

# Discriminative Learning of Markov Random Fields for Segmentation of 3D Scan Data

Anguelov et. al., (CVPR), 2005

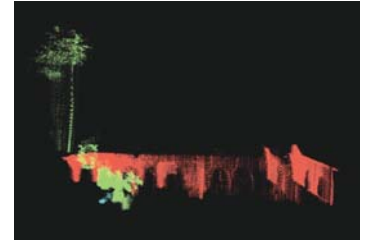
Presentation for CS 6784  
Sarah Iams, 18 Feb 2010

## Intuition for the problem

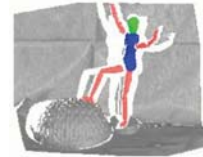
(1) ID vehicles vs background (synthetic data)



(2) Find buildings, trees, shrubs, ground

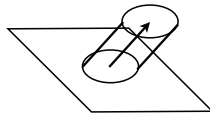
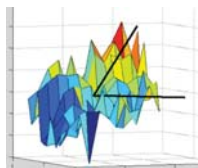


(3) Find head, limbs, torso, background



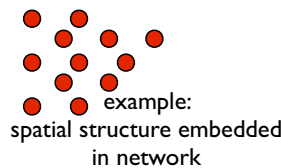
## Features

- How planar is the neighborhood of the point?
- Is a point close to the ground?
- Are there many points nearby?
- What are the principal components of the spin images?



## Capture problem structure

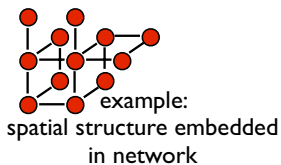
- Markov network captures geometry of the problem
- Scan points are represented by nodes in a graph
- Edges connect nearby scan points
- Each node will eventually have a label,  $Y_i \in \{1, \dots, K\}$
- The entire network is associated with a set of labels,  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_N\}$
- They are interested in a distribution over  $\{1, \dots, K\}^N$  specified by the geometry of the graph



example:  
one possible labeling  
 $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_N\}$

## Capture problem structure

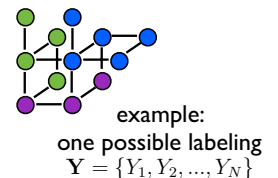
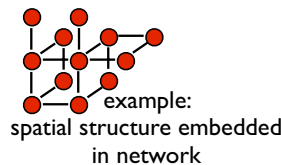
- Markov network captures geometry of the problem
- Scan points are represented by nodes in a graph
- Edges connect nearby scan points
- Each node will eventually have a label,  $Y_i \in \{1, \dots, K\}$
- The entire network is associated with a set of labels,  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_N\}$
- They are interested in a distribution over  $\{1, \dots, K\}^N$  specified by the geometry of the graph



example:  
one possible labeling  
 $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_N\}$

## Capture problem structure

- Markov network captures geometry of the problem
- Scan points are represented by nodes in a graph
- Edges connect nearby scan points
- Each node will eventually have a label,  $Y_i \in \{1, \dots, K\}$
- The entire network is associated with a set of labels,  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_N\}$
- They are interested in a distribution over  $\{1, \dots, K\}^N$  specified by the geometry of the graph



# Pairwise MRF assumption

- pairwise Markov network: nodes and edges are associated with potentials,  $\phi_i(Y_i)$  and  $\phi_{ij}(Y_i, Y_j)$
- all potentials are then multiplied (and normalized) to produce  $P(Y|X)$
- This is identical to saying the logs of the potentials are added to produce  $\log P(Y|X)$
- the feature values,  $\psi_i$ , at each node dictate the values of  $\phi_i(Y_i)$
- the similarity of the prospective labels,  $\psi_{ij}$ , along an edge dictates  $\phi_{ij}(Y_i, Y_j)$

$$P(Y|X) = \frac{1}{Z} \prod_i \phi_i(Y_i) \prod_{ij} \phi_{ij}(Y_i, Y_j)$$

$$\log P(Y|X) = \sum_i \log \phi_i(Y_i) + \sum_{ij} \log \phi_{ij}(Y_i, Y_j) - \log(Z)$$

$$\log \phi_i(k) = \mathbf{w}_n^k \cdot \psi_i$$

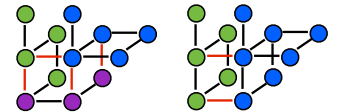
$$\log \phi_{ij}(k, l) = \mathbf{w}_e^{kl} \cdot \psi_{ij}$$

