Deep Hough Voting for 3D Object Detection in Point Clouds

Charles Qi (祁芮中台) GAMES Webinar December 5th, 2019

Joint work with Or Litany, Kaiming He, Leonidas Guibas. ICCV 2019.

3D object detection

Estimate oriented 3D bounding boxes and semantic classes from sensor data.







Prior work relies on 2D object detection

Bird's eye view detector



[MV3D by Chen et al. CVPR 2017]

Frustum-based detector



[F-PointNet by Qi et al. CVPR 2018]

Prior work relies on 2D object detection

3D CNN detector



[Deep Sliding Shapes by Song et al. CVPR 2016]

Observation: 2D v.s. 3D



Our idea: "ask" the points to vote for object centers



GENERALIZING THE HOUGH TRANSFORM TO DETECT ARBITRARY SHAPES*





vote for center of object

From U. Toronto CSC420

Hough voting pipeline (on 2D images):

- Select interest points
- Match patch around each interest point to a training patch (codebook)
- Vote for object center given that training instance



vote for center of object

From U. Toronto CSC420

Hough voting pipeline (on 2D images):

- Select interest points
- Match patch around each interest point to a training patch (codebook)
- Vote for object center given that training instance



vote for center of object

Hough voting pipeline (on 2D images):

- Select interest points
- Match patch around each interest point to a training patch (codebook)
- Vote for object center given that training instance



of course some wrong votes are bound to happen...

Hough voting pipeline (on 2D images):

- Select interest points
- Match patch around each interest point to a training patch (codebook)
- Vote for object center given that training instance



But that's ok. We want only peaks in voting space.

Hough voting pipeline (on 2D images):

- Select interest points
- Match patch around each interest point to a training patch (codebook)
- Vote for object center given that training instance
- Votes clustering to find peaks



Find patches that voted for the peaks (back-projection).

Hough voting pipeline (on 2D images):

- Select interest points
- Match patch around each interest point to a training patch (codebook)
- Vote for object center given that training instance
- Votes clustering to find peaks
- Find patches that voted for the peaks by back-projection



Find full objects based on the back-projected patches.

Hough voting pipeline (on 2D images):

- Select interest points
- Match patch around each interest point to a training patch (codebook)
- Vote for object center given that training instance
- Votes clustering to find peaks
- Find patches that voted for the peaks by back-projection
- Find full objects based on back-projected patches



- + Computation is only on "interest" points instead of on all pixels/voxels.
- + Support "templates" (used in 6DoF pose estimation)
- Not end-to-end optimizable

3D object proposal: A return of hough voting!

Deep hough voting with PointNet++

Interest points \rightarrow seed points sampled from the point clouds

Votes \rightarrow learned mapping from point features to votes

Clustering \rightarrow local pointnet layers to group and aggregate local votes

Object recovery \rightarrow learned bounding box predictor

End-to-end optimizable!

Deep Hough voting: Detection pipeline



Deep Hough voting: Detection pipeline



Results: SUN RGB-D (single depth images)



Results: ScanNet (3D reconstructions)





Comparing with previous methods

	Input	bathtub	bed	bookshel	lf chair	desk	dresser	nightsta	nd sofa	table	toilet	mAP
DSS [42]	Geo + RGB	44.2	78.8	11.9	61.2	20.5	6.4	15.4	53.5	50.3	78.9	42.1
COG [38]	Geo + RGB	58.3	63.7	31.8	62.2	45.2	15.5	27.4	51.0	51.3	70.1	47.6
2D-driven [20]	Geo + RGB	43.5	64.5	31.4	48.3	27.9	25.9	41.9	50.4	37.0	80.4	45.1
F-PointNet [34]	Geo + RGB	43.3	81.1	33.3	64.2	24.7	32.0	58.1	61.1	51.1	90.9	54.0
VoteNet (ours)	Geo only	74.4	83.0	28.8	75.3	22.0	29.8	62.2	64.0	47.3	90.1	57.7

SUN RGB-D: +3.7mAP with just 3D geometry data as input.

Comparing with previous methods

	Input	mAP@0.25	mAP@0.5
DSS [42, 12]	Geo + RGB	15.2	6.8
MRCNN 2D-3D [11, 12]	Geo + RGB	17.3	10.5
F-PointNet [34, 12]	Geo + RGB	19.8	10.8
GSPN [54]	Geo + RGB	30.6	17.7
3D-SIS [12]	Geo + 1 view	35.1	18.7
3D-SIS [12]	Geo + 3 views	36.6	19.0
3D-SIS [12]	Geo + 5 views	40.2	22.5
3D-SIS [12]	Geo only	25.4	14.6
VoteNet (ours)	Geo only	58.6	33.5

ScanNet: +18.3 mAP compared with prior art (3D CNN based method) with 3D & multi-view images.

Can images help the VoteNet detection?





Images are in high resolution, have rich texture, and can even provide useful geometric cues for object localization & shape/pose estimation.

ImVoteNet: Boosting 3D Object Detection in Point Clouds with Image Votes

On-going work with Xinlei Chen, Or Litany and Leonidas Guibas









Geometric cues from images: Lifted image votes





Results on SUN RGB-D

+5.7mAP with lifted image cues for voting



Results on SUN RGB-D



36

Summary

VoteNet: a revival of Hough voting with 3D deep learning.

- End-to-end optimizable hough voting with point cloud deep nets.
- A new detection model with a simple design shows state-of-the-art results on SUN RGB-D and ScanNet with geometry data only.

Code: <u>https://github.com/facebookresearch/votenet</u>

ImVoteNet: boosting 3D detection with lifted image votes.

Many open possibilities to extend the pipeline (e.g. 6D pose estimation, template based detection).