

Multiple object tracking

Different challenges

- Multiple objects of the same type
- Heavy occlusions
- Appearance is often very similar





Tracking-by-detection

- We will focus on algorithms where a set of detections is provided
 - Remember detections are not prefect!



Find detections that match and form a trajectory

Online vs offline tracking

- Online tracking
 - Processes two frames at a time
 - For real-time applications
 - Prone to drifting \rightarrow hard to recover from errors or occlusions
- Offline tracking
 - Processes a batch of frames
 - Good to recover from occlusions (short ones as we will see)
 - Not suitable for real-time applications
 - Suitable for video analysis





• 1. Track initialization (e.g. using a detector)







• 1. Track initialization (e.g. using a detector)



• 2. Matching detections at t with detections at t+1



• 1. Track initialization (e.g. using a detector)



- 2. Matching detections at t with detections at t+1
- Repeat for every pair of frames

 3. Matching detections at t with detections at t+1



Detections

- Bipartite matching
 - Define distances between boxes
 - (e.g., IoU, pixel distance,
 - 3D distance)



Detections

- Bipartite matching
 - Define distances between boxes
 - (e.g., IoU, pixel distance, 3D distance)
 - Solve the unique matching with e.g., the Hungarian algorithm*



*Demo: http://www.hungarianalgorithm.com/solve.php

Detections

- Bipartite matching
 - Define distances between boxes
 - (e.g., IoU, pixel distance, 3D distance)
 - Solve the unique matching with e.g., the Hungarian algorithm*
 - Solutions are the unique assignments that minimize the total cost



*Demo: http://www.hungarianalgorithm.com/solve.php

- Problems with frame-by-frame tracking
 - Cannot recover from errors. If a detection is missing from a frame, we have to end the trajectory.
 - All decisions are essentially local
 - Hard to recover from errors in the detection step

• Solution: find the minimum cost solution for ALL frames and ALL trajectories



Graph-based MOT









• Node = detection

• Edge = flow = trajectory

• 1 unit of flow = 1 pedestrian

• Solving the Minimum Cost Flow Problem

"Determine the minimum cost of shipment of a commodity through a network"



Ravindra K. Ahuja; Thomas L. Magnanti & James B. Orlin. "Network Flows: Theory, Algorithms, and Applications". Prentice-Hall, Inc. 1993

• Solving the Minimum Cost Flow Problem

"Determine the minimum cost of shipment of a commodity through a network"



• Objective function $\mathcal{T}^* = \underset{\mathcal{T}}{\operatorname{arg\,min}} \sum_{i,j} C(i,j) f(i,j)$



• Objective function $\mathcal{T}^* = \underset{\mathcal{T}}{\operatorname{arg\,min}} \sum_{i,j} C(i,j) f(i,j)$







Transition: $cost \propto distance$ between detections

FLOW = TRAJECTORY = PEDESTRIAN



Transition: $cost \propto distance$ between detections

FLOW = TRAJECTORY = PEDESTRIAN



Entrance/exit: cost to start or end a trajectory

Transition: $cost \propto distance$ between detections

Flow can only start at the source node and end at the sink node





CV3DST | Prof. Leal-Taixé

• We need a negative cost

U

Detection edge



f(x) = log(x/(1-x))

Zhang et al. "Global Data Association for Multi-Object Tracking Using Network Flows". CVPR 2008

Complete graph





Connections that allow us to start a trajectory

31

Complete graph

Connections that allow us to end a trajectory





Linear Program MOT formulation

Why a linear program?

