberger, 1965):

**Umformung** eines *algebraischen* Gleichungssystems in ein äquivalentes *algebraisches* System “einfacherer” Gestalt.

$$\begin{array}{l} 0 = p(x), 0 = q(x) \\ \Rightarrow 0 = p(x) \cdot r(x), \quad \forall r \in \mathbf{K}[X_1, X_2, \dots] \\ \Rightarrow 0 = p(x) + q(x). \end{array}$$

# Idee des Buchberger-Algorithmus

Idee des Buchberger-Algorithmus (Buchberger, 1965):

**Umformung** eines *algebraischen* Gleichungssystems in ein äquivalentes *algebraisches* System “einfacherer” Gestalt.

$$\begin{aligned} 0 = p(x), 0 = q(x) &\Rightarrow 0 = p(x) \cdot r(x), \quad \forall r \in \mathbf{K}[X_1, X_2, \dots] \\ &\Rightarrow 0 = p(x) + q(x). \end{aligned}$$

Varietät eines **Ideals**:

$$\begin{aligned} \mathcal{V}(\mathcal{F}) &= \mathcal{V}(\langle \mathcal{F} \rangle) = \mathcal{V}(\langle \mathcal{G} \rangle) = \mathcal{V}(\mathcal{G}) \\ \mathcal{V}(\langle 2x^3 - 5x^2 - 8x + 20, x^3 + x^2 - 4x - 4 \rangle) &= \mathcal{V}(\langle x^2 - 4 \rangle) = \{2, -2\} \end{aligned}$$

# Euklidische Ringe

---

Beispiel: Ideal einer Familie  $\mathcal{F}$  ganzer Zahlen.

$$\langle 198, 108 \rangle = \langle \text{ggT}(198, 108) \rangle = \langle 18 \rangle.$$

**Euklidischer Algorithmus:** Größter gemeinsamer Teiler kann durch sukzessive Division mit Rest (wobei der "Grad" des Restgliedes abnimmt) bestimmt werden.

$$\begin{aligned} 198 &= 108 \cdot 1 + 90 \\ 108 &= 90 \cdot 1 + 18 \\ 90 &= 18 \cdot 5 + 0 \end{aligned}$$

**Hauptidealringe:** *Euklidische Ringe sind Hauptidealringe.* Beispiel: Polynomring in einer Variablen mit "Grad"=höchster vorkommender Exponent des Polynoms.

$$\mathcal{V}(\langle 2x^3 - 5x^2 - 8x + 20, x^3 + x^2 - 4x - 4 \rangle) = \mathcal{V}(\langle x^2 - 4 \rangle) = \{2, -2\}$$

# Polynomringe mehrerer Variablen

---

Polynomringe in *mehreren* Variablen sind *keine* euklidischen Ringe. Sprich: Nicht jedes Ideal wird von einem einzigen Polynom erzeugt.

**Verallgemeinerte Division mit Rest:** Sei  $\mathcal{F} = \{p_1, p_2, \dots, p_n\} \subset \mathbf{K}[X_1, X_2, \dots]$  gegeben, dann existiert zu jedem  $p \in \mathbf{K}[X_1, X_2, \dots]$  eine Zerlegung

$$p = a_1p_1 + a_2p_2 + \dots + a_np_n + r$$

mit Polynomen  $a_1, \dots, a_n, r \in \mathbf{K}[X_1, X_2, \dots]$ .

Und:  $a_i p_i \neq 0 \Rightarrow \text{mulgrad}(p) \geq \text{mulgrad}(a_i p_i)$ .

# Gröbner-Basen

---

$$p = xy^2 - x, \quad \mathcal{F} = \{xy + 1, y^2 - 1\}$$

$$xy^2 - x = y(xy + 1) + 0(y^2 - 1) + (-x - y)$$

$$xy^2 - x = x(y^2 - 1) + 0(xy + 1) + 0$$

Gibt es eine Basis  $\mathcal{G}$  des Ideals  $\langle \mathcal{F} \rangle$ , so dass der Rest bei verallgemeinerter Division eindeutig bestimmt ist?

Antwort ja: "Gröbner-Basen". Buchberger-Algorithmus (verallgemeinert Gauß-Elimination und Euklidischen Algorithmus)

0.  $\mathcal{F}$  Erzeugendensystem eines Ideals.  $\mathcal{G} := \mathcal{F}$ .
1. Erweiterung auf eine Gröbner-Basis:
  - a) Konstruiere Testpolynome aus  $\langle \mathcal{F} \rangle$  (siehe Gauß-Eliminationsschritt) und führe unter  $\mathcal{G}$  die verallgemeinerte Division durch.
  - b) Ist der Rest jeweils 0, dann ist die Gröbner-Basis ausreichend, gehe zu 2. Andernfalls füge die Nicht-Null-Reste zu  $\mathcal{G}$  hinzu und wiederhole 1.
2. Normiere die Polynome aus  $\mathcal{G}$  und reduziere (auf bestimmte Weise) die Gröbner-Basis.

---

## 4. Diskussion bisheriger Resultate

- **Colin Dewey, Peter Huggins, Kevin Woods, Bernd Sturmfels, Li-or Pratcher:** Sensitivitätsanalyse des parametrisierten Sequenz-Alignments durch Polytop-Algebra.
- **Colin Dewey, Kevin Woods:** Biologisch korrektes Alignment, das für keine Parameterwahl des statistischen Modells (Sequenz-Alignment-Graph) optimal ist.
- **Marta Casanellas, Luis David Garcia, Seth Sullivant:** Katalog von Eigenschaften "kleiner" Bäume im Web, Auswahl des richtigen statistischen Modells für phylogenetische Bäume anhand dieser Daten möglich.
- **Sergi Elizade:** "Few inference functions theorem".

Statistisches Problem	Algebraische Formulierung	Lösung
wahrsch. Pfade Modell-Invarianten Max-Likelihood	$(\max, +)$ -Algebra Eliminationstheorie Varietäten	geschicktes Ausklammern Buchberger-Algorithmus Buchberger-Algorithmus

## Verwendete Literatur

---

**Lior Patcher, Bernd Sturmfels (eds.):** *Algebraic Statistics for Computational Biology*, Cambridge University Press, Oktober 2005.

**Hans J. Stetter:** *Numerical Polynomial Algebra*, Society for Industrial and Applied Mathematics, Philadelphia, 2004.

**David Cox, John Little, and Donal O'Shea:** *Ideals, Varieties and Algorithms. An Introduction to Computational Algebraic Geometry And Commutative Algebra*, 2nd ed., Undergraduated Texts in Mathematics, Springer-Verlag, New York, 1997.