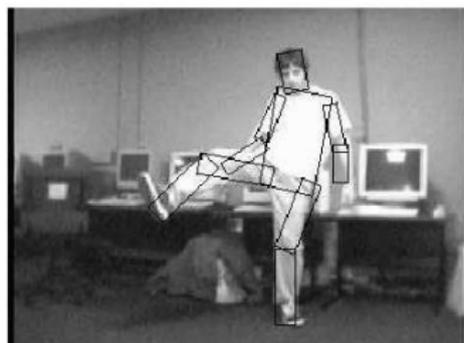
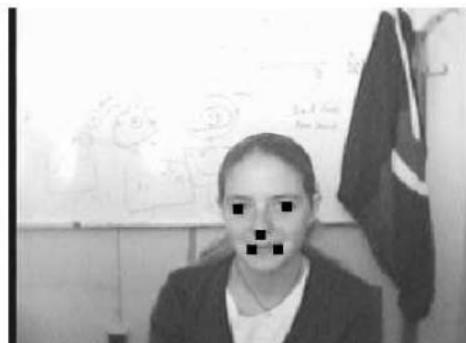


Computer Vision: Pictorial Structures

André Steinborn - TUD/INF/KI/IS/AS

17. November 2011

- Ziel: Lokalisierung von Objekten in Bildern.
- Objektmodell: Menge von Teilen, wobei einige Teile miteinander verbunden sind.



(Felzenszwalb, Huttenlocher 2005)

Literatur: Petro F. Felzenszwalb, Daniel P. Huttenlocher 2005:
Pictorial Structures for Object Recognition.

- Für jeden Anwendungsfall muss definiert werden:
 - Raum der Lokalisationen.
 - Erscheinungsmodell für jedes Teil.
 - Art der Verbindungen zwischen den Teilen.
- Drei Beispiele:
 - Iconic Model.
 - Articulated Models.
 - Boundary Models.

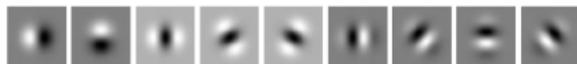
Pictorial Structures

Beispiel: Iconic Model

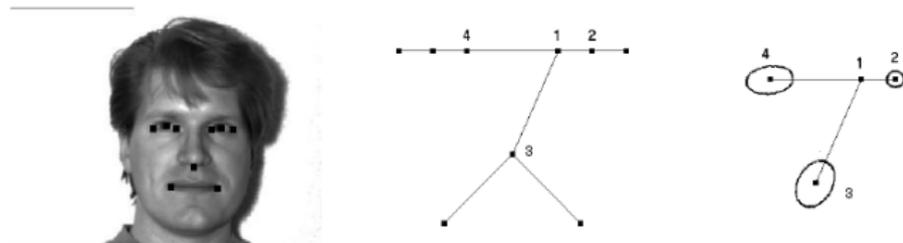
Modell für Gesichtserkennung mit 5 Teilen:



Beobachtungsmodell als Vektor der Filterantworten der Filter:



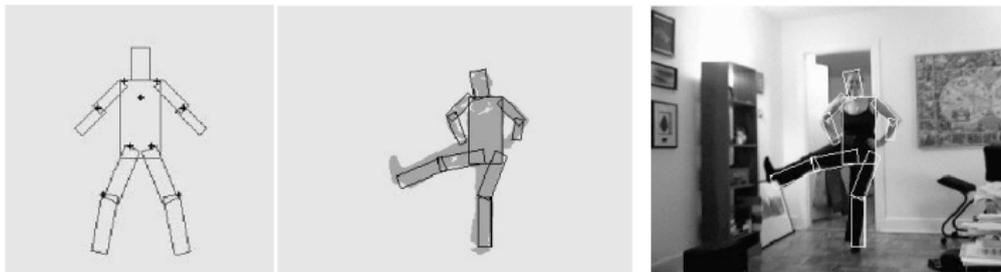
Modell für Gesichtserkennung mit 9 Teilen:



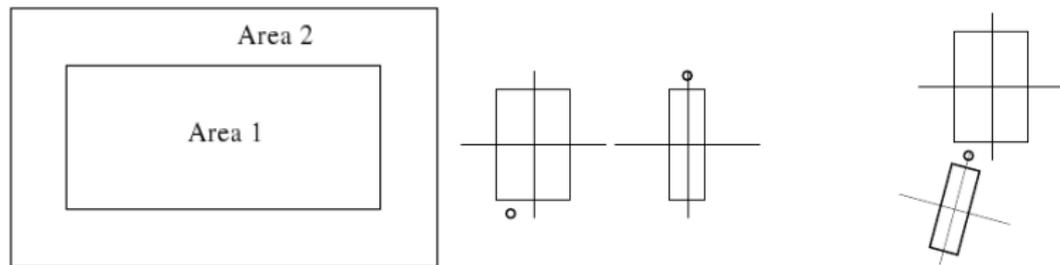
Pictorial Structures

Beispiel: Articulated Models

Modell für verschiedene Posen von Menschen:



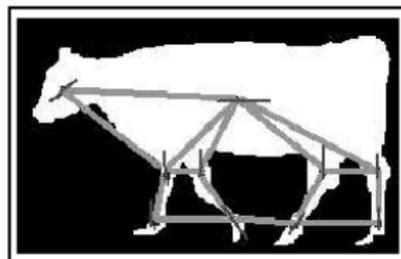
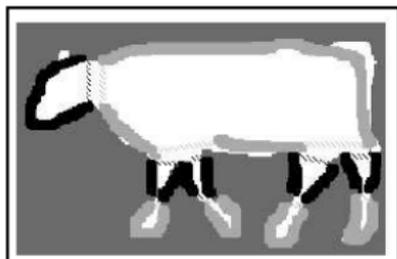
Beobachtungsmodell und Strukturmodell:



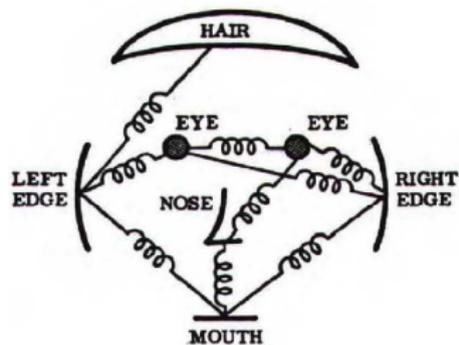
Pictorial Structures

Beispiel: Boundary Models

- Beobachtungsmodell:
 - Form der Umrandung.
 - Textur.
- Strukturmodell:
 - Vollverbundener Graph.
 - Potts-Modell.



Literatur: Kumar, Torr, Zisserman 2004: Extending Pictorial Structures for Object Recognition.



(Fischler, Eischlager 1973)

- Konfiguration $L = (l_1, \dots, l_n)$ ist Positionslabeling, wobei l_i die Position des Teils v_i beschreibt.
- Die Bildbewertungsfunktion $m_i(l_i)$ bewertet die Position l_i des Teils v_i als Grad der Abweichung von einem Idealwert unter Berücksichtigung der üblichen Streuung.
- Die Paarbewertungsfunktion $d_{ij}(l_i, l_j)$ bewertet die relative Lage der Teile v_1 und v_2 , jeweils an den Positionen l_1 und l_2 .

- Sei I ein Bild, L eine Konfiguration der Teile und θ ein Satz von Modellparametern.
- Zu betrachtende Verteilungen:
 - Likelihood: $p(I|L; \theta)$.
 - A-priori Verteilung: $p(L; \theta)$.
 - A-posteriori-Verteilung:

$$p(L|I; \theta) \propto p(I|L; \theta) \cdot p(L; \theta)$$

- Zu betrachtende Probleme:
 - Maximum-A-Posteriori (MAP) -Schätzung (Lokalisierung).
 - Lernen der Modellparameter.

- Modellparameter: $\theta = (u, E, c)$, wobei
 - $u = \{u_1, \dots, u_n\}$ sind Erscheinungsparameter für jedes der n Teile.
 - E ist die Menge der Kanten und beschreibt, welche Teile miteinander verbunden sind.
 - $c = \{c_{ij} \mid (v_i, v_j) \in E\}$ sind die Parameter für die Verbindungen zwischen den Teilen.

Die statistische Modellierung des Bildentstehungsprozesses erfolgt durch die Likelihoodfunktion:

$$p(I|L, \theta) = p(I|L; u) \\ \propto \prod_{i=1}^n p(I|l_i; u_i)$$

wobei die Zahlen $p(I|l_i; u_i)$ die Positionen der Teile auf der Grundlage des Bildes bewerten.
Solange sich die Teile nicht überlappen ist das Produkt über diesen Zahlen eine gute Approximation der gewünschten Wahrscheinlichkeitsverteilung.

