

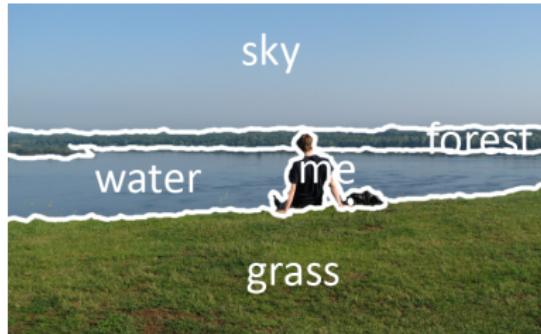
Object Detection with Occlusion

Josef Schulz

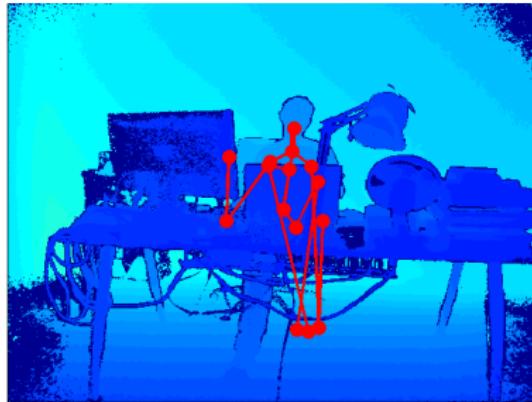
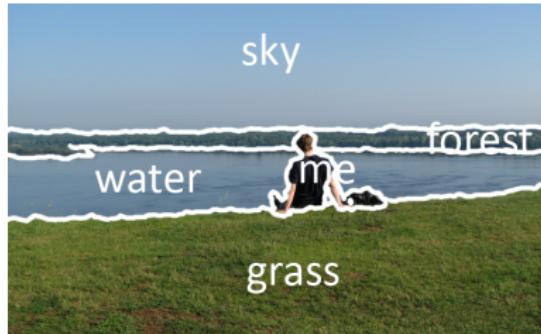
April 7, 2016

Example Problems

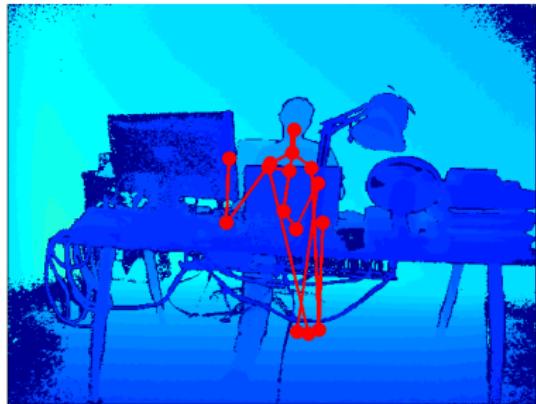
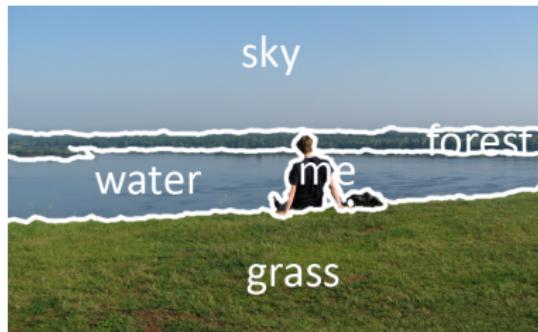
Example Problems



