

Computer Vision II

Speed is Important

Tracking: estimate the number & state of objects in a region of interest

↓

- 1: single target tracking
- 0 or 1: detection & tracking
- N: Multi-target detection & tracking

[x, y]: location
 [x, y, dx, dy]: location, velocity
 [x, y, appearance params]: humans..

- Articulated Tracking: estimate the motion of objects with multiple, coordinate parts
 Ex: face expression reconstruction, human pose estimation
- Active Tracking: moving sensor in response to motion of target need to be real time

+ Elements of Tracking:

- Detection: find object(s) of interest
- Association: which observations come from same object
- Prediction: predict future motion based on the observed motion pattern
 - use this prediction to improve detection & association

+ Visual Odometry: - track camera motion

- also involves a data association problem
- is prone to drift due to local views

- +)Visual SLAM (Simultaneous Localization & Mapping)
- Global & local optimization methods

+)Lecture Overview:

⊗ Exercise:

- 1) Simple Background Models
Statistical Background Models
Generalized Lukas - Kanade Tracking
- 2) Kalman Filter, Extended ICF Tracking
Particle Filter Tracking

Single-Object Tracking

↳ Background Modelling

- General algorithm for all kind of objects
- Real time \Rightarrow need to be simple calculation
- Assume: static cameras, moving objects are important

↳ Simple Background Models:

+ Simple Background Subtraction: $I(t) - B \rightarrow \text{abs} \rightarrow \text{threshold } \lambda \rightarrow M(t)$

- Problems:
 - objects enter scene, then stop \rightarrow still detected
 - assumed-static object starts moving \rightarrow 2 objects
 - illumination
 - global threshold is often suboptimal

+ Single Frame Differencing: $I(t) - B(t-1) \rightarrow \text{abs} \rightarrow T(\lambda) \rightarrow M(t)$

- (+) quick to adapt to change of lightning, camera motion
- (-) few pixels of object are labeled when whole object body
object moving away / toward camera has similar color..

$\Rightarrow +$ Temporal Scale: $D(N) = \|I(t) - I(t+N)\|$

$\Rightarrow +$ Three-Frame Differencing: $D(+N) \xrightarrow{\text{AND}} D(-N) \xrightarrow{\text{AND}} \text{object}$

+ Adaptive Background Subtraction:

$$I(t) \xrightarrow{\oplus} \text{abs} \xrightarrow{\downarrow} M(t)$$

$B(t-1)$
 \uparrow
 $B(t)$
 \uparrow
 \uparrow
 $\alpha \cdot I(t) + (1-\alpha)B(t-1)$

- trails of pixels
- ghosting effects
- center of large object fades

Comments

Simple background models are very heuristic, hard to find good parameters
optimal temporal scale depends on object (size & speed)
global threshold is often suboptimal for parts of the image

⇒ Need better founded approach

⇒ 2) Statistical Background Models:

+ Single Gaussian: $\hat{\mu} = \frac{1}{N} \sum x_n$; $\hat{\sigma}^2 = \frac{1}{N-1} \sum (x_n - \hat{\mu})^2$

- Online Adaptation: μ & σ for each pixel

Running Estimate:

$$\begin{aligned}\hat{\mu}^{(t+1)} &= \hat{\mu}^{(t)} + \frac{1}{N} x^{(t+1)} - \frac{1}{N} x^{(t+1-T)} && \text{add new one} \\ (\hat{\sigma}^2)^{(t+1)} &= (\hat{\sigma}^2)^{(t)} + \frac{1}{N-1} (x^{(t+1)} - \hat{\mu}^{(t+1)})^2 \\ &\quad - \frac{1}{N-1} (x^{(t+1-T)} - \hat{\mu}^{(t+1-T)})^2 && \text{remove oldest one}\end{aligned}$$

Exponential Moving Average Filter:

since new data are usually more valuable

$$\begin{aligned}\hat{\mu}^{(t+1)} &= (1-\alpha) \hat{\mu}^{(t)} + \alpha \cdot x^{(t+1)} \\ (\hat{\sigma}^2)^{(t+1)} &= (1-\alpha) (\hat{\sigma}^2)^{(t)} + \alpha (x^{(t+1)} - \hat{\mu}^{(t+1)})^2\end{aligned}$$

α as fixed learning rate

- Problems: [A single Gaussian is clearly insufficient
Adaptation speed vs Sensitivity]

+ MoG (Mixture of Gaussians):

- EM Algorithm: $p(x) = \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)$

+ Stauffer - Grimson: K Gaussians

$$p(x) = \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \text{ where } \Sigma_k = \sigma_k^2 \cdot I$$

- Check every new pixel value, within $2,5 \sigma_k$ of $\mu_k \Rightarrow$ matched

- If matched \Rightarrow adapt the corresponding component

else, replace least probable component with new value as μ_k
and an initial high σ_k and small π_k

- Order components by π_k/σ_k , select B best components as background
- $$B = \underset{b}{\operatorname{argmin}} \left(\sum_{k=1}^b \frac{\pi_k}{\sigma_k} > \epsilon \right)$$

- Online adapting for matching components:

$M_{k,t} = 1$ iff component k matched, else 0

$$\pi_k^{(t+1)} = (1-\alpha)\pi_k^{(t)} + \alpha \cdot M_{k,t}$$

$$\mu_k^{(t+1)} = (1-\rho)\mu_k^{(t)} + \rho \cdot x^{(t+1)}$$

$$\Sigma_k^{(t+1)} = (1-\rho)\Sigma_k^{(t)} + \rho(\mu^{(t+1)} - \mu^{(t+1)})(x^{(t+1)} - \mu^{(t+1)})^T$$

where $\rho = \alpha \cdot \mathcal{N}(x_t | \mu_k, \Sigma_k)$ update weighted by component likelihood

- # Comments MoG.
- static foreground objects can be integrated
 - ordering by π_k/σ_k : components have more evidence π_k smaller variance σ_k
 - can select only one component if unimodal
 - not flexible with dynamic areas (wind + trees, rippling water..)

+ Kernel background Modelling:

$$p(x^{(t)}) = \frac{1}{N} \sum_{i=1}^N \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{1}{2} \frac{(x_j^{(t)} - x_j^{(i)})^2}{\sigma_j^2}}$$

current pixel value $x^{(t)}$ with N other values at other time $x^{(i)} = \frac{1}{N} \sum_{i=1}^N K(x^{(t)} - x^{(i)})$ with K as Gaussian $\mathcal{N}(0, \Sigma)$ $\Sigma = \text{diag}(\sigma_j)$

\Rightarrow foreground if $p(x^{(t)}) < \epsilon$

can be speed up with look-up table..

can be very flexible, deal with large amounts of background motion

+ Update:

- FIFO: discard oldest
- Selective update: uncontaminated or go to deadlock situation
- Blind update: no deadlock, but more false negatives

\Rightarrow Solution: combine Short-term model with Long-term model
most recent window (Selective update) (Blind update) larger time window

Summary: Background Modelling: fast & simple, if applicable, use it

Stauffer - Grimson model } perform well
 Kernel background model } used extensively

④ Template based tracking:

L3/

1, Lucas - Kanade Optical Flow: Try to measure motion of each pixel in the scene

→ Key assumptions:

- Brightness constancy
- Small motion
- Spatial coherence (points move like their neighbors)

- Brightness constancy:

$$f(x, y, t) = f(x + dx, y + dy, t + dt)$$

⇒ Taylor expansion for approximation

$$f(x, y, t) = f(x, y, t) + \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy + \frac{\partial f}{\partial t} dt$$

$$\Leftrightarrow f_x \cdot dx + f_y \cdot dy + f_t \cdot dt = 0$$

$$\Leftrightarrow \boxed{f_x \cdot u + f_y \cdot v + f_t = 0}$$

$$\Rightarrow v = -\frac{f_x}{f_y} u - \frac{f_t}{f_y} \quad \text{equation of a line}$$

⇒ The Aperture problem

(The Barber Pole Illusion)

can't detect velocity that perpendicular to the gradient

→ Horn & Schunk:

$$\text{minimize } \iint \{ (f_x u + f_y v + f_t)^2 + \lambda (u_x^2 + u_y^2 + v_x^2 + v_y^2) \} dx dy$$

brightness constancy

smoothness constraint

$$\Leftrightarrow \begin{cases} (f_x u + f_y v + f_t) f_x + \lambda (\Delta^2 u) = 0 \\ (f_x u + f_y v + f_t) f_y + \lambda (\Delta^2 v) = 0 \end{cases}$$

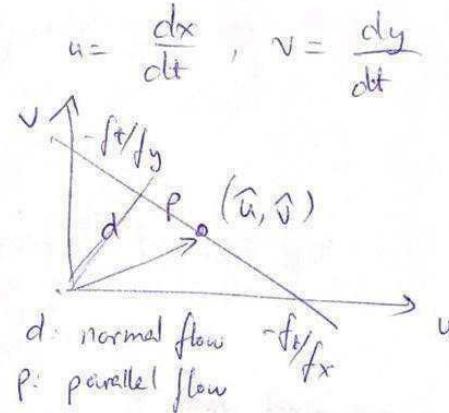
$$\text{discretize: } \begin{cases} (f_x u + f_y v + f_t) f_x + \lambda (u - u_{av}) = 0 \\ (f_x u + f_y v + f_t) f_y + \lambda (v - v_{av}) = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} u = u_{av} - f_x \frac{P}{D} \\ v = v_{av} - f_y \frac{P}{D} \end{cases}$$

$$P = f_x u_{av} + f_y v_{av} + f_t$$

$$D = \lambda + f_x^2 + f_y^2$$

Iterative algorithm till u, v converge



7) Lukas - Kanade (Least-squares)

- Pretend neighbor pixels have the same (u, v)

→ Use 5×5 window $\Rightarrow 25$ equations of brightness constancy:

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = -\begin{bmatrix} I_t(p_1) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

$$A \cdot d = b$$

Since A is not square, can not invert normally

$$\Rightarrow A^T A \cdot d = A^T b$$

$$\Rightarrow d = \underbrace{(A^T A)^{-1}}_{\text{pseudo-inverse}} A^T b$$

$$\Leftrightarrow \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = -\begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

- Remember that we assume all pixels in window have same (u, v)
actually, they don't, thus $f_{xi}u + f_{yi}v + f_t \neq 0$

$$\Rightarrow \text{Minimize } \min_t \sum_i (f_{xi}u + f_{yi}v + f_t)^2$$

$$\Leftrightarrow \begin{cases} \sum (f_{xi}u + f_{yi}v + f_t)f_{xi} = 0 \\ \sum (f_{xi}u + f_{yi}v + f_t)f_{yi} = 0 \end{cases} \Leftrightarrow \begin{cases} \sum f_{xi}^2 u + \sum f_{xi} f_{yi} v = -\sum f_{xi} f_t \\ \sum f_{xi} f_{yi} u + \sum f_{yi}^2 v = -\sum f_{yi} f_t \end{cases}$$

7) Problem: Errors from non-linear

→ Iterative LK Refinement:

Forward warping

H

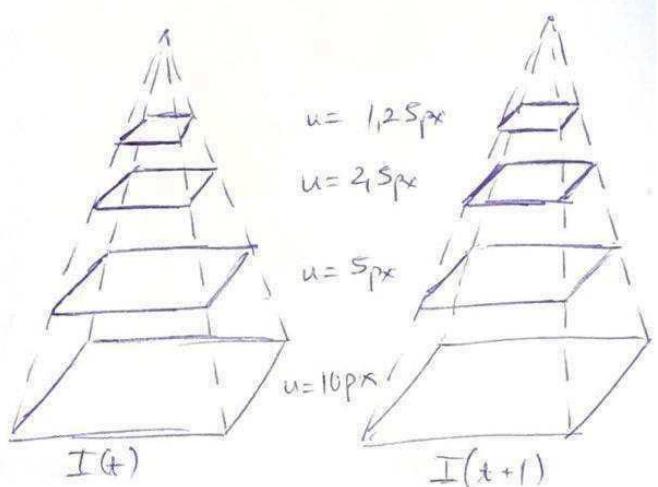
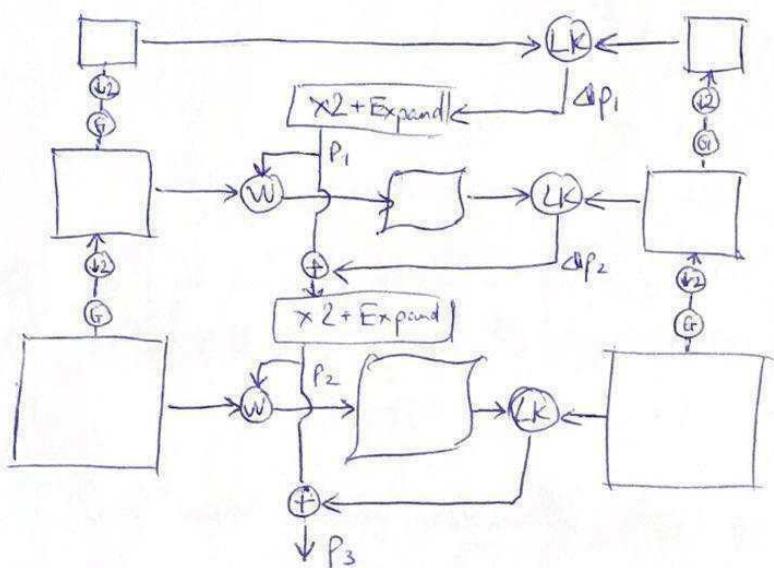
Backward warping
(with bilinear interpolation)

- Estimate velocity at each pixel by solving LK equations.
- Warp $I_t \rightarrow I_{t+1}$ used (u, v) , but we only get I'_{t+1}
- $I'_{t+1} \neq I_{t+1}$
- Repeat (estimate (u, v) from I'_{t+1} and $I_{t+1} \rightarrow$ warp $\Rightarrow \dots$) until convergence

+ Problem 2: Large Motions (more than 1-2 pixels)

→ Coarse-to-Fine Optical Flow

→ Hierarchical UK:



- 1) Image Pyramid with Gaussian blur and down sampling
- 2) UK on level k
- 3) $\times 2$ and Upsampling
- 4) Warp image at level $k-1$
- 5) Refinement → repeat

2) Feature Tracking:

Sparse UK · Shi-Tomasi Feature Tracker (KLT)

Only find optical flow of key feature points

For points in middle of object, we can't estimate (u, v) accurately

In terms of math, A has small eigenvalues

- Comparing to the first frame helps minimize the drift

KLT: first step of SfM / SLAM

LK Template Tracking

L9)

+ Idea: In KLT, we track corner-like feature
(small patches: $5 \times 5 \text{px}$)

compute optical flow

Now, LK Template Tracking, use larger window

$$+ E(u,v) = \sum_x [I(x+u, y+v) - T(x,y)]^2$$

- Problematic Assumption: a constant flow (pure translation)

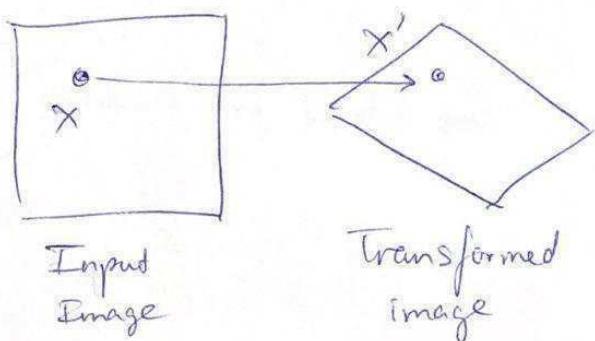
(-) In practice, objects might rotate -

→ Use "warping" function W with parameters p

$$E(p) = \sum_x [I(W([x,y]; p)) - T(x,y)]^2$$

+ Geometric Image Warping:

$$x' = W(x; p)$$



- Translation:

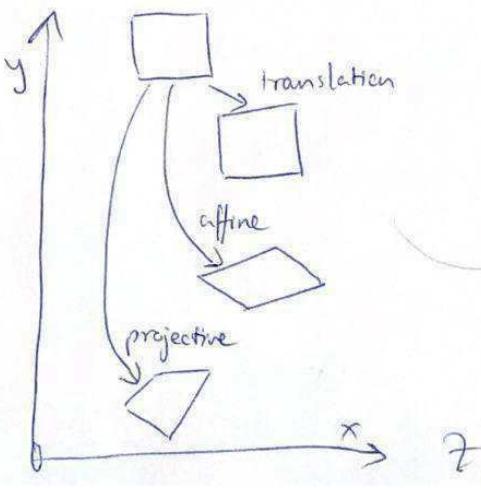
$$\text{2 unknowns } W([x,y]; p) = \begin{bmatrix} x + p_1 \\ y + p_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & p_1 \\ 0 & 1 & p_2 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

- Affine:

$$\text{6 unknowns } W([x,y]; p) = \begin{bmatrix} x + p_1x + p_3y + p_5 \\ y + p_2x + p_4y + p_6 \end{bmatrix} = \begin{bmatrix} 1 + p_1 & p_3 & p_5 \\ p_2 & 1 + p_4 & p_6 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

- Perspective

$$\text{8 unknowns } W([x,y]; p) = \frac{1}{p_7x + p_8y + 1} \begin{bmatrix} x + p_1x + p_3y + p_5 \\ y + p_2x + p_4y + p_6 \end{bmatrix}$$



- With KLT, we kind of already consider warping, but only translation warping, thus $\rho = (u, v)$

\Rightarrow Now, we consider affine warping, ρ has 6 dof

$$I(W(x; \rho + \Delta\rho)) \approx I(W(x; \rho)) + \nabla I \frac{\partial W}{\partial \rho} \Delta\rho$$

- As $\frac{\partial W}{\partial \rho} = \begin{bmatrix} \frac{\partial W_x}{\partial \rho_1} & \frac{\partial W_x}{\partial \rho_2} & \dots & \frac{\partial W_x}{\partial \rho_n} \\ \frac{\partial W_y}{\partial \rho_1} & \frac{\partial W_y}{\partial \rho_2} & \dots & \frac{\partial W_y}{\partial \rho_n} \end{bmatrix}$

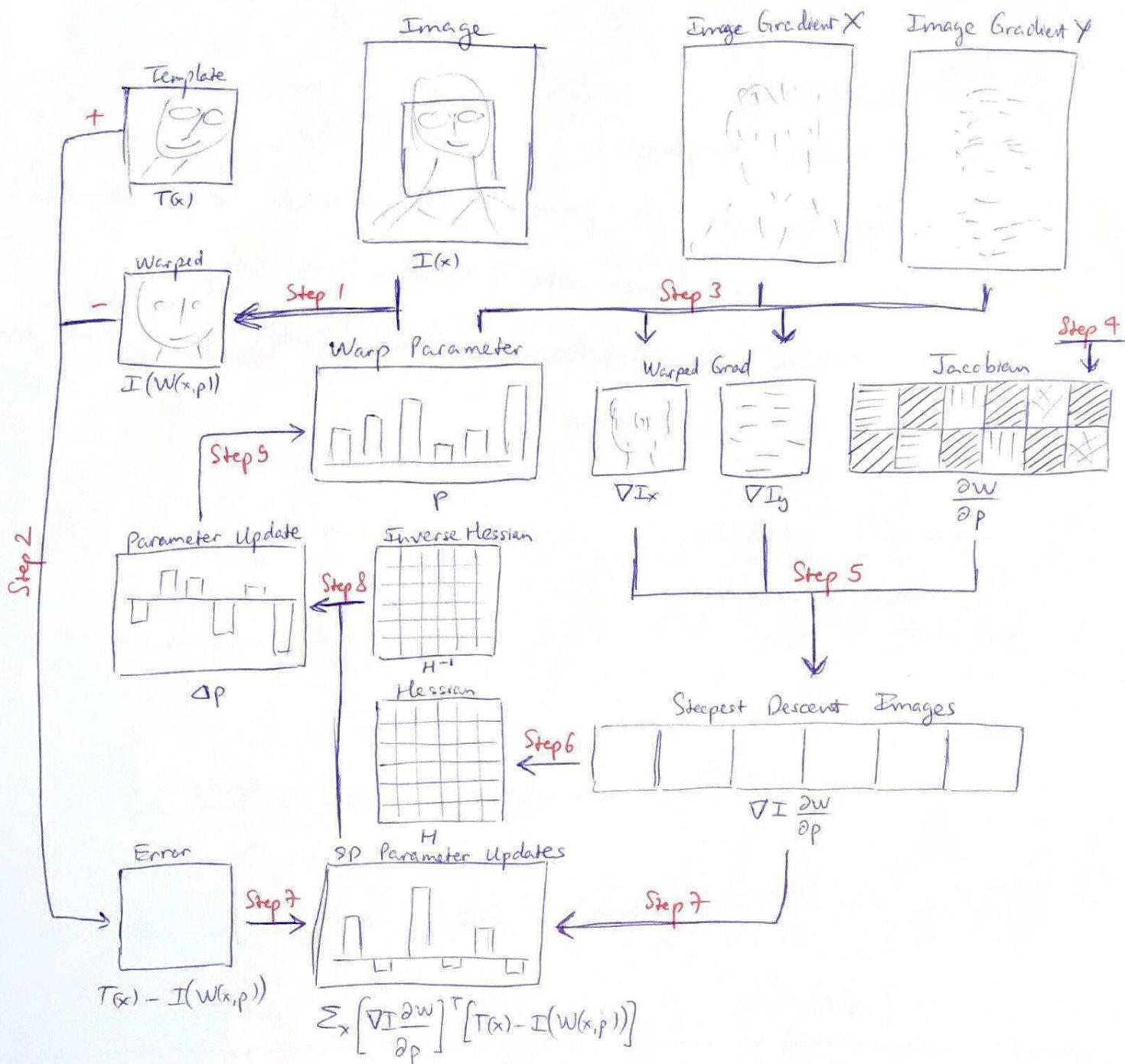
And that: $W([x, y]; \rho) = \begin{bmatrix} 1 + \rho_1 & \rho_3 & \rho_5 \\ \rho_2 & 1 + \rho_4 & \rho_6 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x + \rho_1 x + \rho_3 y + \rho_5 \\ y + \rho_2 x + \rho_4 y + \rho_6 \end{bmatrix}$
for affine if

$$\Rightarrow \frac{\partial W}{\partial \rho} = \begin{bmatrix} x & 0 & y & 0 & 1 & 0 \\ 0 & x & 0 & y & 0 & 1 \end{bmatrix}$$

$$\Rightarrow \text{Problem: } \arg \min_{\Delta\rho} \sum_x \left[I(W(x; \rho)) + \nabla I \frac{\partial W}{\partial \rho} \Delta\rho - T(x) \right]^2$$

$$\Rightarrow \text{Closed form solution: } \boxed{\Delta\rho = H^{-1} \sum_x \left[\nabla I \frac{\partial W}{\partial \rho} \right]^T [T(x) - I(W(x; \rho))]}$$

$$\text{Hessian } H = \sum_x \left[\nabla I \frac{\partial W}{\partial \rho} \right]^T \left[\nabla I \frac{\partial W}{\partial \rho} \right]$$



②, Discussion:

(+) Fast & Simple algorithm

Applicable to
Optical flow estimation
stereo disparity estimation
parametric motion tracking

(-) Prone to local minima

Relatively small movement

→ Good init necessary

* Note:

- When tracking templates..

Typically use the previous frame's result as initialization

Higher the frame rate, smaller the warp

⇒ Tracking becomes easier (& faster) with higher frame rates

Tracking by Online Classification

L5)

Idea:

- As object moves, new objects | scenes will appear
background

- Object Detector \Rightarrow Object Tracker

Fixed training set
General object detector

On-line update
Object vs Background

, Adaboost in Viola - Jones Detector

* From offline \Rightarrow Online Boosting

Offline	Online
<u>Given:</u> <ul style="list-style-type: none"> - Set of labeled samples $X = \{ \langle x_i, y_i \rangle, \dots, \langle x_L, y_L \rangle \mid y_i = \pm 1 \}$ - Weight distribution of samples $D_0 = 1/L$ 	<ul style="list-style-type: none"> - One labeled sample $\langle x, y \rangle \mid y = \pm 1$ - Strong classifiers to update - Initial importance $\lambda = 1$ of sample
<u>Loop</u> : for $n=1$ to N <ul style="list-style-type: none"> - Train a new weak classifier $h_n^{\text{weak}}(x) = L(X, D_{n-1})$ - calculate error e_n - calculate weight for classifier $\alpha_n = f(e_n)$ - update weight dist. of samples D_n 	<u>Loop</u> : <ul style="list-style-type: none"> - Update weak classifiers using samples and their import $h_n^{\text{weak}}(x) = L(h_n^{\text{weak}}, \langle x, y \rangle, \lambda)$ - Update error estimation \hat{e}_n - Update weight of classifier: $\alpha_n = f(\hat{e}_n)$ - Update importance weights of samples: λ
$h^{\text{strong}}(x) = \text{sign} \left(\sum_{n=1}^N \alpha_n h_n^{\text{weak}}(x) \right)$	$h^{\text{strong}}(x) = \text{sign} \left(\sum_{n=1}^N \alpha_n h_n^{\text{weak}}(x) \right)$

- Compare:

Offline Boosting	Online Boosting
At each iteration Pass all samples Add a new classifier	Pass 1 samples Update all classifier
- Given same training set,
Online Boosting will converge to same weak classifiers as
Offline Boosting in limit of $N \rightarrow \infty$ iterations (Oza, 2001)

+ Features:

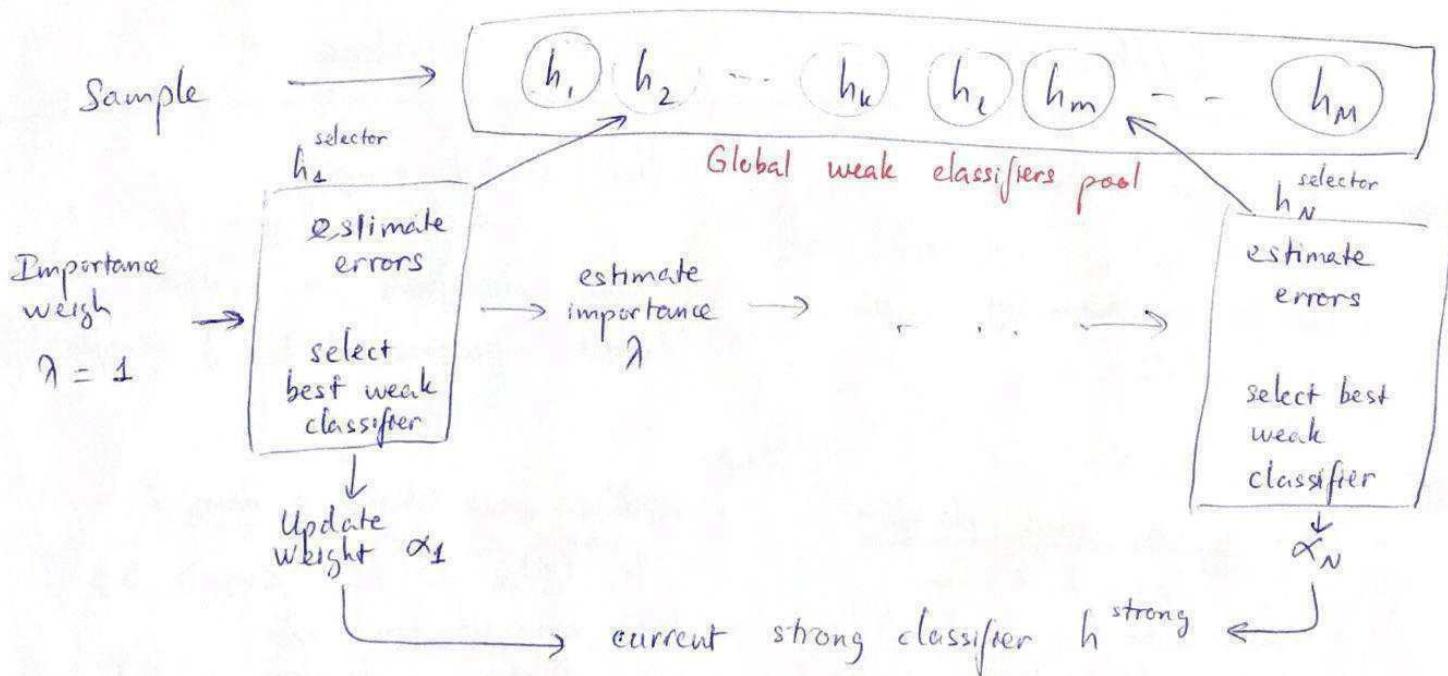
- Haar-like Wavelets
- HOG
- Locally binary pattern (LBP)