• Fast solvers (e.g., Simplex algorithm)

Guaranteed to converge to the global optimum

Minimum cost flow problem

• Objective function

$$\mathcal{T}* = \underset{\mathcal{T}}{\arg\min} \sum_{i,j} C(i,j) f(i,j)$$

Constraints

• Objective function

 $f_{det}(i)$

 $\sum f_t(i,j)$



 $f_{\rm det}(i)$

 $\sum f_t(j,i)$
Constraints

• Objective function

$$\mathcal{T}* = \underset{\mathcal{T}}{\operatorname{arg\,min}} \sum_{i,j} C(i,j) f(i,j)$$

• Subject to

Flow conservation at the nodes

$$f_{in}(i) + f_{det}(i) = \sum_{j} f_t(i,j) \qquad \qquad \sum_{j} f_t(j,i) = f_{out}(i) + f_{det}(i)$$

 $\label{eq:fin} \begin{array}{ll} \mbox{Edge capacities} & \mbox{NP-hard!!} \\ f_{\rm in}(\mathfrak{i}) + f_{\rm det}(\mathfrak{i}) \in \{0,1\} & f_{\rm out}(\mathfrak{i}) + f_{\rm det}(\mathfrak{i}) \in \{0,1\} & f \in \{0,1\} \end{array}$

LP relaxation

• Objective function

$$\mathcal{T}* = \underset{\mathcal{T}}{\arg\min} \sum_{i,j} C(i,j) f(i,j)$$

• Subject to

Flow conservation at the nodes

$$f_{\mathrm{in}}(\mathfrak{i}) + f_{\mathrm{det}}(\mathfrak{i}) = \sum_{\mathfrak{j}} f_{\mathrm{t}}(\mathfrak{i},\mathfrak{j}) \qquad \qquad \sum_{\mathfrak{j}} f_{\mathrm{t}}(\mathfrak{j},\mathfrak{i}) = f_{\mathrm{out}}(\mathfrak{i}) + f_{\mathrm{det}}(\mathfrak{i})$$

 $\begin{array}{ll} \mbox{Edge capacities} & \mbox{LP-relaxation} \\ 0 \leq f_{\rm in}(\mathfrak{i}) + f_{\rm det}(\mathfrak{i}) \leq 1 & \mbox{0} \leq f_{\rm out}(\mathfrak{i}) + f_{\rm det}(\mathfrak{i}) \leq 1 & \mbox{0} \leq f \leq 1 \end{array}$

Solver

• Objective function

$$\mathcal{T}* = \underset{\mathcal{T}}{\arg\min} \sum_{i,j} C(i,j) f(i,j)$$

Given the shape of the constraints (total unimodularity), we solve the relaxed problem and still get integer solutions.

Objective function

• Objective function

$$\mathcal{T}* = \underset{\mathcal{T}}{\arg\min} \sum_{i,j} C(i,j) f(i,j)$$

Objective function

• Objective function

$$\mathcal{T}* = \underset{\mathcal{T}}{\arg\min} \sum_{i} C_{in}(i) f_{in}(i) + \sum_{i,j} C_{t}(i,j) f_{t}(i,j)$$
$$+ \sum_{i} C_{det}(i) f_{det}(i) + \sum_{i} C_{out}(i) f_{out}(i)$$
$$C(i) = -\log P(i)$$

Equivalent to Maximum a-posteriori tracking formulation

$$\mathcal{T}* = \underset{\mathcal{T}}{\arg \max} \prod_{j} P(\mathbf{o}_{j}|\mathcal{T})P(\mathcal{T})$$

Two ways forward

- 1. Improving the costs (aka more learning)
 - L. Leal-Taixé et al. "Learning an image-based motion context for multiple people tracking". CVPR 2014.
 - L. Leal-Taixé et al. "Everybody needs somebody: Modeling social and grouping behavior on a linear programming multiple people tracker". ICCVW 2011
 - L. Leal-Taixé et al. "Learning by tracking: Siamese CNN for robust target association". CVPRW 2016.
 - S. Schulter et al. "Deep network flow for multi-object tracking". CVPR 2017.
 - J. Son at al. "Multi-object tracking with quadruplet convolutional neural networks". CVPR 2017.
 - Ristani and Tomasi. "Features for multi-target multi-camera tracking and re-identification". CVPR 2018
 - J. Xu et al. "Spatial-temporal relation networks for multi-object tracking". ICCV 2019.

Two ways forward

- 2. Making the graph more complex (including more connections)
 - M. Keuper et al. "Motion segmentation and multiple object tracking by correlation co-clustering". PAMI 2018.
 - S. Tang et al. "Subgraph decomposition for multi-target tracking:. CVPR 2015.
 - S. Tang et al. "Multiple people tracking by lifted multicut and person reidentification". CVPR 2017

End-to-end learning?