Example Problems



Example Problems



Content

1 Examples

2 Algorithms

- Semantic Occlusion Model
- Occlusion Patterns
- Robust Instance Recognition

3 Conclusions

4 Discussion

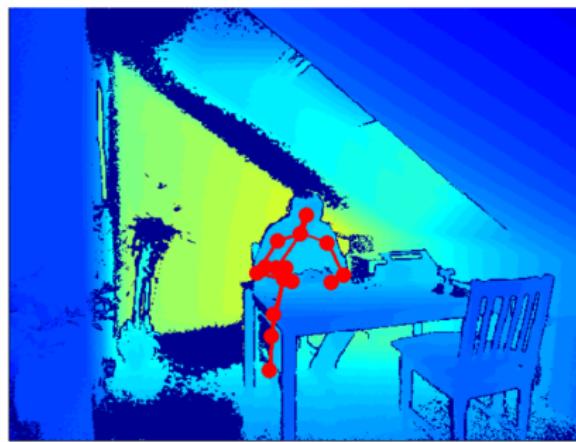
5 References

A Semantic Occlusion Model for Human Pose Estimation

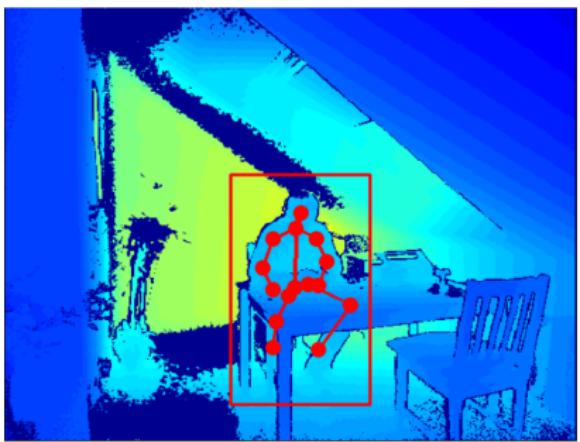
Input : Single Depthimage

Output : Estimated Poses of all Parts

3D Human Pose Estimation



Kinect2 SDK



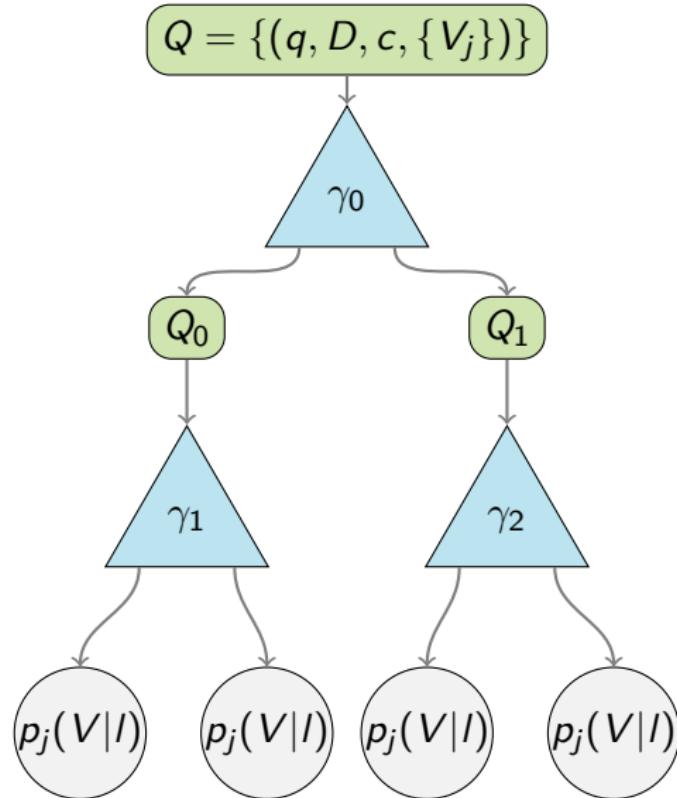
OARF

Trainingsset

$$\{(q, D, c, \{V_j\}), \dots\}$$

- ▶ q is a pixel position.
- ▶ D reference image x .
- ▶ c class label, according to a limb
- ▶ $\{V_j\}$ is set of vectors $V_j = q_j - q$

Regression Tree



slitnode

$$\Phi_{\gamma}(q, D) \mapsto \{0, 1\}$$

$$\Phi_{\gamma}(q, D) = \begin{cases} 1 & \text{if } D(q + \frac{\textcolor{blue}{u}}{D(q)}) - D(q + \frac{\textcolor{blue}{v}}{D(q)}) > \tau \\ 0 & \text{else} \end{cases}$$

$$\gamma = (\textcolor{blue}{u}, \textcolor{blue}{v}, \tau)$$

choose splitting function

$$\Phi^* = \arg \max_{\Phi} g(\Phi)$$

$$g(\Phi) = H(Q) - \sum_{s \in \{0,1\}} \frac{|Q_s(\Phi)|}{|Q|} H(Q_s(\Phi))$$

$$H(Q) = - \sum_c p(c|Q) \log(p(c|Q))$$

leafnode

A leafnode I stores $p_j(V, I)$

V_j is clustered by mean-shift with Gaussian Kernel with bandwidth b

$$p_j(V|I) \propto \sum_{k \in K} w_{Ijk} \cdot \exp\left(-\left\|\frac{V - V_{Ijk}}{b}\right\|_2^2\right)$$

w_{Ijk} is determined by offset vectors ended in the cluster k .

Pose Estimation

a set of pixels q are sampled in D

$$x_j = q + V_{ijk}$$

$$w_j = w_{ijk} \cdot D^2(q)$$

$$X_j = \{(x_j, w_j)\}$$

$$p_j(x|D) \propto \sum_{(x_j, w_j) \in X_j} w_j \cdot \exp\left(-\left\|\frac{x - x_j}{b_j}\right\|_2^2\right)$$

Votes are clustered and only the cluster with the highest summed weights w_j are used for prediction.