Fast computation using efficient data structures: image integral histograms

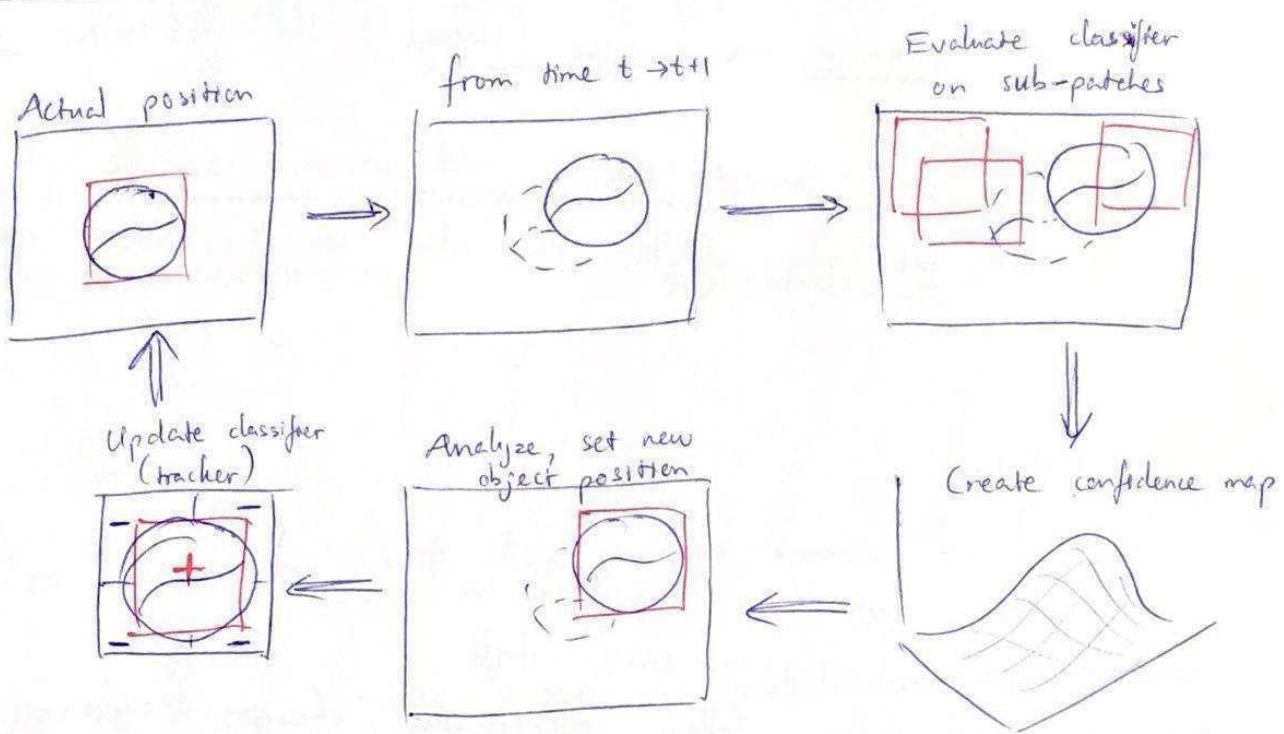
+ Selector

choose 1 feature for each weak classifier
(because there are many feature choices)

⊗ Direct Feature Selection



⊗ Tracking by Online Classification



- Add partial overlap samples as negative samples
⇒ to penalize if misalign object center
- Online feature exchange:
Combine difference feature type Haar, LBP ..

⊗ Summary

- Performs Boosting on Selectors, not weak classifiers
(what features go with what classifiers)
- Use shared feature pool across classifiers
(all N classifiers have same choice from M feature)
not each classifier has different M feature choices

② Problem: Drift.

- If the change is slow enough, the classifier adapt to new objects
- Because we might slowly add wrong sample
- It will not only drift, but also stay confident about it

⇒ Solution:

1) Match against initialization

(+) robust to catch drift

(-) can't follow appearance changes (lightning, rotation ...)

2) Semi supervised learning

Object Detector



Object Tracker

Fixed Prior for updating
an adaptive on-line classifier
(should still be in same class)

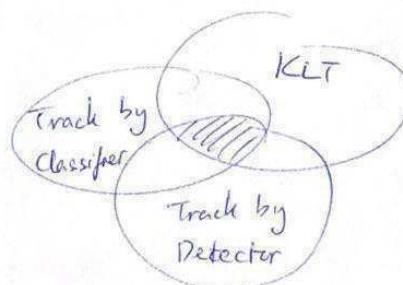
3) Using additional cues:

+/-) Tracking - Learning - Detection (TLD)

- Combine with informations from

[KLT

Tracking by Detection ..



Tracking by Detection

L6)

currently the basis for Multi Object tracking
surveillance

→ Idea:

Detection \Rightarrow Data Association \Rightarrow Prediction

- Apply generic object detector
 - Use appearance models of those objects
 - Even for back-ground segmentations for better appearance model
 - Predict position of objects in next time step
- Objects

→ SVM based Detectors:

- Sliding - window approach
- HoG, SVM, Image Pyramid, Non-Maximum Suppression

→ Deformable Part-based Model (DPM)

- Idea: Some objects are vary in form, pose ...

- Response of i -th part filter, in l -th level:

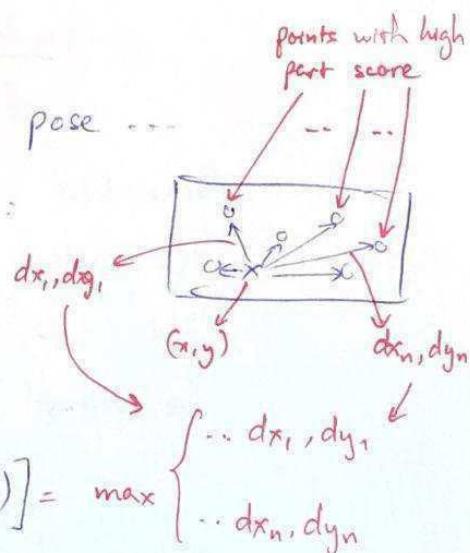
$$\text{From part score} \rightarrow R_{i,l}(x,y) = F_i \cdot \phi(H_l(x,y,l))$$

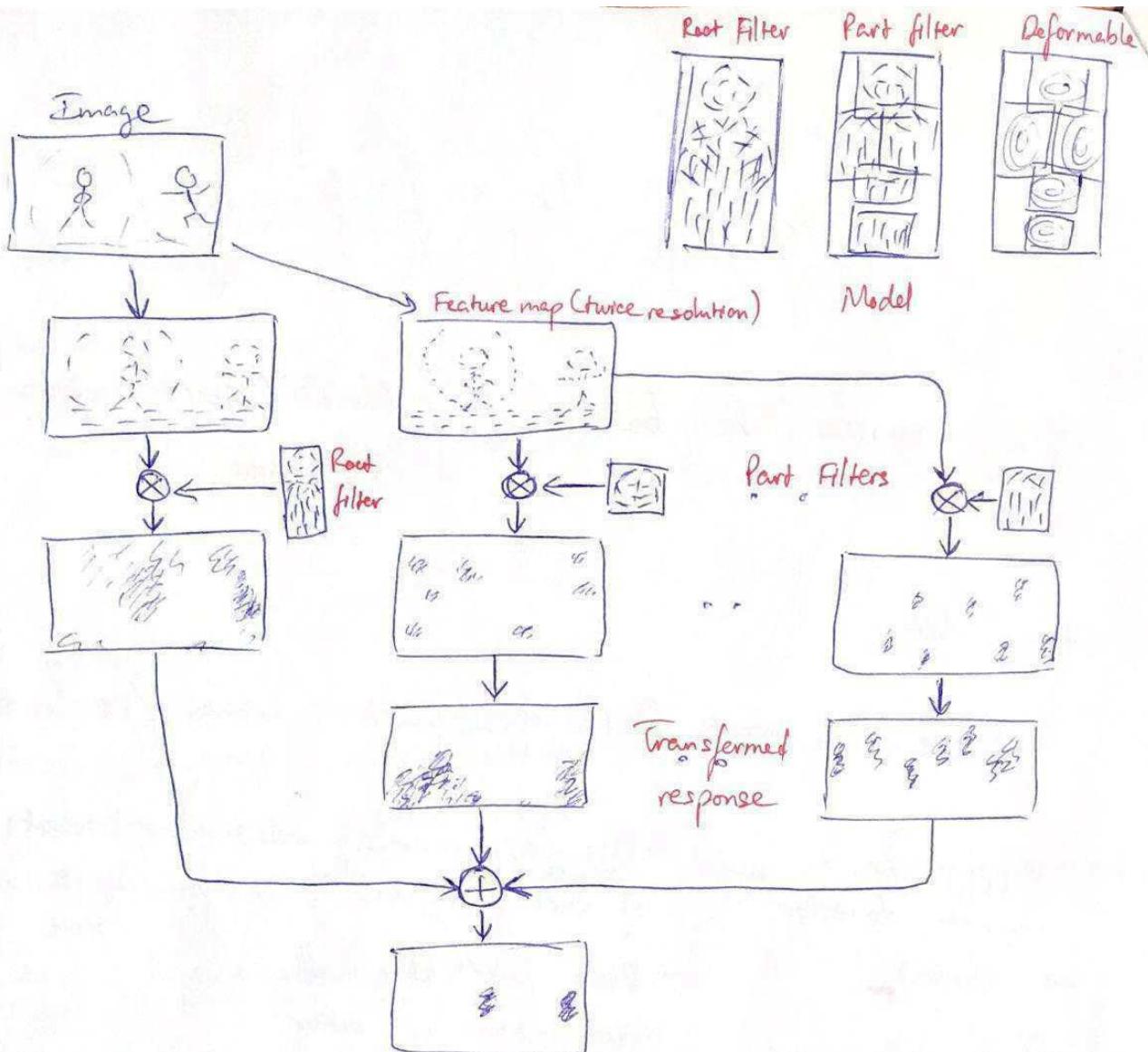
to root score

② Transformed response

$$D_{i,l}(x,y) = \max_{dx,dy} [R_{i,l}(x+dx, y+dy) - d_i \cdot \phi_d(dx,dy)] = \max \begin{cases} \dots dx_1, dy_1 \\ \dots dx_n, dy_n \end{cases}$$

dx, dy (the displacement between root & part)
has to make sense for the deformable loss to be small





2) AdaBoost-based Detectors

- Viola - Jones: Integral Images + Cascading Classifiers

⇒ Integral Channel Features

- Idea: Generalization of Haar Wavelet
Combine other feature channels (HOG, color, texture)
 6 orientation bins + Grayscale magnitude + 3 LUV color channels
- Generalize also block computation
 - 1st order features: sum of pixels in rectangular region
 - 2nd _____: Haar-like difference over blocks
 - Generalized Haar: sum with weights

+ Very Fast / Roerel

- $\begin{cases} 1 \text{ model} \\ 50 \text{ image scales} \end{cases} \Rightarrow \begin{cases} 50 \text{ models} \\ 1 \text{ image scale} \end{cases} \Rightarrow \begin{cases} 5 \text{ models} \\ 1 \text{ image scale} \end{cases}$
Then by using feature interpolation
 \Rightarrow Transfer test time computation to train time
- Ensemble of short trees, learned by AdaBoost

3) CNN-based Detectors

- Faster RCNN
- Mask RCNN
- YOLO, SSD

Bayesian Filtering

L7)

1) Tracking with Dynamics

+ let state be X , measurement by Y

⇒ Prediction: past measurements → current state

$$P(X_t | Y_0 = y_0, \dots, Y_{t-1} = y_{t-1})$$

Correction: Updated state from prediction & new measurement

$$P(X_t | Y_0 = y_0, \dots, Y_{t-1} = y_{t-1}, Y_t = y_t)$$

+ Simplify assumptions:

- Only immediate past matters:

$$P(X_t | X_0, \dots, X_{t-1}) = P(X_t | X_{t-1})$$

- Measurements depend only on current state

$$P(Y_t | X_0, \dots, X_t) = P(Y_t | X_t)$$

$$P(A) = \int P(A, B) dB$$

$$\int P(X_t, X_{t-1} | y_0, \dots, y_{t-1}) dX_t$$

⇒ Formulation:

$$\text{Prediction } P(X_t | y_0, \dots, y_{t-1}) = \int P(X_t | X_{t-1}) P(X_{t-1} | y_0, \dots, y_{t-1}) dX_{t-1} \stackrel{\substack{\text{Dynamics model} \\ \text{corrected estimate}}}{=} \int P(X_t | X_{t-1}) P(X_t | y_0, \dots, y_{t-1}) dX_t$$

$$\text{Correction } P(X_t | y_0, \dots, y_t) = \frac{P(y_t | X_t, y_0, \dots, y_{t-1}) P(X_t | y_0, \dots, y_{t-1})}{P(y_t | y_0, \dots, y_{t-1})} \stackrel{P(A|B) = \frac{P(B|A)P(A)}{P(B)}}{=} \frac{P(y_t | X_t) P(X_t | y_0, \dots, y_{t-1})}{\int P(y_t | X_t) P(X_t | y_0, \dots, y_{t-1}) dX_t}$$

measurement model

$$= \frac{P(y_t | X_t) \cdot P(X_t | y_0, \dots, y_{t-1}) - \text{predicted state}}{\int P(y_t | X_t) P(X_t | y_0, \dots, y_{t-1}) dX_t}$$

$$= \eta \cdot P(y_t | X_t) \cdot P(X_t | y_0, \dots, y_{t-1})$$

2) Linear Dynamics Models.