# AMN assumption

- want to find the  $\mathbf{Y}$  that maximized  $P(Y|X)$ . Note maximizing  $P(Y|X)$  is identical to maximizing  $\log P(Y|X)$
- they make one more assumption to simplify the optimization problem: edge weight is 0 when an edge connects nodes with different labels. Otherwise, the weight is non-negative.
- This is the associative Markov network assumption.

$$\arg \max_{\mathbf{Y}} \log P(Y|X) = \arg \max_{\mathbf{Y}} \left( \sum_i \log \phi_i(Y_i) + \sum_{ij} \log \phi_{ij}(Y_i, Y_j) - \log(Z) \right)$$

$$\begin{aligned} \text{---} &= 0 && \log \phi_{ij}(k, l) = 0 \text{ for } (k \neq l) \\ \text{---} &\geq 0 && \log \phi_{ij}(k, k) = \mathbf{w}_e^k \cdot \psi_{ij} \geq 0 \end{aligned}$$



# Optimization problem

$$\arg \max_{\mathbf{Y}} \log P(Y|X) = \arg \max_{\mathbf{Y}} \left( \sum_i \log \phi_i(Y_i) + \sum_{ij} \log \phi_{ij}(Y_i, Y_j) - \log(Z) \right)$$

$$\log \phi_i(k) = \mathbf{w}_n^k \cdot \psi_i$$

$$\log \phi_{ij}(k, l) = 0 \text{ for } (k \neq l)$$

$$\log \phi_{ij}(k, k) = \mathbf{w}_e^k \cdot \psi_{ij} \geq 0$$

$$\arg \max_{\mathbf{Y}} \log P(Y|X) = \arg \max_{\mathbf{Y}} \left( \sum_i \sum_k (\mathbf{w}_n^k \cdot \psi_i) y_i^k + \sum_{ij} \sum_k (\mathbf{w}_e^k \cdot \psi_{ij}) y_{ij}^k \right)$$

- Given weights, we can solve this (min-cut algorithm)
- Or (evidently), we can reformulate as integer program & relax to linear program: they choose this route because this arg max will reappear in the course of their learning method!

# Learning method

$$\begin{aligned} \arg \max_{\mathbf{Y}} \log P(Y|X) &= \arg \max_{\mathbf{Y}} \left( \sum_i \sum_k (\mathbf{w}_n^k \cdot \psi_i) y_i^k + \sum_{ij} \sum_k (\mathbf{w}_e^k \cdot \psi_{ij}) y_{ij}^k \right) \\ &= \arg \max_{\mathbf{Y}} \mathbf{wXy} \end{aligned}$$

- Switch to vector notation (all those subscripted w's,  $\psi$ 's & y's become vectors in a natural way, with  $\psi \rightarrow \mathbf{X}$ )
- They take a single training scene.
- Could train weights to maximize  $P(Y_{\text{correct}}|X)$
- Instead, maximize confidence in correct answer:  $P(Y_{\text{correct}}|X) - P(Y|X)$  (where  $Y_{\text{correct}}$  is the true label, and  $Y$  is any other labeling - this is maximum margin for the Markov network)
- Advantages: allows some kernelization later on
- Evidently pretty accurate

# M<sup>3</sup>N problem

$$\begin{aligned} \max_{\mathbf{y}} \text{ s.t. } \quad & \mathbf{wX}(\mathbf{y}_{\text{correct}} - \mathbf{y}) \geq \gamma \Delta(\mathbf{y}_{\text{correct}}, \mathbf{y}); \|\mathbf{w}\|^2 \leq 1 \\ & \Delta(\mathbf{y}_{\text{correct}}, \mathbf{y}) = N - \mathbf{y}_{\text{correct}}^T \mathbf{y} \end{aligned}$$