- Can we learn:
 - Features for multi-object tracking (e.g., costs)
 - To do data association, i.e., find a solution on the graph

• We will exploit the graph structure we have just seen and perform end-to-end learning



Message Passing Networks

Deep learning on graphs

- Several works have proposed generalizations of Neural Networks that can operate on graph-structured domains:
 - Scarselli et al. "The Graph Neural Network Model". IEEE Trans. Neur. Net 2009.
 - Kipf et al. "Semi-Supervised Classification with Graph Convolutional Networks. ICLR 2016.
 - Gilmer et al. "Neural Message Passing for Quantum Chemistry". ICML 2017
 - Battaglia et al. "Relational inductive biases, deep learning, and graph networks". arXiv 2018 (review paper)
- Key challenges:
 - Variable sized inputs (number of nodes and edges)
 - Need invariance to node permutations

General Idea



General Idea



Learning to propagate information

• We can divide the propagation process in two steps: 'node to edge' and 'edge to node' updates.



'Edge to Node' Update

Repeat these two updates for a fixed number of iterations (message passing steps) in order to encode context into embeddings

Node embeddings Edge embeddings

- Notation:
 - connects nodes I and j

Embedding of edge that

- Graph: G = (V, E)- Initial embeddings: $h_{(i,j)}^{(0)}$, $(i,j) \in E$ $h_i^{(0)}$, $i \in V$ - Embeddings after l steps: $h_{(i,j)}^{(l)}$, $(i,j) \in E$ $h_i^{(l)}$, $i \in V$

Embedding of

node i

- Notation:
 - Graph: G = (V, E)- Initial embeddings: $h_{(i,j)}^{(0)}, (i,j) \in E$ $h_i^{(0)}, i \in V$ - Embeddings after l steps: $h_{(i,j)}^{(l)}, (i,j) \in E$ $h_i^{(l)}, i \in V$
- At every message passing step l , first do:

 $h_{(i,j)}^{(l)} = \mathcal{N}_e\left(\left[h_i^{(l-1)}, h_{(i,j)}^{(l-1)}, h_j^{(l-1)}\right]\right)$ Embedding of node j in Embedding of edge (i,j) Embedding of node i in i in the precious the precious message the precious message passing step message passing step passing step CV3DST | Prof. Leal-Taixé 53

- Notation:
 - Graph: G = (V, E)- Initial embeddings: $h_{(i,j)}^{(0)}, (i,j) \in E$ $h_i^{(0)}, i \in V$ - Embeddings after l steps: $h_{(i,j)}^{(l)}, (i,j) \in E$ $h_i^{(l)}, i \in V$
- At every message passing step l , first do:



Combine node and edge embeddings

- Notation:

 - Graph: G = (V, E)- Initial embeddings: $h_{(i,j)}^{(0)}, (i,j) \in E$ $h_i^{(0)}, i \in V$ Embeddings after l steps: $h_{(i,j)}^{(l)}, (i,j) \in E$ $h_i^{(l)}, i \in V$
- At every message passing step l , first do:



CV3DST | Prof. Leal-Taixé

'Edge to node' updates

• After a round of edge updates, each edge embedding contains information about its pair of incident nodes

'Edge to node' updates

- After a round of edge updates, each edge embedding contains information about its pair of incident nodes
- Then, edge embeddings are used to update nodes:

$$m_{(i,j)}^{(l)} = \mathcal{N}_v \left([h_i^{(l-1)}, h_{(i,j)}^{(l)}] \right)$$

Learnable function (e.g. MLP) with shared weights across the entire graph



'Edge to node' updates

- After a round of edge updates, each edge embedding contains information about its pair of incident nodes
- Then, edge embeddings are used to update nodes:

$$m_{(i,j)}^{(l)} = \mathcal{N}_v \left([h_i^{(l-1)}, h_{(i,j)}^{(l)}] \right)$$
$$h_i^{(l)} = \Phi\left(\left\{ m_{(i,j)}^{(l)} \right\}_{j \in N_i} \right)$$