- Dynamics Model $x_t \sim N(D_t x_{t-1}, \Sigma_{dt})$

- Observation Model $y_t \sim N(M_t x_t, \Sigma_{mt})$

- Ex: Constant Velocity (10 points)

$$x_t = \begin{bmatrix} p_t \\ v_t \\ a_t \end{bmatrix} = D_t x_{t-1} + \Sigma_{dt} = \begin{bmatrix} 1 & \Delta t & \frac{1}{2} \Delta t^2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix} \cdot x_{t-1} ; \quad y_t = M_t x_t = \begin{bmatrix} 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} p_t \\ v_t \\ a_t \end{bmatrix} = p_t$$

Ex Constant Acceleration (10 points)

$$x_t = \begin{bmatrix} p_t \\ v_t \\ a_t \end{bmatrix} = D_t x_{t-1} = \begin{bmatrix} 1 & \Delta t & \frac{1}{2} \Delta t^2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} p_{t-1} \\ v_{t-1} \\ a_{t-1} \end{bmatrix} ; \quad y_t = M_t x_t = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} p_t \\ v_t \\ a_t \end{bmatrix} = p_t$$

3) Kalman Filter only linear transition model

$$x_t \sim N(D_t x_{t-1}, \Sigma_{dt})$$

Prediction

$$\bar{x}_t = D_t \cdot x_{t-1}^+$$

$$\bar{\Sigma}_t = D_t \cdot \Sigma_{t-1}^+ \cdot D_t^T + \Sigma_{dt}$$

$$y_t \sim N(M_t x_t, \Sigma_{mt})$$

Correction

Kalman gain

$$K_t = \bar{\Sigma}_t^- M_t^T (M_t \bar{\Sigma}_t^- M_t^T + \Sigma_{mt})^{-1}$$

$$x_t^+ = \bar{x}_t^- + K_t (\underbrace{y_t - M_t \bar{x}_t^-}_{\text{residual}})$$

$$\bar{\Sigma}_t^+ = (I - K_t M_t) \bar{\Sigma}_t^-$$

Cons: restricted to linear model class of motions
 \Rightarrow EKF

4) EKF (Extended Kalman filter)

L8, +) Non linear model:

$$x_t = g(x_{t-1}) + \varepsilon_t$$

$$y_t = h(x_t) + \delta_t$$

+) Prediction step:

$$\bar{x}_t = g(\bar{x}_{t-1}^+)$$

$$\bar{\Sigma}_t^- = G_t \bar{\Sigma}_{t-1}^+ G_t^\top + \bar{\Sigma}_{dt}$$

+), Correction step:

$$K_t = \bar{\Sigma}_t^- H_t^\top (H_t \bar{\Sigma}_t^- H_t^\top + \bar{\Sigma}_{mt})^{-1}$$

$$x_t^+ = \bar{x}_t^- + K_t (y_t - h(\bar{x}_t^-))$$

$$\Sigma_t^+ = (I - K_t H_t) \bar{\Sigma}_t^-$$

Kalman
gain

The Jacobians

$$G_t = \left. \frac{\partial g(x)}{\partial x} \right|_{x=x_{t-1}^+}$$

$$H_t = \left. \frac{\partial h(x)}{\partial x} \right|_{x=\bar{x}_t^-}$$

- Note:
- EKF essentially linearizes around current estimation
 - Generally not an optimal estimator
 - If initial estimate is wrong \Rightarrow the filter will quickly diverge
 - Still, is the de-facto standard in many applications

L9, Particle Filter

To represent posterior pdf by randomly chosen weighted samples
particles

+), Monte Carlo Sampling:

$$\mathbb{E}[f] = \int f(z) p(z) dz$$

Draw L samples $z^{(l)}$ from $p(z)$

\Rightarrow Can approximate $\mathbb{E}[\hat{f}] = \mathbb{E}[f]$

$$\text{with } \hat{f} = \frac{1}{L} \sum_l f(z^{(l)})$$

→ Monte Carlo integration: $F = \int_D f(x) dx$

$$= \int_D \frac{f(x)}{p(x)} p(x) dx = \mathbb{E}_{x \sim p} \left[\frac{f(x)}{p(x)} \right]$$

Again, choose N i.i.d samples from p : x_1, \dots, x_N

$$\Rightarrow F_N = \frac{1}{N} \sum_i^N \frac{f(x_i)}{p(x_i)} \quad \text{as } \lim_{N \rightarrow \infty} F_N = F$$

→ Idea: for Importance Sampling:

$$\begin{aligned} \mathbb{E}[f] &= \int f(z) p(z) dz = \int f(z) \frac{p(z)}{q(z)} q(z) dz \\ &\approx \frac{1}{L} \sum_e^L \underbrace{\frac{p(z^e)}{q(z^e)} f(z^e)}_{\text{Importance weight}} \quad \text{with } z^e \sim q(z) \end{aligned}$$

Approximate f with a given distribution $p(z)$

by finite samples from a new distribution $q(z)$

+ Sequential Importance Sampling Algorithm (SIS)

$$\text{function } \left[\{x_t^i, w_t^i\}_{i=1}^N \right] = \text{SIS} \left[\{x_{t-1}^i, w_{t-1}^i\}_{i=1}^N, y_t \right]$$

$$\eta = 0 \quad (\text{normalize factor})$$

for $i = 1:N$ (N particles)

$$x_t^i \sim q(x_t | x_{t-1}^i, y_t)$$

$$w_t^i = w_{t-1}^i \frac{p(y_t | x_t^i) p(x_t^i | x_{t-1}^i)}{q(x_t | x_{t-1}^i, y_t)}$$

$$\eta = \eta + w_t^i$$

sample from proposal distribution pdf

end

for $i = 1:N$

$$w_t^i = w_t^i / N \quad (\text{normalize the weights})$$

end

$$w_t^i = w_{t-1}^i \cdot p(y_t | x_t^i)$$

+ choose prior pdf $q(x_t | x_{t-1}^i, y_t)$ as $p(x_t | x_{t-1})$ (transitional prior)

→ SIS Algorithm with Transitional Prior

- Sampling step: $x_t^i = g(x_{t-1}^i) + \varepsilon_t^i$

+ The Degeneracy Phenomenon

After few iterations, most particles have negligible weights

\Rightarrow Waste computational effort to update these particles

\rightarrow Need to measure degeneracy : $N_{\text{eff}} = \frac{1}{\sum_{i=1}^N (\omega_t^i)^2}$

[If uniform: $N_{\text{eff}} = N$

Severe degeneracy: $N_{\text{eff}} = 1$

\Rightarrow Resampling according to prior weights

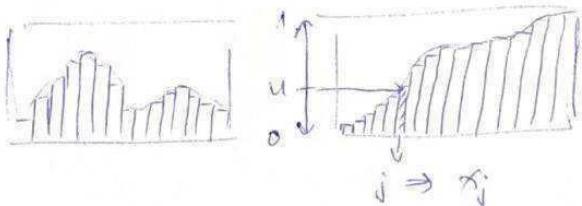
$$\left\{ x_t^i, w_t^i \right\}_{i=1}^N \rightarrow \left\{ x_t^i, \frac{1}{N} \right\}_{i=1}^N$$

+ Inverse Transform Sampling:

(CDF)

- Cumulative distribution function:

$$C(k) = \sum_{i=1}^k w_i / \sum_{i=1}^N w_i$$



Arulampalam's Resampling Algorithm

Given the
weights w
of the values x

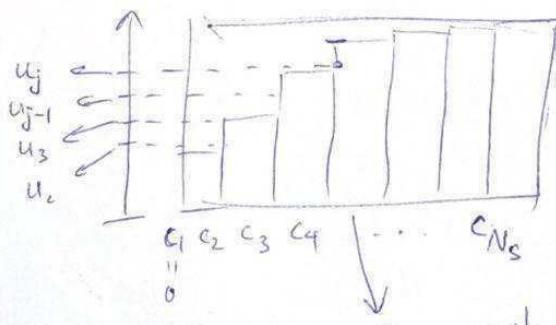
$$\left[\left\{ x_k^{j*}, w_k^j, i^j \right\}_{j=1}^{N_s} \right] = \text{RE SAMPLE} \left[\left\{ x_k^i, w_k^i \right\}_{i=1}^{N_s} \right]$$

How to
sample the
values?

Init CDF : $c_1 = 0$
for $i = 2: N_s$
| $c_i = c_{i-1} + w_k^i$

(Input weights are different to each other)

end
 $i=1$
 $u_1 \sim U[0, N_s^{-1}]$



for $j = 1: N_s$
| $u_j = u_1 + N_s^{-1}(j-1)$

while $u_j > c_i$
| $i++$

end

$x_k^{j*} = x_k^i$

$w_k^j = N_s^{-1}$

$i^j = i$

end

output weights are all equal

~~④~~ Generic Particle Filter:

$$\text{function} \left[\{x_t^i, w_t^i\}_{i=1}^N \right] = PF \left[\{x_{t-1}^i, w_{t-1}^i\}_{i=1}^N, y_t \right]$$

$$\left[\{x_t^i, w_t^i\}_{i=1}^N \right] = SIS \left[\{x_{t-1}^i, w_{t-1}^i\}_{i=1}^N, y_t \right]$$

Calculate N_{eff}

if $N_{eff} < N_{thr}$ (threshold)

$$\left[\{x_t^i, w_t^i\}_{i=1}^N \right] = RESAMPLE \left[\{x_t^i, w_t^i\}_{i=1}^N \right]$$

end

only resample
when needed

avoid drift

④ Another variant: Sampling Importance Resampling (SIR)

$$\text{function} [X_t] = SIR[X_{t-1}, Y_t]$$

$$\bar{X}_t = X_t = \emptyset$$

for $i = 1:N$

$$\left| \text{sample } x_t^i \sim p(x_t | x_{t-1}^i) \right.$$

$$\left| w_t^i = p(y_t | x_t^i) \right.$$

end

for $i = 1:N$

$\left| \text{draw } i \text{ with probability } \propto w_t^i \right.$

$\left| \text{add } x_t^i \text{ to } X_t \right.$

end

sampling every iteration

prone to drift

but reflect change of particles..

④ Application: can track head, body parts

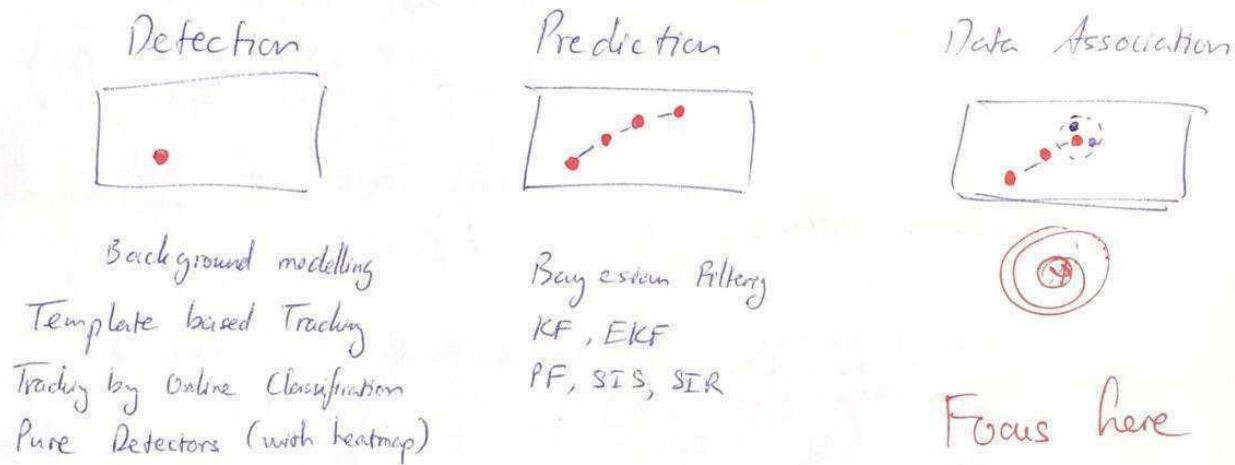
but as single object tracking problem

worst case complexity grows exponentially with dimensions

Multi Object Tracking

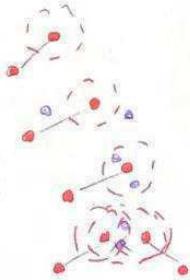
L16,

Idea: Elements of Tracking



⇒ Motion correspondence problem:

- +/- Ambiguities:
- disappear objects, or obscured
 - unexpected measurements
 - one prediction - many measurements
 - many predictions - one measurement



Simple Approaches

+/- Formalization:

- State vector: $\underline{x} = [x, y, v_x, v_y]^T$
- at time t : $\underline{x}^{(k)} = [x^{(k)}, y^{(k)}, v_x^{(k)}, v_y^{(k)}]^T$
- Observation at each time step: $\underline{y}^{(k)} = \{y_1^{(k)}, \dots, y_{M_k}^{(k)}\}$

→ We need to make association between:

tracks: $\{\underline{x}_1^{(k)}, \dots, \underline{x}_{N_k}^{(k)}\}$ and observation $\{y_1^{(k)}, \dots, y_{M_k}^{(k)}\}$

$\underline{x}_l^{(k)} = j$ iff $y_j^{(k)}$ is associated with $\underline{x}_l^{(k)}$

+ Gating: only consider measurements within a certain area around predicted location

+ Nearest-Neighbor filter: closest to the prediction \hat{x}_p

$$z_e^{(k)} = \arg\min_j (\hat{x}_{p,l}^{(k)} - y_j^{(k)})^T (\hat{x}_{p,l}^{(k)} - y_j^{(k)})$$

$$\text{Or better: } z_e^{(k)} = \arg\max_j N(y_j^{(k)}; \hat{x}_{p,l}^{(k)}, \Sigma_{p,l}^{(k)})$$

equivalent to: $z_e^{(k)} = \arg\min_j (\hat{x}_{p,l}^{(k)} - y_j^{(k)})^T \Sigma_{p,l}^{(k)-1} (\hat{x}_{p,l}^{(k)} - y_j^{(k)})$
 (Mahalanobis distance) - makes sense with Kalman filter
 not Euclidean distance