Modellierung des Vorwissens durch die A-priori-Verteilung der Konfigurationen.
Für azyklische Graphen $G = (V, E)$ ergibt sich:

$$\begin{aligned} p(L; \theta) &= p(L; E, c) \\ &= \frac{\prod_{(v_i, v_j) \in E} p(l_i, l_j; c_{ij})}{\prod_{v_i \in V} p(l_i; c)^{\deg(v_i) - 1}} \\ &\propto \prod_{(v_i, v_j) \in E} p(l_i, l_j; c_{ij}) \end{aligned}$$

Die Proportionalitätsrelation in der letzten Zeile gilt, weil kein Vorwissen über die absoluten Positionen l_i modelliert wird und somit $p(l_i; c)$ gleichverteilt sein muss, weswegen der Nenner in der mittleren Zeile konstant ist.

Beachte: Sowohl $p(l_i, l_j; c_{ij})$ als auch $p(L; E, c)$ sind uneigentliche A-priori-Verteilungen, ergeben aber in Verbindung mit der Likelihoodfunktion eine eigentliche A-posteriori-Verteilung.

Die Wahrscheinlichkeiten der Konfigurationen L nachdem ein Bild I beobachtet wurde wird durch die A-posteriori-Verteilung beschrieben:

$$\begin{aligned} p(L|I; \theta) &= \frac{p(I|L; \theta) \cdot p(L; \theta)}{\sum_{L'} p(I|L'; \theta) \cdot p(L'; \theta)} \\ &= \frac{p(I|L; \theta) \cdot p(L; \theta)}{p(I; \theta)} \\ &\propto p(I|L; \theta) \cdot p(L; \theta) \\ &= \prod_{i=1}^n p(I|l_i; u_i) \cdot \prod_{(v_i, v_j) \in E} p(l_i, l_j; c_{ij}) \end{aligned}$$

- Sei I ein Bild, L eine Konfiguration der Teile und θ ein Satz von Modellparametern.
- Zu betrachtende Verteilungen:
 - Likelihood: $p(I|L; \theta)$.
 - A-priori Verteilung: $p(L; \theta)$.
 - A-posteriori-Verteilung:

$$p(L|I; \theta) \propto p(I|L; \theta) \cdot p(L; \theta)$$

- Zu betrachtende Probleme:
 - Maximum-A-Posteriori (MAP) -Schätzung (Lokalisierung).
 - Lernen der Modellparameter.

Maximum-a-posteriori (MAP) Schätzung:

$$\begin{aligned} L^* &= \arg \max_L \left[\prod_{i=1}^n p(I|l_i; u_i) \cdot \prod_{(v_i, v_j) \in E} p(l_i, l_j; c_{ij}) \right] \\ &= \arg \min_L - \ln \left[\prod_{i=1}^n p(I|l_i; u_i) \cdot \prod_{(v_i, v_j) \in E} p(l_i, l_j; c_{ij}) \right] \\ &= \arg \min_L \left[\sum_{i=1}^n \underbrace{-\ln p(I|l_i; u_i)}_{m_i(l_i)} + \sum_{(v_i, v_j) \in E} \underbrace{-\ln p(l_i, l_j; c_{ij})}_{d_{ij}(l_i, l_j) - \eta_{ij}} \right] \\ &= \arg \min_L \left[\sum_{i=1}^n m_i(l_i) + \sum_{(v_i, v_j) \in E} d_{ij}(l_i, l_j) \right] \end{aligned}$$

Die resultierende Aufgabe nennt man (\min, \sum) -Aufgabe oder Energieminimierungsaufgabe.

Unter Ausnutzung der Baumstruktur kann die Dynamische Programmierung angewandt werden:

$$l_j^* = \arg \min_{l_j} \begin{cases} m_j(l_j) + \sum_{v_c \in C_j} B_c(l_j) & \text{Wurzel} \\ m_j(l_j) + d_{ij}(l_i, l_j) + \sum_{v_c \in C_j} B_c(l_j) & \text{Knoten} \\ m_j(l_j) + d_{ij}(l_i, l_j) & \text{Blatt} \end{cases}$$

wobei

$$B_j(l_i) = \begin{cases} \min_{l_j} (m_j(l_j) + d_{ij}(l_i, l_j)) & \text{Blatt} \\ \min_{l_j} (m_j(l_j) + d_{ij}(l_i, l_j) + \sum_{v_c \in C_j} B_c(l_j)) & \text{Knoten} \end{cases}$$