Permutation invariant operation (e.g. sum, mean, max)

Neighbors of node i

The aggregation provides each node embedding with contextual information about its neighbors

Remarks

- Main goal: obtaining node and edge embeddings that contain *context information* encoding graph topology and neighbor's feature information.
- After repeating the node and edge updates for l steps, each node (resp. edge) embedding contains information about all nodes (resp. edge) at distance l (resp. l − 1) → Think of iterations as layers in a CNN
- Observe that all operations used are differentiable, hence, MPNs can be used within end-to-end pipelines
- There is vast literature on different instantiations, as well as variations of the MPN framework we presented. See Battaglia et al. for an extensive review.



MOT with Message Passing Networks

Overview



61



Encode appearance and scene geometry cues into node and edge embeddings







(a) Input









(e) Output

(c) Neural Message Passing

(d) Edge Classification

Overview

Propagate cues across the entire graph with





Learn to directly predict solutions of the Min-Cost Flow problem by classifying edge embeddings





CV3DST | Prof. Leal-Taixé



G. Brasó and L. Leal-Taixé. "Learning a Neural Solver for Multiple Object Tracking", arXiv 2019

• Appearance and geometry encodings



Edge embeddings

• Appearance and geometry encodings



• Appearance and geometry encodings



- Instead of defining pairwise costs for edges and unary costs for nodes (classical setting), feature vectors encoding appearance and geometry cues are used
- Goal: propagate these embeddings across the entire graph in order to obtain new embeddings encoding high-order information among detections

Time-aware Message Passing



(b) Vanilla node update

All node embeddings are aggregated at once



(c) Time-aware node update

Aggregation of nodes is separated between past / future frames

Classifying edges

- After several iterations of message passing, each edge embedding contains high-order information about other detections
- We feed the embeddings to an MLP that predicts whether an edge is active/inactive



Obtaining final solutions

- After classifying edges, we get a prediction between 0 and 1 for each edge in the graph.
- Directly thresholding solutions could yield infeasible solutions (not frequent in practice)
 - Recall the Linear Program constraints!!
- A simple rounding scheme (greedy or LP based) can be used to obtain final trajectories

Some results

• In practice, around 99% of constraints are automatically satisfied, and rounding takes negligible time

• The overall method is fast (~12 fps) and achieves SOTA in the MOT Challenge by a significant margin

74



MOT evaluation
Evaluation metrics

- Compute a set of measures per frame
 - Perform matching between predictions and ground truth (we will use exactly the same Hungarian algorithm)
 - FP = False positives
 - FN = False negatives (missing detections)
 - IDsw: identity switches

Evaluation metrics

• How do we compute ID switches?



(a) An ID switch is counted because the ground truth track is assigned first to red, then to blue.

(b) You count both an ID switch (red and blue both assigned to the same ground truth), but also a fragmentation (Frag) because the ground truth coverage was cut.(c) Identity is preserved. If two trajectories overlap with a ground truth trajectory (within a threshold), the one that forces least ID switches is chosen (the red one).

Leal-Taixé et al.. MOTChallenge 2015: Towards a Benchmark for Multi-Target Tracking. arXiv:1504.01942 2015.

Evaluation metrics

- Compute a set of measures per frame
 - Perform matching between predictions and ground truth (we will use exactly the same Hungarian algorithm)
 - FP = False positives
 - FN = False negatives (missing detections)
 - IDsw: identity switches

Multi-object
tracking accuracy MOTA =
$$1 - \frac{\sum_{t} (FN_t + FP_t + IDSW_t)}{\sum_{t} GT_t}$$
, Ground truth

Datasets

- MOTChallenge: <u>www.motchallenge.net</u> (people)
 - Several challenges from less to more crowded



- KITTI benchmark: http://www.cvlibs.net/datasets/kitti/ (vehicles)
- UA-Detrac: <u>http://detrac-db.rit.albany.edu</u> (vehicles)

Next lectures

• Trajectory prediction

• Next lecture on December 20th.