Occlusion Aware Regression Forests

$Q_{ext} = Q \cup Q_{occ}$ with

$$Q_{occ} = \{(q_{occ}, D, c_{occ}, \{v_{jocc}\})\}$$

$$p_j(V_{occ}|I) \propto \sum_{k \in K} w_{ljk} \cdot \exp(-\|\frac{V_{occ} - V_{ljkocc}}{b}\|_2^2)$$

$$\mathcal{X}_j = \{(x_j, w_j)\} \cup \{(x_{jocc}, w_{jocc})\}$$

Training Data

Synthetic Data (552 images)

- ▶ Human Poses from CMU-Database
- ▶ body part labels for each pixel

Real Data (552 images)

- ▶ Kinect2 SDK, all fails are discarded

Regression Forest

- ▶ 3 regression trees
- ▶ max depth = 20
- ▶ $b = 0.05m$

- ▶ 1000 samples per image
- ▶ samples 2000 splitting functions per node

Occlusion Patterns for Object Class Detection

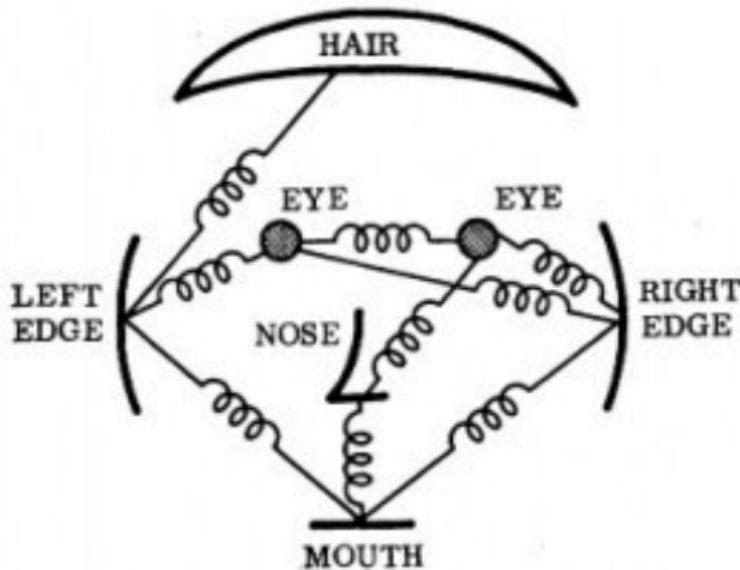
Input : Single RGB-Image

Output : Object-Boundingboxes

Deformable Models approach

- ▶ Consider each object as a deformed version of a template
- ▶ Compact representation

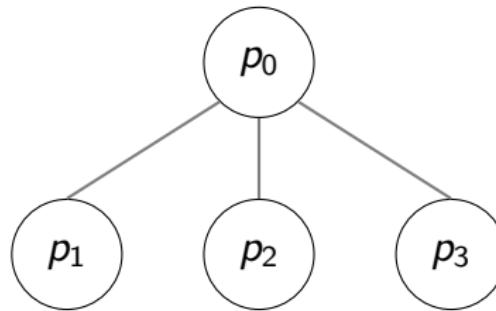
Matching model to image involves joint optimization of part locations "stretch and fit"



Model

Model is represented by a Graph $G = (V, E)$

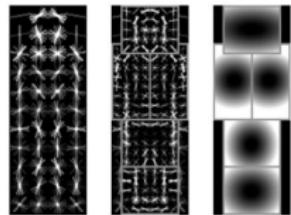
- ▶ $V = p = \{p_0, \dots, p_M\}$ are the parts
- ▶ p_i is parameterized through their bounding box (l_i, r_i, t_i, b_i)
- ▶ $(p_i, p_j) \in E$



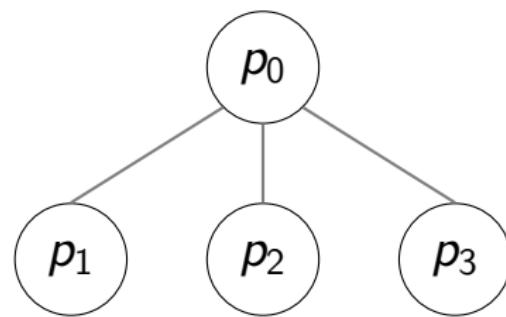
$$E_c(p; I) = \sum_{i=0}^M \langle v_i^c, \Phi(p_i; I) \rangle + \sum_{i=1}^M \langle w_i^c, \Phi(p_0, p_i) \rangle$$