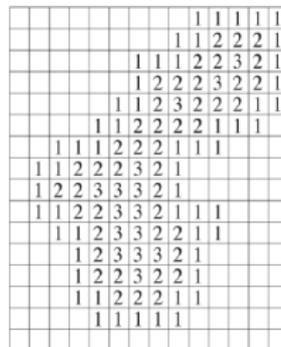
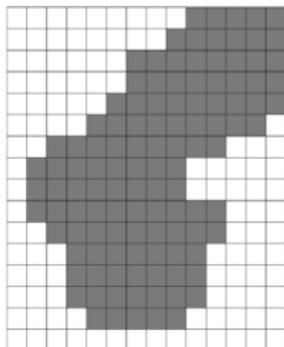
Der Viterbi-Algorithmus läuft in $O(h^2 \cdot n)$, wobei h die Anzahl der möglichen Positionen der Teile ist.

Distance Transformation I

Klassische Distance Transformation für binäre Bilder $f: \mathcal{G} \rightarrow \{0, \infty\}$ über einem Gitter \mathcal{G} für eine Punktmenge $\{p | f(p) = 0, p \in \mathcal{G}\} = P \subseteq \mathcal{G}$ (weiße Pixel):

$$\begin{aligned} \mathcal{D}_P(p) &= \min_{q \in P} d(p, q) \\ &= \min_{q \in \mathcal{G}} (d(p, q) + \mathbf{1}(q)) \end{aligned}$$

wobei $\mathbf{1}(q) = q \in P ? 0 : \infty$.



(Klette, Rosenfeld 2004)

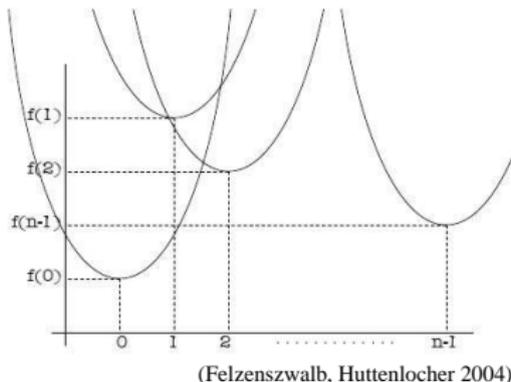
Literatur: Reinhard Klette, Azriel Rosenfeld 2004:
Digital Geometry - Geometric Methods for Digital Picture Analysis.

Verallgemeinerte Distance Transformation

Verallgemeinerung der Distance Transformation für beliebige Funktionen $f: \mathcal{G} \rightarrow \mathbb{R}$ über einem Gitter \mathcal{G} :

$$\mathcal{D}_f(p) = \min_{q \in \mathcal{G}} (d(p, q) + f(q))$$

wobei man für binäre Funktionen $f: \mathcal{G} \rightarrow \{0, \infty\}$ und $f(q) = \mathbf{1}(q)$ die klassische Distance Transformation erhält.



Die Zeitkomplexität ist $O(h)$, wenn h die Anzahl der Gitterpunkte in \mathcal{G} ist.

Literatur: Petro F. Felzenszwalb, Daniel P. Huttenlocher 2004:
Distance Transforms of Sampled Functions.

Verallgemeinerte Distance Transformation:

$$\mathcal{D}_f(p) = \min_{q \in \mathcal{G}} (d(p, q) + f(q))$$

Für die MAP-Schätzung braucht man:

$$\begin{aligned} B_j(l_i) &= \min_{l_j} \left(\underbrace{d_{ij}(l_i, l_j)}_{d(l_i, l_j)} + \underbrace{m_j(l_j) + \sum_{v_c \in C_j} B_c(l_j)}_{f(l_j)} \right) \\ &= \min_{l_j} (d(l_i, l_j) + f(l_j)) \end{aligned}$$

Der Algorithmus läuft in $O(h' \cdot n)$, wobei h' die Anzahl der möglichen Positionen im transformierten Raum ist.