- Note that  $\mathbf{y}$  is an indicator vector, so when  $\mathbf{y}_{\text{correct}}$  and  $\mathbf{y}$  agree on a node label, that contributes to their dot product. When they disagree, it contributes 0 to the dot product.
- They define the loss function to count how many times  $\mathbf{y}$  is wrong on labeling the nodes. (Note M<sup>3</sup>N was approached without a loss function restriction on Tuesday).
- As usual, next they'll divide through by the margin ( $\gamma$ ) and add a slack variable (in case the data isn't separable)

# Primal formulation

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C\xi \text{ s.t. } \mathbf{wX}(\mathbf{y}_{\text{correct}} - \mathbf{y}) \geq N - \mathbf{y}_{\text{correct}}^T \mathbf{y} - \xi \quad \forall \mathbf{y}$$

- this is a quadratic program
- exponentially many constraints
- we can replace the constraints with a single constraint over a quadratic program!

$$\begin{aligned} \mathbf{wX}(\mathbf{y}_{\text{correct}} - \mathbf{y}) &\geq N - \mathbf{y}_{\text{correct}}^T \mathbf{y} - \xi \quad \forall \mathbf{y} \\ \Rightarrow \mathbf{wXy}_{\text{correct}} - N + \xi &\geq \mathbf{wXy} - \mathbf{y}_{\text{correct}}^T \mathbf{y} \quad \forall \mathbf{y} \\ \Rightarrow \mathbf{wXy}_{\text{correct}} - N + \xi &\geq \max_{\mathbf{y}} \mathbf{wXy} - \mathbf{y}_{\text{correct}}^T \mathbf{y} \end{aligned}$$

- we recognize this quadratic program from before
- Recall:  $\arg \max_{\mathbf{Y}} \log P(Y|X) = \arg \max_{\mathbf{Y}} \mathbf{wXy}$

# Switch to dual (twice)

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C\xi \quad \text{s.t.} \quad \mathbf{w}^T \mathbf{X} \mathbf{y}_{\text{correct}} - N + \xi \geq \max_{\mathbf{y}_{\text{correct, nodes}}} \mathbf{w}^T \mathbf{X} \mathbf{y} - \mathbf{y}_{\text{correct, nodes}}^T$$

- They switch to the dual problem in the constraint.

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C\xi \quad \text{s.t.} \quad \mathbf{w}^T \mathbf{X} \mathbf{y}_{\text{correct}} - N - \xi \geq \sum_{i=1}^N \alpha_i; \mathbf{w}_e \geq 0; \alpha_i - \sum_{ij} \alpha_{ij}^k \geq w_n^k \cdot \psi_i - y_{\text{correct},i}^k$$

$$\alpha_{ij}^k + \alpha_{ji}^k \geq w_e^k \cdot \psi_{ij}; \alpha_{ij}^k, \alpha_{ji}^k \geq 0$$

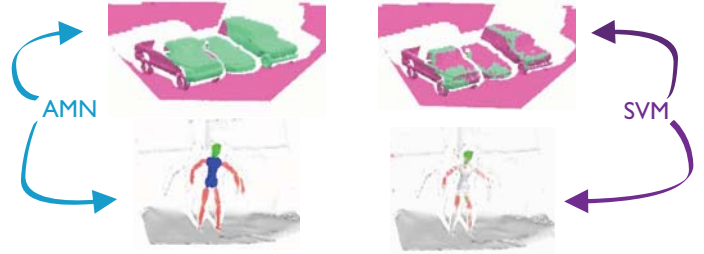
- Then they switch to the dual in the overall problem. (I am not including the dual here.) The primal and dual are related as follows:

$$w_n^k = \sum_{i=1}^N \psi_i (C y_{\text{correct},i}^k - \mu_i^k) \quad w_e^k = f(\phi_{ij}^k) + \sum_{ij} \psi_{ij} (C y_{\text{correct},ij}^k - \mu_{ij}^k)$$

- Since  $w_n^k$  is a sum over  $\psi_i$  multiplied by constants,  $w_n^k \psi$  can be kernelized. The edge potentials cannot be, however, because of the constant term added to the sum.

# Testing the AMN

- The associative Markov network ensures nearby points have the same label (SVM does not do this)
- After five training scenes:



# Testing the AMN