Training

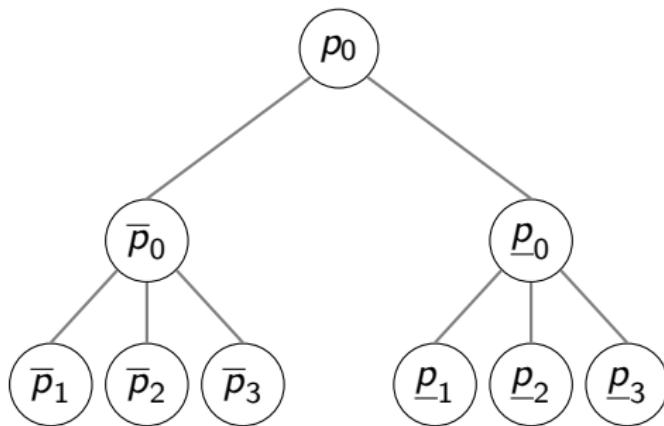


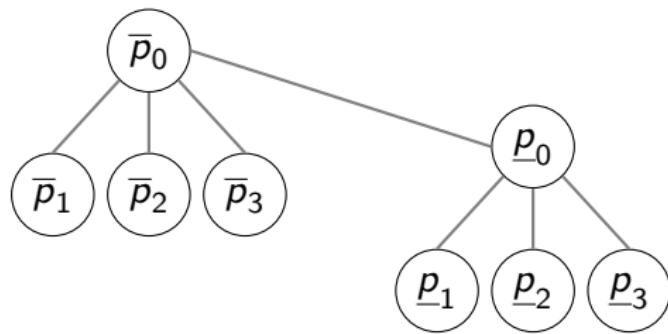
$$C = \{1, \dots, C_{visible}\} \cup C_{invisible}$$

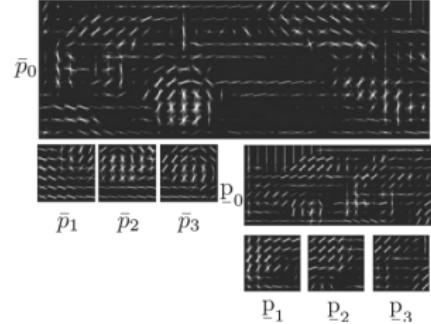
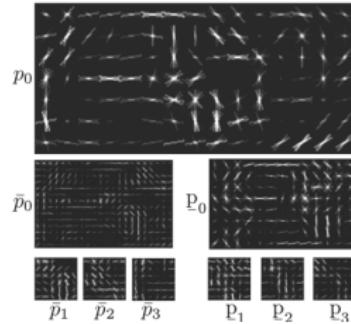
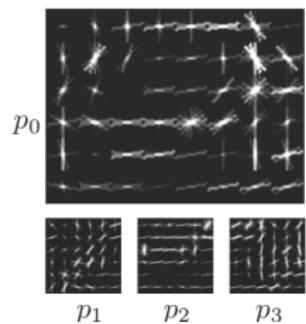
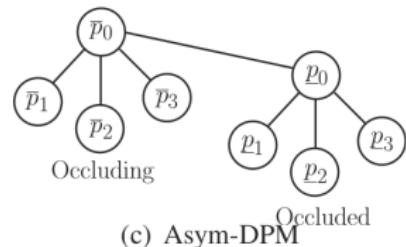
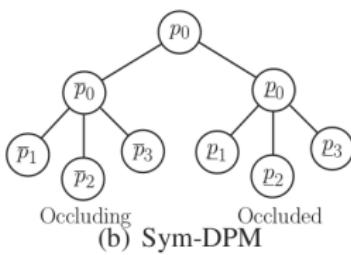
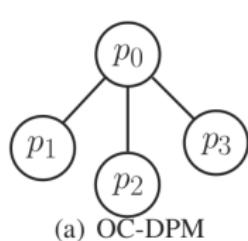


SYM-DPM

$$E'_c(p; I) = \langle v_{joint}^c, \Phi(p_0; I) \rangle + \langle \bar{w}, \Phi(\bar{p}_0, p_{joint}) \rangle + \langle \underline{w}, \Phi(\underline{p}_0, p_{joint}) \rangle \\ + E_c(\bar{p}_0; I) + E_c(\underline{p}_0; I)$$







KITTI

KITTI contains 7481 images

	#objects	#occluded objects	%
Car	28521	15231	53.4
Pedest.	4445	1805	40.6
Cycles	1612	772	44.5

components:

visible 6

occluded 16 – 15

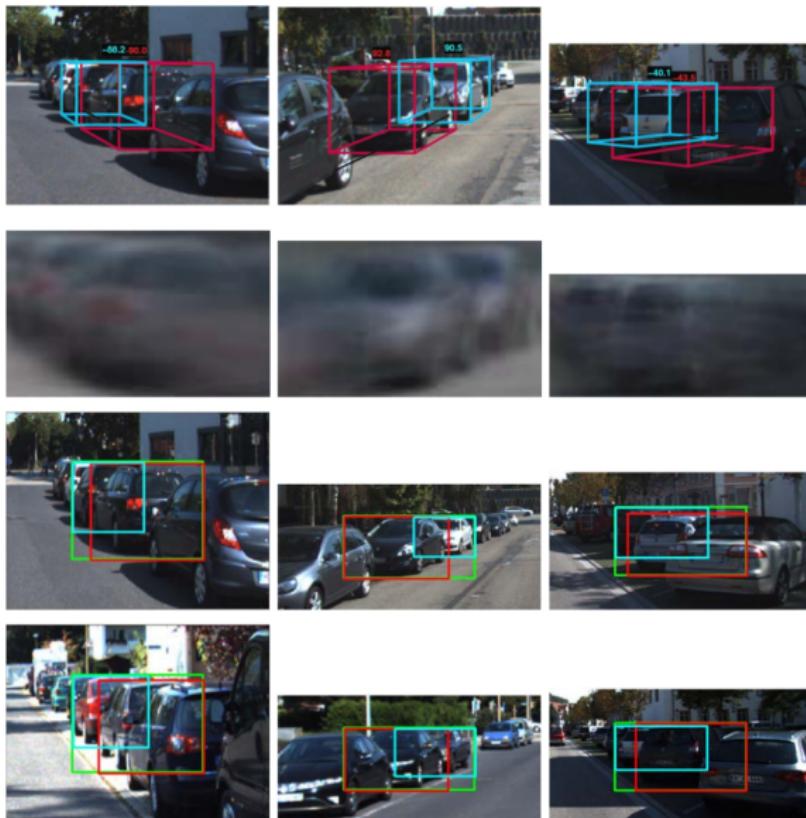
Mining Trainingsdata

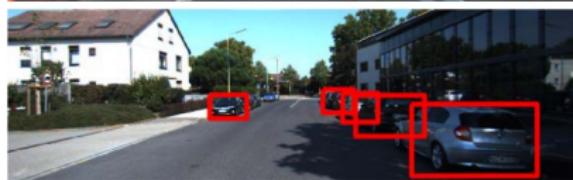
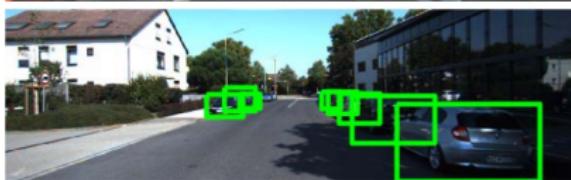
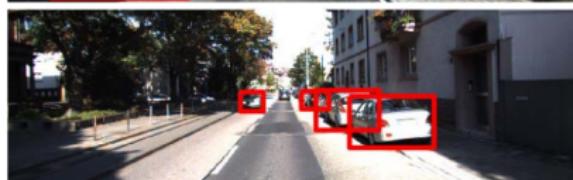
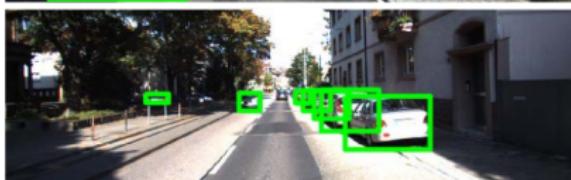
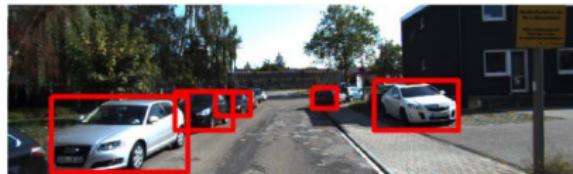
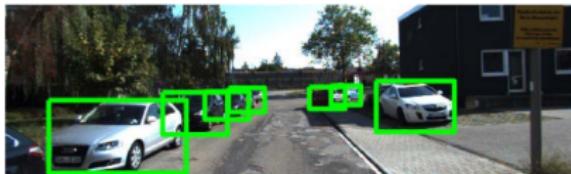
Feature Space:

- i occluder left/right of occludee
- ii orientation of occluder/occludee
- iii occluder is/is not occluded
- iv degree of occlusion of occludee