+ Extra notation:

- The innovation that measurement $y_j^{(k)}$ bring to track x_e at time k

$$v_{j,l}^{(k)} = (y_j^{(k)} - \hat{x}_{p,l}^{(k)})$$

$$\Rightarrow p(y_j^{(k)} | \hat{x}_e^{(k)}) \sim \exp\left\{-\frac{1}{2} v_{j,e}^{(k)T} \Sigma_{p,l}^{(k)-1} v_{j,e}^{(k)}\right\}$$

and the ellipsoidal gating / validation volume

$$V^{(k)}(\gamma) = \{y | (y - \hat{x}_{p,l}^{(k)})^T \Sigma_{p,l}^{(k)-1} (y - \hat{x}_{p,l}^{(k)}) \leq \gamma\}$$

⊕ Problems: hard assignment - finite chance that the association is incorrect
 (of nearest neighbor)

⇒ 2) Track-Splitting Filter: as soft/alternate assignments

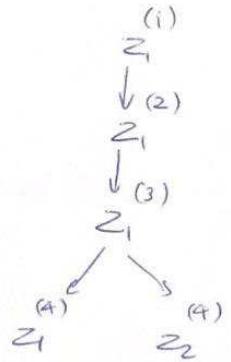
→ postpone decision until additional measurements have been gathered

+ Track likelihood:

Given track l , denote by $\theta_{k,l}$ the event that the sequence of assignments from time step $1 \rightarrow k$ belongs to same object.

$$Z_{k,l} = \{z_{i_1,l}^{(1)}, \dots, z_{i_k,l}^{(k)}\}$$

$$L(\theta_{k,l}) = \prod_{j=1}^k p(z_{i_j,l}^{(k)} | Z_{j-1,l}, \theta_{k,l})$$



- If assume Gaussian obs. likelihood

$$\Rightarrow L(\theta_{k,l}) = \prod_{j=1}^k \frac{1}{(2\pi)^{d/2} |\Sigma_l^{(j)}|^{1/2}} \exp \left[-\frac{1}{2} \sum_{j=1}^k v_{ij,l}^{(j)T} \Sigma_l^{(j)-1} v_{ij,l}^{(j)} \right]$$

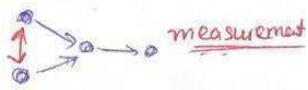
\Rightarrow The modified log-likelihood λ_l of track l

$$\begin{aligned}\lambda_l(k) &= -2 \log \left[\frac{L(\theta_{k,l})}{\sum_{j=1}^k (2\pi)^{d/2} |\Sigma_l^{(j)}|^{-1/2}} \right] \\ &= \sum_{j=1}^k v_{ij,l}^{(j)T} \Sigma_l^{(j)-1} v_{ij,l}^{(j)} \\ &= \lambda_l(k-1) + v_{ik,l}^{(k)T} \Sigma_l^{(k)-1} v_{ik,l}^{(k)}\end{aligned}$$

\Rightarrow Recursive calculation, sum of Mahalanobis distances
(the track with shortest "Mahalanobis" trajectories)

+ Pruning strategies:

- Deleting unlikely tracks: compare $\lambda_l(k)$ with threshold
- Use sliding window or exponential decay term to deal with long tracks (reduce the effect of very old observation that might overrule new observation)
- Merging track nodes: If state estimates are similar
 - both tracks validate identical subsequent measurements
 - rank the tracks based on $\lambda_l(k)$



L11) 3) Multi-Hypothesis Tracking (MHT)

- Target-oriented approaches

evaluate probabilities of

measurement belongs to established objects

$$P(Y|X)$$

- Measurement oriented approaches

established object | belongs to a measurement sequence
new target | give rise

$$P(X|Y)$$

MHT

7

+ Hypothesis generation:

- Set of hypotheses at time k :

$$\underline{\Omega}^{(k)} = \{\underline{\Omega}_j^{(k)}\}$$

- Latest set of measurements

$$\underline{Y}^{(k)} = \{y_1^{(k)}, y_2^{(k)}, \dots, y_{Mk}^{(k)}\}$$

- Feasible associations:

- [] A continuation of previous track
- [] A false alarm
- [] A new target

\Rightarrow Hypothesis matrix of feasible associations:

$$\Theta = \begin{bmatrix} x_1 & x_2 & \dots & x_{fa} & x_{nt} \\ 1 & 0 & \dots & 1 & 1 \\ 1 & 1 & \dots & 1 & 1 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 1 & 1 \end{bmatrix} \begin{matrix} y_1 \\ y_2 \\ \vdots \\ y_{Mk} \end{matrix}$$

$\Rightarrow \Theta(i,j) = 1$ if y_i is contained in validation region of x_j

\Rightarrow Hypothesis $\underline{\Omega}_j^{(k)} =$

Assignment Z_j	x_1	x_2	\dots	x_{fa}	x_{nt}
y_1	0	1		0	0
y_2	1	0		0	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
y_{Mk}	0	0		1	0

each row, only one "1"
each column x_n , only one "1"
except x_{fa} & x_{nt}

$\underline{\Omega}_{pj}^{(k-1)}$ is parent hypothesis of $\underline{\Omega}_j^{(k)}$

Z_j : the set of assignments $\underline{\Omega}_j^{(k)}$

\Rightarrow Hypothesis probabilities:

$$\begin{aligned}
 p(\underline{\Omega}_j^{(k)} | \underline{Y}^{(k)}) &= p(Z_j^{(k)}, \underline{\Omega}_{pj}^{(k-1)} | \underline{Y}^{(k)}) \\
 &\stackrel{\text{Bayes}}{=} n \cdot p(\underline{Y}^{(k)} | Z_j^{(k)}, \underline{\Omega}_{pj}^{(k-1)}) \cdot p(Z_j^{(k)}, \underline{\Omega}_{pj}^{(k-1)}) \\
 &= n \cdot p(\underline{Y}^{(k)} | Z_j^{(k)}, \underline{\Omega}_{pj}^{(k-1)}) \cdot p(Z_j^{(k)} | \underline{\Omega}_{pj}^{(k-1)}) \cdot p(\underline{\Omega}_{pj}^{(k-1)})
 \end{aligned}$$

normalization factor measurement likelihood prob. of assignment set prob. of parent

$$+ M_k = N_{\text{det}} + N_{\text{fail}} + N_{\text{new}}$$

+ Measurement likelihood: $\begin{cases} \text{assume Gaussian dist of } y_i^{(k)} \text{ around } \hat{x}_j^{(k)} \text{ with } \hat{\Sigma}_j^{(k)} \\ \text{probabilities of new track over volume } W. \quad W^{-1} \\ \text{false alarm} \end{cases}$

$$\Rightarrow p(\underline{y}^{(k)} | \underline{z}_j^{(k)}, \Omega_{p(j)}^{(k-1)}) = W^{-(N_{\text{fail}} + N_{\text{new}})} \prod_{i=1}^{N_{\text{det}}} \mathcal{N}(y_i^{(k)}; \hat{x}_j^{(k)}, \hat{\Sigma}_j^{(k)})$$

\Rightarrow Probabilities of an Assignment set.

① Probabilities of number of tracks $N_{\text{det}}, N_{\text{fail}}, N_{\text{new}}$

N_{det} : Binomial Dist. while $N_{\text{fail}}, N_{\text{new}}$: Poisson dist

N is no. of tracks in the parent hypothesis

$$\Rightarrow p(N_{\text{det}} | \Omega_{p(j)}^{(k-1)}) = \binom{N}{N_{\text{det}}} \cdot p_{\text{det}}^{N_{\text{det}}} \cdot (1-p_{\text{det}})^{(N-N_{\text{det}})}$$

$$\Rightarrow p(N_{\text{det}}, N_{\text{fail}}, N_{\text{new}} | \Omega_{p(j)}^{(k-1)}) = \binom{N}{N_{\text{det}}} \cdot p_{\text{det}}^{N_{\text{det}}} \cdot (1-p_{\text{det}})^{(N-N_{\text{det}})} \cdot \mu(N_{\text{fail}}; \lambda_{\text{fail}}, W) \cdot \mu(N_{\text{new}}; \lambda_{\text{new}})$$

⊗ ~~$\binom{N}{N_{\text{det}}} = \frac{N!}{N_{\text{det}}!(N-N_{\text{det}})!}$~~ $N_{\text{det}}^N C_{N_{\text{det}}} = \frac{N!}{N_{\text{det}}!(N-N_{\text{det}})!}$ No. of combinations
order doesn't matter

$$\mu(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$$

② Probability of specific assignments of measurements

$$1 / \binom{M_k}{N_{\text{det}}} \binom{M_k - N_{\text{det}}}{N_{\text{fail}}} \underbrace{\binom{M_k - N_{\text{det}} - N_{\text{fail}}}{N_{\text{new}}}}_{=1} = \frac{1}{\binom{M_k}{N_{\text{det}}} \cdot \binom{M_k - N_{\text{det}}}{N_{\text{fail}}}}$$

③ Probability of a specific assignment of tracks

$$1 / \frac{N!}{(N - N_{\text{det}})!} \underbrace{\binom{N - N_{\text{det}}}{N_{\text{det}}}}_{=1} = \frac{(N - N_{\text{det}})!}{N!}$$

Example: a lots in military scenarios
detect, track unauthorized vehicles

- ⊗ - Also works for object with similar appearance model (ants, plane)
observation input (from radar..)
- Fit well with particle filter

$$P(Z_j^{(k)} | \Omega_{p(j)}^{(k-1)}) = \frac{N_{\text{fail}}! \cdot N_{\text{new}}!}{M_k!} \cdot P_{\text{det}}^{N_{\text{det}}} \cdot (1 - P_{\text{det}})^{(N - N_{\text{det}})} \cdot \mu(N_{\text{fail}}; \lambda_{\text{fail}} W) \cdot \mu(N_{\text{new}}; \lambda_{\text{new}} W)$$

+ Practical Issues - Exponential Complexity

- Need to apply heuristics pruning

+ Pros: when combine with current online-learning appearance model
(Initially just popular with radar target tracking)

+ Tricks: use MWIS (Maximum Weight Independent Set) as efficient way to organize MHT's hypotheses

From here down, consider Multi-Object Tracking purely as Data Association Prob.

4) Data Association as Linear Assignment Problem

- Similarity scores based from
 - motion prediction (Mahalanobis dist.)
 - appearance
 - both

⇒ matrix of pairwise similarity scores

⇒ +) Linear Assignment Problem: (LAP)

- maximize $\sum_{i=1}^N \sum_{j=1}^M w_{ij} z_{ij}$
- ↓ ↓
weight assignment
- N tracks such $\sum_{j=1}^N z_{ij} = 1$
 M observations $\sum_{i=1}^M z_{ij} = 1$
 (1 obs per track
 1 track per obs)

- Also could instead minimize cost:

$$\underset{z_{ij}}{\text{argmin}} \sum_{i=1}^N \sum_{j=1}^M c_{ij} z_{ij}$$

⇒ Greedy algorithm: [find largest
 remove scores in same row & column
 pretty good, but not optimal solution]

+) Hungarian algorithm Use this, not Greedy Algorithm

- The Munkres Assignment algorithm:

① Row reduction: minus each row with the min of that row

② Column reduction: column column

③ Draw λ lines to cover all zeros

minimum
If the no. of lines = no. of rows / columns

\Rightarrow optimal solution exists

\Rightarrow Find 0 in way each row & column has only one 0

④ If not, shift zero:

Find smallest uncover value, subtract from all uncover values
add to intersections of lines

⑤ Make final assignments ..

If no. rows \neq no. columns \Rightarrow add rows/ columns of 0

But remove solution in these rows/ columns

- The overall complexity $\Theta(n^3)$

L12/

② LAP as 2-frame problems

for multi-frame problems

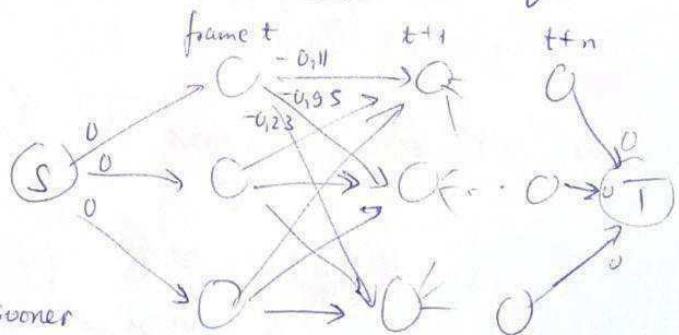
	frame $t+1$		
	1	2	3
a	0,11	0,95	0,23
b	0,35	0,25	0,85
c	0,90	0,12	0,81

6) Tracking as Network Flow Optimization

+) Network flow formulation:

- Min-cost flow
with weights transform to costs

Complication: object appear later, end sooner



\Rightarrow source, single connect to all intermediate nodes

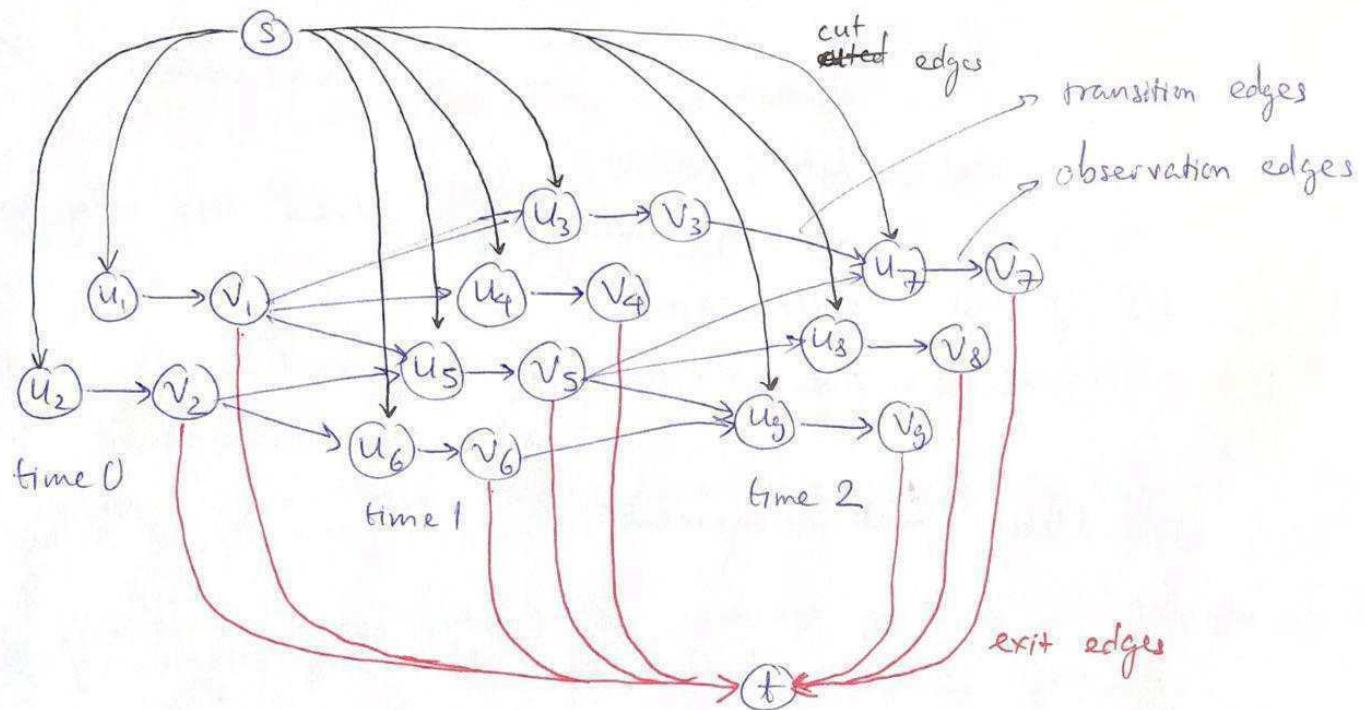
+ Solution: Divide each detection into 2 nodes



detection edge :

$$c_i = \log \frac{\beta_i}{1-\beta_i}$$

prob. that detection
i is a false alarm



+ Objective function:

$$\mathcal{T}^* = \arg\min_{\mathcal{T}} \sum_{\text{source}} C_{in,i} \cdot f_{in,i} + \sum_{\text{sink}} C_{i,out} \cdot f_{i,out} + \sum_{ij} C_{i,j} \cdot f_{i,j} + \sum_i C_{fi} \cdot f_{fi}$$

subject to:

$$\begin{cases} f_{in,i} + \sum_j f_{0,j} = f_i = f_{out,i} + \sum_j f_{i,j} & \forall i \\ f_i \leq 1 \end{cases}$$

$$C_i = -\log(p_i)$$

$$\Leftrightarrow \mathcal{T}^* = \arg\max_{\mathcal{T}} \prod_i P(o_i | \mathcal{T}) \cdot P(\mathcal{T})$$

$$P(\mathcal{T}) = \prod_{T_k \in \mathcal{T}} P(T_k) \quad (\text{Markov})$$

7) Comments

- Pros:
 - [Clear framework, equivalence to probabilistic formulation]
 - [Well understood LP optimization prob]
 - [Globally optimal solution]
- Cons:
 - [not possible to encode exclusion constraints]
 - [only with zero-velocity model, not constant-velocity model]
 - C_{in}, C_{out} are fiddly to set in practice
 - [too low \Rightarrow fragmentations]
 - [too high \Rightarrow ID switches]

Visual Odometry

43)

- Track motion of camera from its images
- Also involves a data association problem

- ~~Complement other motion sensors~~

GPS
 Inertial Measurement Unit (IMU)
 Wheel odometry
 etc

- More accurate than wheel odometry, not prone to slippage
- Important in GPS-denied environments (indoors, close to buildings, Mars)

1) Geometric Background:

+ Geometric Point Primitives

- Point: $\bar{x} = \begin{bmatrix} x \\ y \end{bmatrix} \in \mathbb{R}^2$ $[x \ y \ z]^T \in \mathbb{R}^3$

- Augmented vector: $\bar{x} = [x \ y \ 1]^T \in \mathbb{R}^4$ $(x \ y \ z \ 1) \in \mathbb{R}^4$

- Homogeneous coordinates $\bar{x} = [x \ y \ w]^T \in \mathbb{P}^2$ $[x \ y \ z \ w] \in \mathbb{P}^3$
 $\tilde{x} = \tilde{w} \cdot \bar{x}$

+ Euclidean transformations:

$$\bar{x}' = R\bar{x} + t \quad \bar{x}' = \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix} \cdot \bar{x}$$

+ Rotation:

- Special Orthogonal Group $SO(n)$:
associative
but not commutative
inverse & neutral element exist
 - if R is orthonormal matrix
 - $R \in SO(n) \subset \mathbb{R}^{n \times n}$
 - $\det(R) = 1$, $RR^T = I$

+/-) 3D Rotation Representations

- $R = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \in \mathbb{R}^{3 \times 3}$, orthonormal matrix

\Rightarrow 3 parameters for 3 DoF

\Rightarrow - Euler Angles: minimal with 3 parameters
But with singularities (gimbal lock)

\Rightarrow - Axis & Angle: rotate along axis $n \in \mathbb{R}^3$ by angle $\theta \in \mathbb{R}$

$$R(n, \theta) = I + \sin(\theta)\hat{n} + (1 - \cos(\theta))\hat{n}^2 \quad \|n\|_2 = 1$$

(unit norm constraint)

$$\hat{x} = \begin{bmatrix} 0 & -z & y \\ z & 0 & -x \\ -y & x & 0 \end{bmatrix} \quad \hat{x}y = \hat{x} \times \hat{y} \quad (n, \theta): 4 \text{ params}$$

$w = \theta n: 3 \text{ params}$

+/-) Quaternions

$$q = [q_x \ q_y \ q_z \ q_w]^T \in \mathbb{R}^4, \|q\|_2 = 1$$

- Relation with Axis & Angle:

$$- q(n, \theta) = [n^T \sin(\theta/2), \cos(\theta/2)]$$

$$- n(q) = \begin{cases} (q_x \ q_y \ q_z)^T / \sin(\theta/2) & \theta \neq 0 \\ 0 & \theta = 0 \end{cases}$$

$$\theta = 2 \arccos(q_w)$$

axis angle \Rightarrow quaternion

quaternion \Rightarrow axis & angle

+/-) Special Euclidean Group $SE(3)$

$$T = \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix} \in SE(3) \subset \mathbb{R}^{4 \times 4}$$

$$SE(3) \times SE(3) \rightarrow SE(3)$$

$$T_B^A \cdot T_C^B \rightarrow T_C^A$$

→ Sensor types for Visual Odometry:

- Monocular cameras

(+) low power, light weight, low-cost, simple

Input

$$I_{0:t} = \{I_0, \dots, I_t\}$$

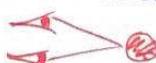
(-) require motion parallax & texture
scale not observable



- Stereo cameras: (left + right)

(+) depth without motion, less power than active sensor

(-) require texture



accuracy depends on baseline, calibration..

have to invest in computation hardware

$$I_{0:t}^l = \{I_0^l, \dots, I_t^l\}$$

$$I_{0:t}^r = \{I_0^r, \dots, I_t^r\}$$

- Active RGB-D sensors

(+) No texture needed



$$I_{0:t} = \{I_0, \dots, I_t\}$$

(-) consume power, blackbox depth estimation

$$Z_{0:t} = \{Z_0, \dots, Z_t\}$$

Mostly indoor, because outdoor, the sun can overwhelm the infrared from our sensor,

⇒ Output of Visual odometry:

relative transformation estimates $T_t^{t-1} \in SE(3)$
(pose at t to $t-1$)

- $T_t \in SE(3)$ is the camera pose at time t in world frame

$$T_t = T_0 T_1^0 \dots T_t^{t-1}$$

→ Direct vs Indirect Methods:

Direct Methods

Photometric error

$$E(\xi) = \int_{u \in \Omega} |I_1(u) - I_2(u, \xi)| du$$

of the whole image

like Lukas-Kanade, but with 6 dof in 3D

Indirect Methods

Reprojection error of geometric primitives
(points, lines...)

$$E(\xi) = \sum_i |y_{1,i} - a(y_{2,i}, \xi)|$$

f

L14

2) Point-based Visual Odometry as indirect method

2D to 2D Motion Estimation

match points from 2D image
to 2D points of previous image
observations

- Given corresponding image point

$$Y_t = \{y_{t,1}, y_{t,2}, \dots, y_{t,N}\} \text{ & } Y_{t-1} = \{y_{t-1,1}, y_{t-1,2}, \dots, y_{t-1,N}\}$$

of unknown 3D points $X = \{x_1, x_2, \dots, x_N\}$

- Goal: determine relative motion T_t^{t-1} between frames

- Reprojection error: nonlinear, non-convex

$$\arg \min_{T_t^{t-1}} E(T_t^{t-1}, X) = \sum_{i=1}^N \| \tilde{y}_{t,i} - \pi(x_i) \|_2^2 + \| \tilde{y}_{t-1,i} - \pi(T_t^{t-1} \tilde{x}_i) \|_2^2$$

\Rightarrow not convex
 \Rightarrow thus leads to multiple solution
 \Rightarrow thus use good initialization first

- Essential matrix: $E = \hat{\tau} R = [t]_{\times} \cdot R$
 $\tilde{y}^T \cdot E \cdot \tilde{y} = 0$

As a linear approach to estimate T_t^{t-1}

② Eight-point Algorithm for Essential Matrix

- ① Rewrite epipolar constraints as a linear system of equations

- For each pair of correspondence \tilde{y}_i & \tilde{y}'_i : $\tilde{y}_i^T E \tilde{y}'_i = 0 = a_i E_s$
 \Rightarrow All pairs: $A \cdot E_s = 0$ with $E_s = [e_{11}, e_{12}, \dots, e_{33}]^T$

$$A = [a_1^T, a_2^T, \dots, a_N^T]$$

- ② Apply SVD to A : $A = U_A \Sigma_A V_A^T$

$A \cdot E_s = 0 \Rightarrow E_s$ is the last column of V_A (the 9th column)

- ③ Project approximate \tilde{E} into the normalized essential space

$$\tilde{E} = U \cdot \text{diag}(\tilde{\sigma}_1, \tilde{\sigma}_2, \tilde{\sigma}_3) \cdot V^T \text{ with } U, V \in \text{SO}(3)$$

Replace $\tilde{\sigma}_1, \tilde{\sigma}_2, \tilde{\sigma}_3$ with 1, 1, 0 to enforce rank 2

$$E = U \cdot \text{diag}(1, 1, 0) \cdot V^T$$

- ④ Decomposition of E to R and T

$$R = U \cdot R_z\left(\pm \frac{\pi}{2}\right) \cdot V^T$$

$$\hat{\tau} = U \cdot R_z\left(\pm \frac{\pi}{2}\right) \cdot \text{diag}(1, 1, 0) \cdot U^T$$

there will be 2 solutions
 $2R, 2T \Rightarrow 4$ sols. choose 1 that the point is in front of both cam
 \Rightarrow the depth > 0 :

$R_z(\varphi)$: as the rotation matrix around z-axis to angle φ $\tilde{y}_i = \lambda' R \tilde{y}'_i + \gamma^T$
 $\Rightarrow \lambda, \lambda' > 0$

④ Remark: - Eight-Point Algorithm is used as initialization for non-linear least-square reprojection error.

- Normalized essential matrix $\|E\| = 1$
(stack to column vector, then normalized with L_2)
- There is also 5-points algorithm (nonlinear)
- Will have degenerated solution if all are on same plane
- Normalized Eight-Point Algorithm: better accuracy

④ Triangulation (repeated)

- Goal: reconstruct 3D points: $\tilde{x} = [x \ y \ z \ w]^T \in \mathbb{R}^4$
from 2D observation $\{y_1, \dots, y_N\}$ and known camera pose (T_1, \dots, T_N)
- Reprojection: $y'_i = \pi(T_i \tilde{x}) = \begin{cases} r_{11}x + r_{12}y + r_{13}z + t_x w \\ r_{21}x + r_{22}y + r_{23}z + t_y w \\ r_{31}x + r_{32}y + r_{33}z + t_z w \end{cases}$

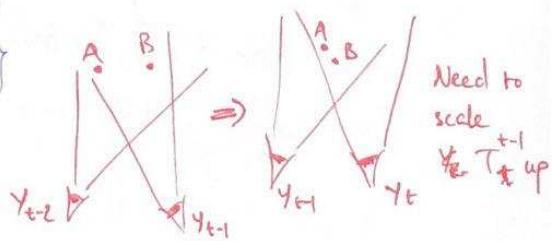
+ Linear solution: each image: $y_i - y'_i = 0$

$$\Rightarrow A \cdot \tilde{x} = 0$$

\Rightarrow Set $w=1 \Rightarrow$ solve non-homogeneous system

[Find null space of A with SVD]

+ Non-linear solution: $\min_{\tilde{x}} \left\{ \sum_{i=1}^N \|y_i - y'_i\|_2^2 \right\}$



+ Relative Scale Recovery:

- Triangulation overlapping points y_{t-2}, y_{t-1}, y_t
 \Rightarrow points: $x_{t-2,t-1}, x_{t-1,t}$

\Rightarrow Rescale the translation of current relative pose estimate to match reconstruction scale with the distance ratio between corresponding pairs

$$r_{i,j} = \frac{\|x_{t-2,t-1,i} - x_{t-2,t-1,j}\|_2}{\|x_{t-1,t,i} - x_{t-1,t,j}\|_2} \quad (\text{of 1 pair})$$