- Wegen Baumstruktur:
 - MAP-Schätzung mit dynamischer Programmierung möglich.
 - Zeitkomplexität: $O(h^2 \cdot n)$.
- Wegen Mahalanobisabstand:
 - Berechnen der Zahlen $B_j(l_i)$ als Distance Transformation (Zeitkomplexität: $O(h')$) möglich.
 - Zeitkomplexität für MAP-Schätzung: $O(h' \cdot n)$.

- Gegeben:
 - Eine Menge von repräsentativen Beispielbildern: $\{I^1, \dots, I^m\}$
 - sowie korrespondierende Konfigurationen: $\{L^1, \dots, L^m\}$.
- Gesucht: Modellparameter $\theta = (u, E, c)$.
- Maximum-Likelihood (ML) - Schätzung:

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p(I^1, \dots, I^m, L^1, \dots, L^m; \theta) \\ &= \arg \max_{\theta} \prod_{k=1}^m p(I^k, L^k; \theta) \\ &= \arg \max_{\theta} \prod_{k=1}^m p(I^k | L^k; \theta) \cdot \prod_{k=1}^m p(L^k; \theta) \\ &= \arg \max_{\theta} \left[\sum_{k=1}^m \ln p(I^k | L^k; u) + \sum_{k=1}^m \ln p(L^k; E, c) \right]\end{aligned}$$

Wegen der Unabhängigkeit der beiden Summanden:

$$\begin{aligned}u^* &= \arg \max_u \sum_{k=1}^m \ln p(I^k | L^k; u) \\ (E, c)^* &= \arg \max_{E, c} \sum_{k=1}^m \ln p(L^k; E, c)\end{aligned}$$

ML-Schätzung für die Erscheinungsparameter:

$$\begin{aligned}u^* &= \arg \max_u \sum_{k=1}^m \ln p \left(I^k | L^k; u \right) \\&= \arg \max_u \sum_{k=1}^m \ln \prod_{i=1}^n p \left(I^k | l_i^k; u_i \right) \\&= \arg \max_u \sum_{k=1}^m \sum_{i=1}^n \ln p \left(I^k | l_i^k; u_i \right) \\&= \arg \max_u \sum_{i=1}^n \sum_{k=1}^m \ln p \left(I^k | l_i^k; u_i \right)\end{aligned}$$

Weswegen die Erscheinungsparameter für jedes Teil unabhängig voneinander gelernt werden können:

$$u_i^* = \arg \max_{u_i} \sum_{k=1}^m \ln p \left(I^k | l_i^k; u_i \right)$$

ML-Schätzung für die Verbindungsparameter c_{ij} :

$$(E, c)^* = \arg \max_{E, c} \sum_{k=1}^m \ln p \left(L^k; E, c \right)$$

Unter Ausnutzung der Baumstruktur erhält man

$$\begin{aligned} (E, c)^* &= \arg \max_{E, c} \sum_{k=1}^m \ln \prod_{(v_i, v_j) \in E} p \left(l_i^k, l_j^k; c_{ij} \right) \\ &= \arg \max_{E, c} \sum_{(v_i, v_j) \in E} \sum_{k=1}^m \ln p \left(l_i^k, l_j^k; c_{ij} \right) \end{aligned}$$

Somit können die Parameter c_{ij} für jede mögliche Verbindung bereits ohne Kenntnis von E gelernt werden:

$$c_{ij}^* = \arg \max_{c_{ij}} \sum_{k=1}^m \ln p \left(l_i^k, l_j^k; c_{ij} \right)$$

ML-Schätzung für die Kantenmenge E :

$$(E, c)^* = \arg \max_{E, c} \sum_{(v_i, v_j) \in E} \sum_{k=1}^m \ln p(l_i^k, l_j^k; c_{ij})$$

Für bekannte c_{ij}^* erhält man für die Kantenmenge

$$\begin{aligned} E^* &= \arg \max_E \sum_{(v_i, v_j) \in E} \sum_{k=1}^m \ln p(l_i^k, l_j^k; c_{ij}^*) \\ &= \arg \max_E \sum_{(v_i, v_j) \in E} \underbrace{\ln \prod_{k=1}^m p(l_i^k, l_j^k; c_{ij}^*)}_{-q(v_i, v_j)} \\ &= \arg \min_E \sum_{(v_i, v_j) \in E} q(v_i, v_j) \end{aligned}$$

Man erhält das bekannte Problem des minimal aufspannenden Baumes (MST). Dieses kann mit dem bekannten Kruskal-Algorithmus in $O(n^2 \cdot \ln n)$ gelöst werden.