Rule-based clustering

Mining Trainingsdata



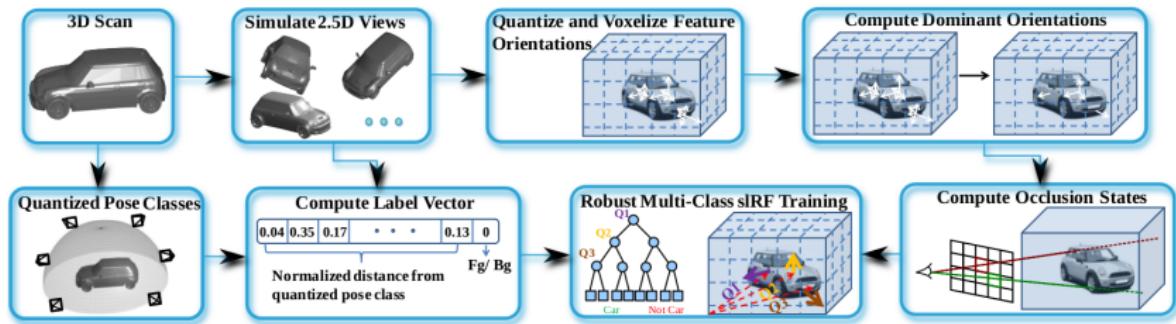


Robust Instance Recognition in Presence of Occlusion and Clutter

Input : 5-10 consecutive frames as one Pointcloud

Output : 6D-Object-Pose

Overview



Edgelet

N points per Pointcloud j

FOR ALL $i \in \{1, \dots, N\}$

calc λ_1 and λ_2

$$r = \frac{\lambda_1}{\lambda_2}$$

$r \rightarrow \text{curvatureMap}$

`hysteresisThresholding(curvatureMap);`

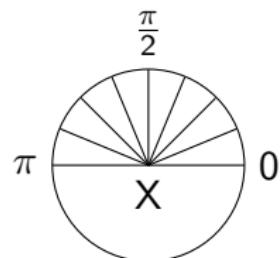
`nonmaximalSuppression(curvatureMap);`

`hysteresisThresholding(depthMap);`

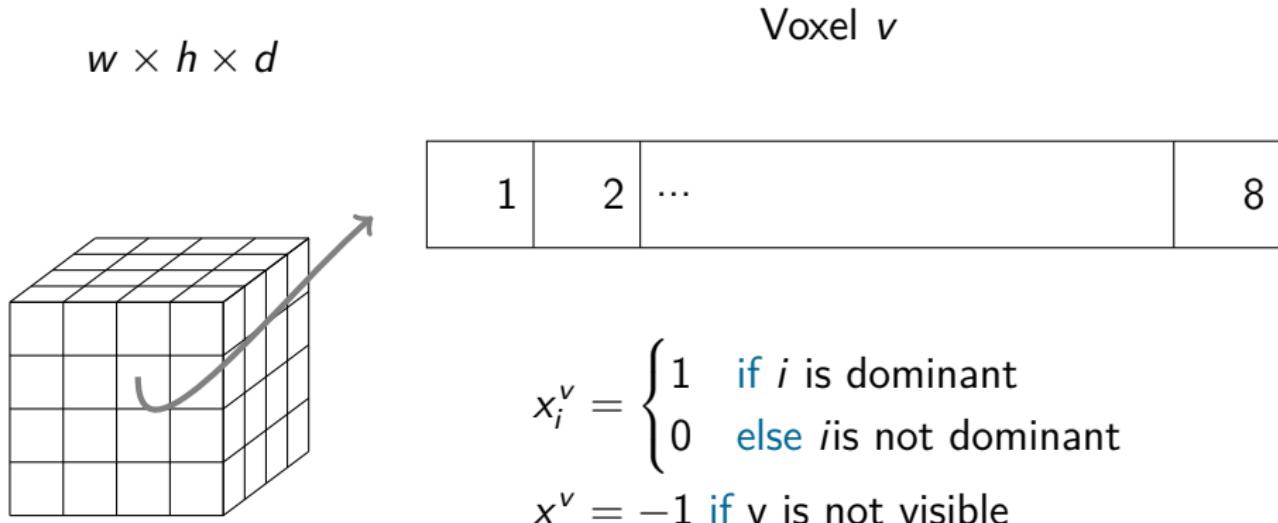
`projectToPointcloud(curvatureMap, depthMap);`

RANSAC line fitting

orientation to 8 bins // (direction % π)



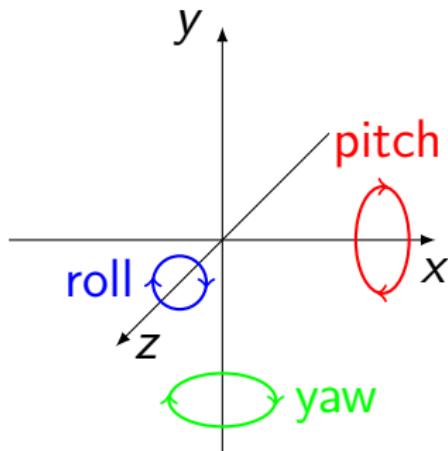
Feature Vector



The resulting feature vector is the concatenation of all voxels:

$$w \times h \times d \times 8$$

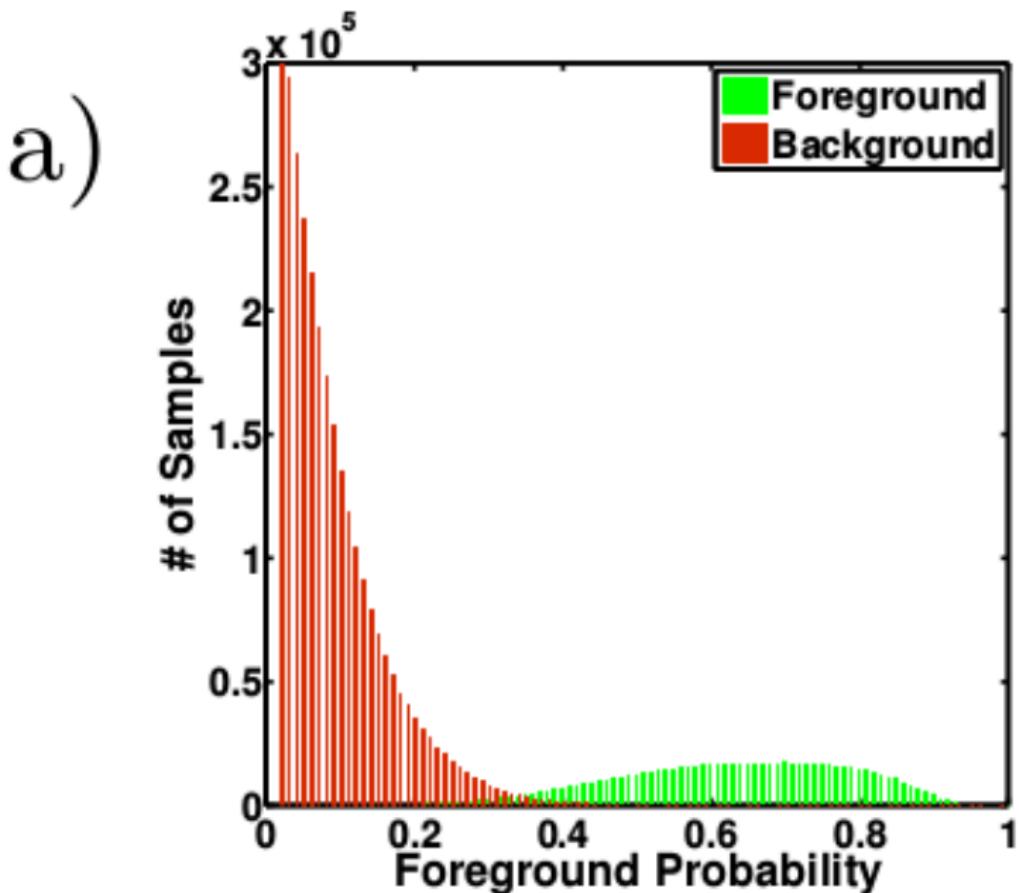
Soft Label Random Forest



- ▶ 16 pose classes
- ▶ +1 class = $\begin{cases} 1 & \text{if bg} \\ 0 & \text{else} \end{cases}$
- ▶ $d_j^i = \|I - R_j^i\|_F$

$$l_j^i = \exp(-d_j^{i2}), i \in \{1, \dots, 16\}$$

IF fg THEN $1 = \sum_{i=1}^{16} l_j^i$ ELSE 0

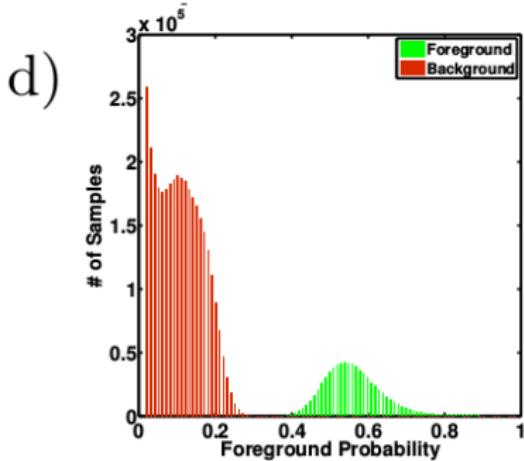
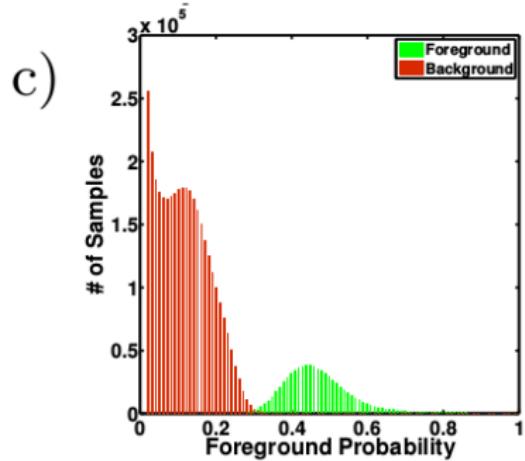
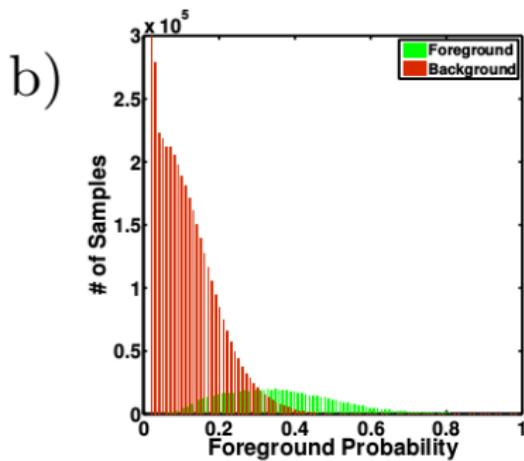
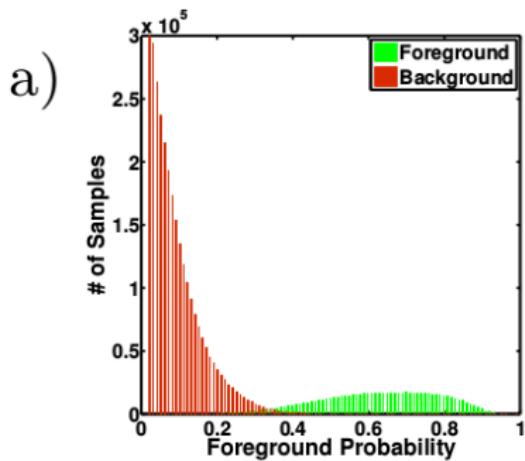