\Rightarrow Use mean / robust median over all available pairs

+ 2D to 3D Motion Estimation

- Given a local set of 3D points: $X = \{x_1, x_2, \dots, x_N\}$ and corresponding image observations $Y_t = \{y_{t,1}, y_{t,2}, \dots, y_{t,N}\}$
- Goal: determine T_t within local map
- Minimize geometric reprojection error

$$E(T_t) = \sum_{i=1}^N \|y_{t,i} - \pi(T_t x_i)\|_2^2 \quad (\text{PnP})$$

⇒ Direct Linear Transform (DLT) one method of Perspective n Point pose

- Goal: determine reprojection matrix $P = [R \ t] \in \mathbb{R}^{3 \times 4} = \begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix}$

2D: $\tilde{y}_i = [x_i \ y_i \ w_i]^T \rightarrow 3D: \tilde{x}_i = [x_i \ y_i \ z_i \ w_i]^T \in \mathbb{P}^3$

$$\Rightarrow 2 \text{ constraints: } \begin{bmatrix} 0 & -w_i \tilde{x}_i^T & y_i \tilde{x}_i^T \\ w_i \tilde{x}_i^T & 0 & -x_i \tilde{x}_i^T \end{bmatrix} \cdot \begin{bmatrix} P_1^T \\ P_2^T \\ P_3^T \end{bmatrix} = 0 \quad \text{from } \tilde{y}_i \cdot (P \tilde{x}_i) = 1$$

⇒ Linear system of equations: $A p = 0$ with $p = \begin{bmatrix} P_1^T \\ P_2^T \\ P_3^T \end{bmatrix} \in \mathbb{R}^9$ (?)

from $N \geq 6$ correspondences

L15) + 3D to 3D Motion Estimation

- Given 3D point coordinates of corresponding points in 2 frames $X_{t-1} = \{x_{t-1,1}, \dots, x_{t-1,N}\}$ & $X_t = \{x_{t,1}, \dots, x_{t,N}\}$
- Goal: determine relative camera pose T_t^{t-1}
- Geometric least square error:

$$E(T_t^{t-1}) = \sum_{i=1}^N \|\bar{x}_{t,i} - T_t^{t-1} \bar{x}_{t-1,i}\|_2^2 \quad N \geq 3$$

Closed form solution available

- Means of 3D points. $\mu_{t-1} = \frac{1}{N} \sum_{i=1}^N x_{t-1,i}$, $\mu_t = \frac{1}{N} \sum_{i=1}^N x_{t,i}$

$$\Rightarrow \text{Set } A = \sum_{i=1}^N (x_{t-1,i} - \mu_{t-1})(x_{t,i} - \mu_t)^T, \quad A = U \Sigma V^T$$

$$\Rightarrow R_{t-1}^t = V \cdot U^T \quad \text{and} \quad t_{t-1}^t = \mu_t - R_{t-1}^t \mu_{t-1}$$

Algorithm: 2D to 2D Visual Odometry:

- Input: image sequence $I_{0:t}$
- Output: aggregated camera pose $T_{0:t}$

- For each I_k :

1) Extract & match key points between $I_{k-1} \times I_k$

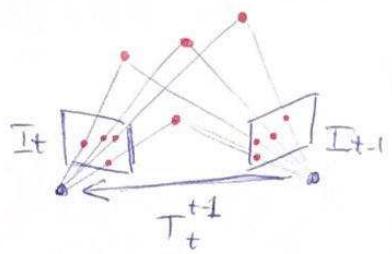
2) Compute T_k^{k-1} from E between I_{k-1}, I_k

3) Compute relative scale and rescale translation of T_k^{k-1}

4) Aggregate $T_k = T_{k-1} T_k^{k-1}$

$$E(T_t^t, X) = \sum_{i=1}^N \|y_{t,i} - \pi(X_i)\|_2^2 + \|y_{t-1,i} - \pi(T_t^{t-1} \bar{x}_t)\|_2^2$$

8 pts
algorithm



Algorithm: 2D to 3D Visual Odometry:

- Input: image sequence $I_{0:t}$
- Output: aggregated camera pose $T_{0:t}$

- Initialize:

1) Extract & match key points I_0 and I_1

2) Determine camera pose E and triangulate 3D point X_1

- For each current I_k :

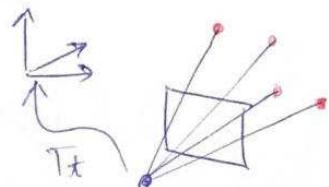
1) Extract & match key point between $I_{k-1} \times I_k$

2) Compute T_k using PnP from 2D-to-3D matches using \bar{x}_k from X_{k-1}

3) Triangulate all new key point matches between $I_{k-1} \times I_k$
and add to local map X_k

$$E(T_t) = \sum_{i=1}^N \|y_{t,i} - \pi(T_t \bar{x}_i)\|_2^2$$

X_k is not really accurate here
In SLAM, we will reestimate position
of points



Algorithm 3D to 3D Visual Odometry:

- Input: stereo image sequence: $I_{0:t}^l, I_{0:t}^r$
- Output: aggregated camera pose $T_{0:t}$

- For each current stereo image I_k^l, I_k^r

1) Extract and match key points between I_k^l and I_{k-1}^l

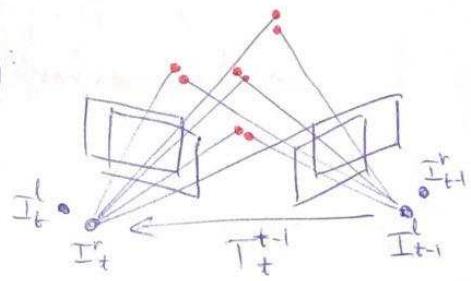
2) Triangulate 3D points X_k between I_k^l and I_k^r

3) Compute camera pose T_k^{t-1} from 3D-to-3D point matches X_k to X_{k-1}

4) Aggregate $T_k = T_{k-1} \cdot T_k^{t-1}$

$$E(T_t^{t-1}) = \sum_{i=1}^N \|\bar{x}_{t-1,i} - T_t^{t-1} \bar{x}_{t,i}\|_2^2$$

Arun's
Method



→ Further Consideration

- Key point Detectors : Harris, FAST
LoG, DoG, SIFT, SURF

- Data association : $\begin{cases} \text{match by reprojection error} \\ \text{detect-then-track KLT Tracker} \\ \text{matching by descriptor} \\ \text{robustness through RANSAC} \end{cases}$

- RANSAC

$$\# \text{iteration } N = \frac{\log(1-p)}{\log(1-(1-\epsilon)^s)}$$

p : success prob
 s : required data points
 ϵ : % outlier

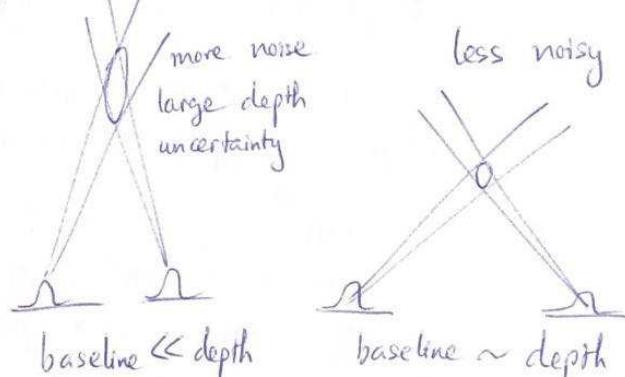
- Probabilistic Modelling:

$$p(X, \xi | Y) \propto p(Y | X, \xi) \cdot p(X, \xi)$$

$$= p(X, \xi) \prod_{i=1}^N p(y_i | x_i, \xi)$$

$$E(X, \xi) = \text{const} + (\xi - \mu_{\xi, 0})^T \sum_{\xi, 0}^{-1} (\xi - \mu_{\xi, 0}) \quad \text{baseline} \ll \text{depth}$$

$$+ \sum_{i=1}^N (x_i - \mu_{x_i, 0})^T \sum_{x_i, 0}^{-1} (x_i - \mu_{x_i, 0}) + (y_i - \pi(T(\xi)x_i))^T \sum_{y_i}^{-1} (y_i - \pi(T(\xi)x_i))$$



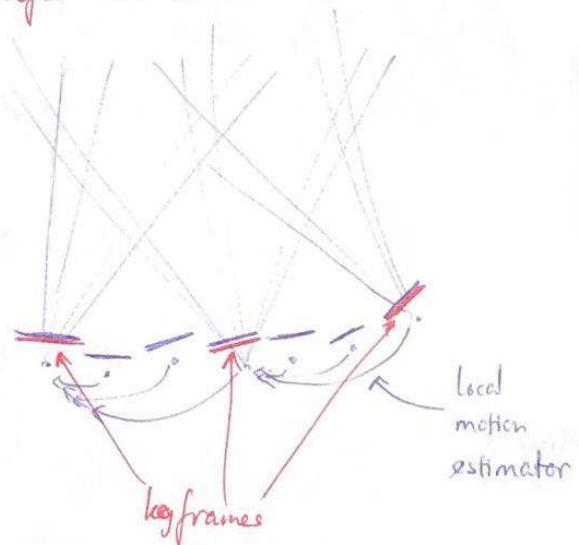
Note: 3D-to-3D use stereo camera \Rightarrow obviously baseline \ll depth

\Rightarrow actually less accurate than 2D-to-3D NO

\Rightarrow Thus, even though 3D-to-3D is simpler, still prefer 2D-to-3D more

→ → Key frames

- Have a local motion estimate
- But only when the motion is great enough that the baseline is sufficient
- \Rightarrow create new key frame
- Only aggregate motion considering key frames



+ Local optimization windows

Optimize over more than 2 Images

- Bundle adjustment

$$E(X_{t-k:t}, \xi_{t-k:t}) = \sum_{j=0}^k \sum_{i=0}^{N_{t-j}} \|y_{t-j,i} - \pi(T(\xi_{t-j}), x_{t-j,i})\|_2^2$$

\Rightarrow minimize a bundle of transformation, not just 1 at a time

avoid manually designed key points

3) Direct Methods

+ Warping requires depth

\Rightarrow need camera:

- RGB-D
- Fixed-baseline stereo
- Temporal stereo, ...

$$E(\xi) = \int_{u \in \Omega} |I_1(u) - I_2(\omega(u, \xi))| du$$

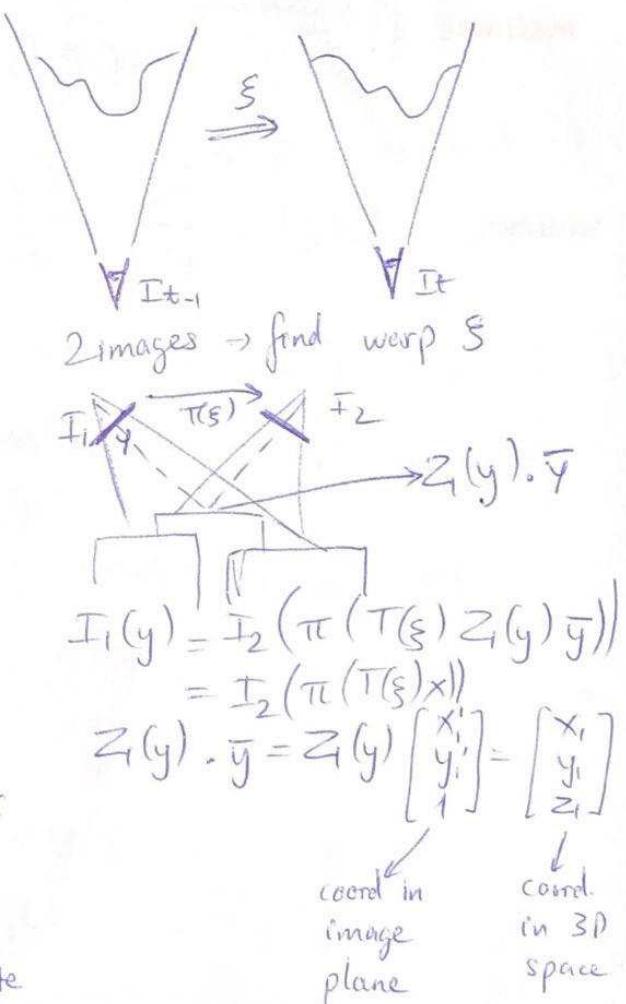
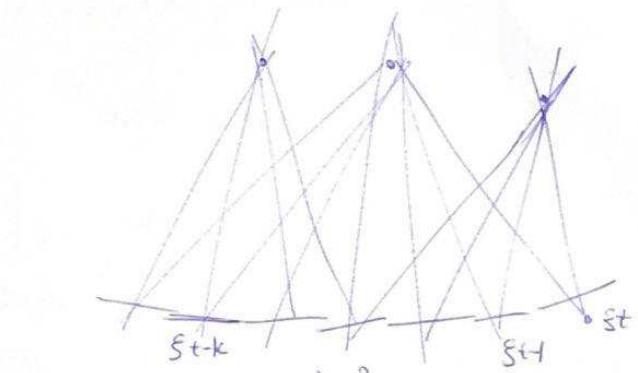
(over the whole image)

$$+, I_1(y) = I_2(\pi(T(\xi)Z_1(y), \bar{y}))$$

$$Z_1(y) \cdot \bar{y} = Z_1(y) \cdot \begin{bmatrix} x'_1 \\ y'_1 \\ 1 \end{bmatrix} = \begin{bmatrix} x_1 \\ y_1 \\ z_1 \end{bmatrix} = x$$

coordinates in image plane coordinate in 3D space

$I_2(\underbrace{\pi(T(\xi)x)}_{\text{reprojection of } T(\xi)x \text{ to image plane}}) = \text{image at that pose}$



$$I_1(y) = I_2(\pi(T(\xi)Z_1(y), \bar{y}))$$

$$= I_2(\pi(T(\xi)x))$$

$$Z_1(y) \cdot \bar{y} = Z_1(y) \begin{bmatrix} x'_1 \\ y'_1 \\ 1 \end{bmatrix} = \begin{bmatrix} x_1 \\ y_1 \\ z_1 \end{bmatrix}$$

coord.
in
image
plane

coord.
in
3D
space

+ Photometric error: $I_1(y) = I_2(\pi(T(\xi) \cdot Z_1(y) \cdot \bar{y}))$

Geometric error: $[T(\xi) Z_1(y) \bar{y}]_z = Z_2(\pi(T(\xi) Z_1(y) \bar{y}))$

depth of points after
apply transformation $T(\xi)$ depth in image 2 of ..