training scheme for sLRF

Input: $X = \{x_j, I_j\}$

Output: Learnt sLRF classifier

1. $X_s \subset X$, $|X_s| = |X|/20$
2. Train sLRF with X_s , compute p_{fg} with X
3. add borderline positive (low p_{fg}) and borderline negative (high p_{fg})
4. add confusing samples
5. compute d_L for all positive samples, add samples with high d_L
6. repeat 2-5 till p_{fg} for all positive data is greater than p_{fg} for all negative data.



Conclusion

- ▶ redundanze
- ▶ training with occlusion
- ▶ semantic (scene understanding)

Discussion

?



Ujwal Bonde, Vijay Badrinarayanan, and Roberto Cipolla.
“Computer Vision – ECCV 2014: 13th European Conference,
Zurich, Switzerland, September 6-12, 2014, Proceedings, Part
II”. In: ed. by David Fleet et al. Cham: Springer International
Publishing, 2014. Chap. Robust Instance Recognition in
Presence of Occlusion and Clutter, pp. 520–535. ISBN:
978-3-319-10605-2. DOI: [10.1007/978-3-319-10605-2_34](https://doi.org/10.1007/978-3-319-10605-2_34).
URL:
http://dx.doi.org/10.1007/978-3-319-10605-2_34.



Andreas Geiger. “Are We Ready for Autonomous Driving?
The KITTI Vision Benchmark Suite”. In: *Proceedings of the
2012 IEEE Conference on Computer Vision and Pattern
Recognition (CVPR)*. CVPR '12. Washington, DC, USA:
IEEE Computer Society, 2012, pp. 3354–3361. ISBN:
978-1-4673-1226-4. URL: <http://dl.acm.org/citation.cfm?id=2354409.2354978>.



Stefan Hinterstoisser et al. “Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes.” In: *ICCV*. Ed. by Dimitris N. Metaxas et al. IEEE Computer Society, 2011, pp. 858–865. ISBN: 978-1-4577-1101-5. URL: <http://dblp.uni-trier.de/db/conf/iccv/iccv2011.html#HinterstoisserHCIKNL11>.



Alexander Krull et al. “6-DOF Model Based Tracking via Object Coordinate Regression.” In: *ACCV (4)*. Ed. by Daniel Cremers et al. Vol. 9006. Lecture Notes in Computer Science. Springer, 2014, pp. 384–399. ISBN: 978-3-319-16816-6. URL: <http://dblp.uni-trier.de/db/conf/accv/accv2014-4.html#KrullMBGIR14>.



Bojan Pepik et al. “Occlusion Patterns for Object Class Detection.” In: *CVPR*. IEEE, 2013, pp. 3286–3293. URL: <http://dblp.uni-trier.de/db/conf/cvpr/cvpr2013.html#PepikSGS13>.

-  Umer Rafi, Juergen Gall, and Bastian Leibe. "A Semantic Occlusion Model for Human Pose Estimation From a Single Depth Image". In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. June 2015.
-  Training Deformable Partmodels. *unkown*. [Online; accessed March 18, 2016]. 2016. URL:
<http://i.stack.imgur.com/Awffz.jpg>.
-  unkown. *unkown*. [Online; accessed April 4, 2016]. unkown. URL:
http://3.bp.blogspot.com/_IDEWIOP9RbA/TB6VYXUN_CI/AAAAAAAACA/D1lZ1rwWktA/s1600/semsegm.png.
-  unkown. *unkown*. [Online; accessed April 4, 2016]. unkown. URL: https://farm6.staticflickr.com/5668/21114849206_759853e42a_z.jpg.