+ With noise: $I_1(y) = I_2(\pi(T(\xi) Z_1(y) \bar{y})) + \varepsilon$
 $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

- If assume that pixel measurements are stochastically independent

maximize $p(\xi | I_1, I_2) \propto p(I_1 | \xi, I_2) \cdot p(\xi)$
 $\propto p(\xi) \prod_{y \in \Omega} \mathcal{N}(I_1(y) - I_2(\pi(T(\xi) \cdot Z_1(y) \cdot \bar{y})), 0, \sigma^2)$

minimize $E(\xi) = \sum_{y \in \Omega} \frac{r(y, \xi)^2}{\sigma^2}$ with $r(y, \xi) = I_1(y) - I_2(\pi(T(\xi) \cdot Z_1(y) \cdot \bar{y}))$

- Stack into vector: $E(\xi) = r(\xi)^T W r(\xi)$

→ This is nonlinear least square problem

+ Gauss-Newton: to solve nonlinear least square

Linearize $\tilde{r}(\xi) = r(\xi_i) + \nabla_{\xi} r(\xi_i)(\xi - \xi_i)$ $J_i := \nabla_{\xi} r(\xi_i) \in \mathbb{R}^{\dim(r) \times \dim(\xi)}$

$\xi = \underset{\xi}{\operatorname{argmin}} \tilde{E}(\xi) = \frac{1}{2} \tilde{r}(\xi)^T W \tilde{r}(\xi)$

$\nabla_{\xi} \tilde{E}(\xi) = J_i^T W \tilde{r}(\xi)$

$\nabla_{\xi}^2 \tilde{E}(\xi) = J_i^T W J_i = H_i \in \mathbb{R}^{\dim(\xi) \times \dim(\xi)}$

Find minimum of linearized system, linearize and set $\nabla_{\xi} \tilde{E}(\xi) = 0$

$$\xi_{i+1} = \xi_i - (\nabla_{\xi}^2 \tilde{E}(\xi_i))^{-1} \nabla_{\xi} \tilde{E}(\xi_i)$$

$$= \xi_i - H_i^{-1} J_i^T W r(\xi_i)$$

→ Levenberg - Marquadt Method

$$\xi_{i+1} = \xi_i - (H_i + \lambda I)^{-1} J_i^T W r(\xi_i)$$

- If error increase, reject update, increase $\lambda \Rightarrow$ act like gradient descent
- If error decrease, decrease $\lambda \Rightarrow$ act like Gauss-Newton

→ Practical implementation:

H_i and b_i are summed over individual pixels

$$H_i = \sum_{y \in \Omega} \frac{w(y, \xi_i)}{\sigma_i^2} J_{i,y}^T J_{i,y} \quad ; \quad b_i = \sum_{y \in \Omega} J_{i,y} \frac{w(y, \xi_i)}{\sigma_i^2} r(y, \xi_i)$$

$$J_{i,y} = \nabla_{\xi} r(y, \xi \oplus \xi_i)$$

allow parallel processing with GPU

→ Lie Algebra:

$$\xi = \begin{bmatrix} \omega \\ v \end{bmatrix} \in \mathbb{R}^6 \quad \omega \in \mathbb{R}^3$$

$$\hat{\xi} = \begin{bmatrix} \hat{\omega} & v \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{4 \times 4}$$

From a twist in Lie Algebra to transformation

$$T = \exp(\hat{\xi}) = \begin{bmatrix} \exp(\hat{\omega}) & Av \\ 0 & 1 \end{bmatrix}$$

$$\text{with } \exp(\hat{\omega}) = I + \frac{\sin|\omega|}{|\omega|} \hat{\omega} + \frac{1 - \cos|\omega|}{|\omega|^2} \hat{\omega}^2$$

$$A = I + \frac{1 - \cos|\omega|}{|\omega|^2} \hat{\omega} + \frac{1 - \cos|\omega|}{|\omega|^3} \hat{\omega}^2$$

- From a transformation to a twist

$$\xi = \log(T) = \begin{bmatrix} \log R & A^{-1}t \\ 0 & 0 \end{bmatrix}, \quad T = \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix}$$

$$\log(R) = \frac{1\omega_1}{2\sin|\omega_1|} (R - R^T), \quad |\omega_1| = \cos^{-1}\left(\frac{\text{tr}(R) - 1}{2}\right)$$

- Inversion of hat operator

$$\begin{bmatrix} 0 & -\omega_3 & \omega_2 & v_1 \\ \omega_3 & 0 & -\omega_1 & v_2 \\ -\omega_2 & \omega_1 & 0 & v_3 \\ 0 & 0 & 0 & 0 \end{bmatrix}^V = [\omega_1 \ \omega_2 \ \omega_3 \ v_1 \ v_2 \ v_3]^T$$

$$\Rightarrow \xi(T) = (\log(T))^V, \quad T(\xi) = \exp(\hat{\xi})$$

- Pose inversion:

$$\xi^{-1} = \log(T(\xi)^{-1}) = -\xi$$

- Pose concatenation

$$\xi_1 \oplus \xi_2 = (\log(T(\xi_2)T(\xi_1)))^V$$

- Pose difference

$$\xi_1 \ominus \xi_2 = (\log(T(\xi_2)^+T(\xi_1)))^V$$

+ Practical Considerations

- Coarse-to-fine Optimization
- Residual distributions

SLAM

Z16)

- Simultaneous Localization and Mapping
 - Visual SLAM: SLAM with vision sensors
- To ensure map consistency → global / local optimization methods

Global

Bundle adjustment

Pose-graph optimization

Offline (create a map
then optimize)

Local

Incremental tracking & mapping

Visual Odometry with local maps

Online (not as good as
global one, but faster)

Hybrids

Realtime local SLAM

+
Global optimization
in a slower parallel process

- +) Challenges - wrong correspondences → divergence
- important to model uncertainties of observations

Formulation:

- Input: images at discrete time steps t
control input $u_{1:t}$
- Output: camera pose $\mathbf{g}_t = \mathbf{g}(T_t)$ $T_t \in \text{SE}(3)$
environment map M
- Data association: $c_{t,i} = j$, $1 \leq i \leq N$, $1 \leq j \leq S$
→ link $y_{t,i}$ with m_j
- Probability:
 - SLAM posterior prob. $p(\mathbf{g}_{0:t}, M | Y_{0:t}, U_{1:t})$
 - Observation likelihood $p(Y_t | \mathbf{g}_t, M)$
 - State transition prob. $p(\mathbf{g}_t | \mathbf{g}_{t-1}, U_t)$

27

→ Application: Robotics, Augmented / Virtual Reality
 Programming, Autonomous Driving

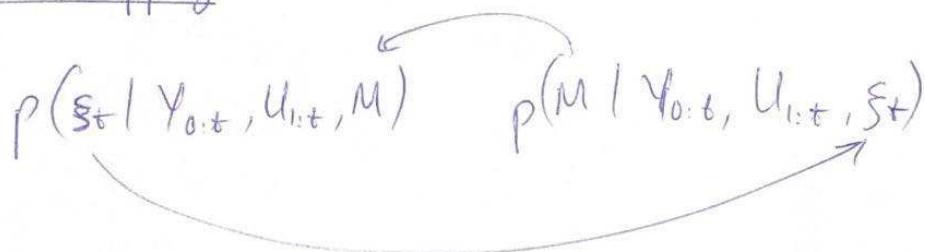
4.2) ↴ Online SLAM:

- Marginalize out previous pose

$$p(\xi_t, M | \gamma_{0:t}, u_{1:t}) = \int \dots \int p(\xi_{0:t}, M | \gamma_{0:t}, u_{1:t}) d\xi_{t-1} \dots d\xi_0$$

- [Tracking-and-Mapping: Alternating pose and map estimation]
- [probabilistic filters: EKF-SLAM]

↳ Tracking & Mapping



- Examples: Semi-dense direct visual odometry
 Kinect Fusion
 PTAM (Parallel Tracking and Mapping)

↳ Filter-based Visual SLAM

- EKF SLAM
- Mono SLAM

P_{Car}

2) EKF SLAM: include landmarks into state variables

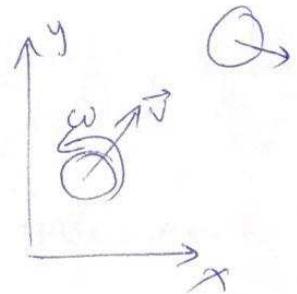
$$x_t = [\xi_t \ m_{t,1} \ \dots \ m_{t,s}]^T, \quad \Sigma_t = \begin{bmatrix} \Sigma_{t,\xi\xi} & \Sigma_{t,\xi m_1} & \dots & \Sigma_{t,\xi m_s} \\ \Sigma_{t,m_1\xi} & \Sigma_{t,m_1m_1} & \dots & \Sigma_{t,m_1m_s} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{t,m_s\xi} & \Sigma_{t,m_sm_1} & \dots & \Sigma_{t,m_sm_s} \end{bmatrix}$$

$$= \begin{bmatrix} \Sigma_{t,\xi\xi} & \Sigma_{t,\xi m} \\ \Sigma_{t,m\xi} & \Sigma_{t,mm} \end{bmatrix}$$

④ Examples EKF SLAM in 2D World:

$$\xi_t = [x_t \ y_t \ \theta_t]^T \quad m_{t,j} = [m_{t,j,x} \ m_{t,j,y}]^T$$

$$u_t = [v_t, \omega_t]^T$$



$$\xi_t = g_\xi(\xi_{t-1}, u_t) + \epsilon_{\xi,t}, \quad \epsilon_{\xi,t} \sim \mathcal{N}(0, \Sigma_{dt,\xi})$$

Transition model

$$g_\xi(\xi_{t-1}, u_t) = \begin{bmatrix} x_{t-1} \\ y_{t-1} \\ \theta_{t-1} \end{bmatrix} + \begin{bmatrix} -v_t/\omega_t \cdot \sin \theta_{t-1} + v_t/\omega_t \cdot \sin(\theta_t + \omega_t \Delta t) \\ v_t/\omega_t \cdot \cos \theta_{t-1} - v_t/\omega_t \cdot \cos(\theta_t + \omega_t \Delta t) \\ \omega_t \cdot \Delta t \end{bmatrix}$$

$$m_t = g_m(m_{t-1}) = m_{t-1}$$

$$\Rightarrow x_t = g(x_{t-1}, u_t) + \epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, \Sigma_{dt})$$

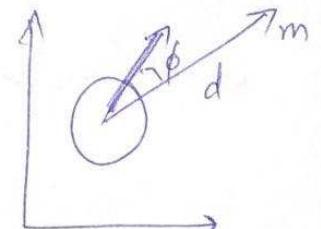
$$g(x_{t-1}, u_t) = \begin{bmatrix} g_\xi(\xi_{t-1}, u_t) \\ g_m(m_{t-1}) \end{bmatrix}, \quad \Sigma_{dt} = \begin{bmatrix} \Sigma_{dt,\xi} & 0 \\ 0 & 0 \end{bmatrix}$$

Observation model

$$y_t = h(\xi_t, m_{t,c_t}) + \delta_t, \quad \delta_t \sim \mathcal{N}(0, \Sigma_{mt})$$

$$= [d_t, \phi_t]^T$$

$$h(\xi_t, m_{t,c_t}) = \begin{bmatrix} \|m_{t,c_t}^{\text{rel}}\|_2 \\ \text{atan}2(m_{t,c_t,y}^{\text{rel}}, m_{t,c_t,x}^{\text{rel}}) \end{bmatrix}$$



$$m_{t,c_t}^{\text{rel}} = R(-\theta_t)(m_{t,c_t} - [x_t \ y_t])$$

- State initialization: $\bar{x}_0 = 0$

$$\bar{\Sigma}_{0,ss} = 0$$

New landmark: $\bar{\Sigma}_{0,sm} = \bar{\Sigma}_{0,mg} = 0$

$$\bar{\Sigma}_{0,mm} = \infty I$$

- On Prediction: covariances are transformed to the new pose
- On Correction: update all state covariances related to the pose and new landmark

+ Problems - Full covariance matrix is huge, expensive

- Fast SLAM: consider landmarks to be independent but introduce more drift (because there still actual correlation)
- Use key frames → still have correlation but less (only between key frames)

3) Loop closure:

- Old landmarks get reobserved

→ Add strong correlations

→ Pose & Landmarks are connected

→ Reduce uncertainty

⊗ Again, wrong correspondences lead to divergence

④ Example Mono SLAM

↳ - Image Patches 9×9 , 15×15

- Any view point with only 12 landmarks

- Map initialized with known pattern

9) Full SLAM:

- Optimize for whole trajectory and all landmarks at once
- Uses future measurements as well to update past poses
- Allow relinearization of all state transitions & measurements at each optimization step

+ SLAM Graph Optimization:

- Just like Bundle Adjustment, but with control inputs
- Common map element observation \rightarrow constraints between poses

+ Pose Graph Optimization:

- Relative pose constraints from image observations

Deep learning for Video Analysis

48)

1) CNNs for Video Analysis

Video classification
Image tagging
Video to Text Description

- Architecture: different ways to fuse features from multiple frame
 - single frame
 - Early fusion
 - Late fusion
- Recap - RNN: $y_n = W_h y \cdot h_n$
 $h_n = \max \{ 0, W_{xh} x_n + W_{hh} h_{n-1} \}$
- LSTM: each module has 9 layers

49/2) CNNs for motion estimation

⊗ Matching & Correspondence estimation

- Learning similarity functions with Siamese Network
 - use in
 - Face Identification
 - Stereo estimation
 - Optical flow
- Metric learning:
 - Contrastive loss: either pull or push
 $\|f(x) - f(x_+)\| \rightarrow 0, \|f(x) - f(x_-)\| \geq m$
 - Triplet loss: pull & push at the same time
 $\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2$

+ Spatial Transformer Network : for patch normalization

Input : a patch \rightarrow output

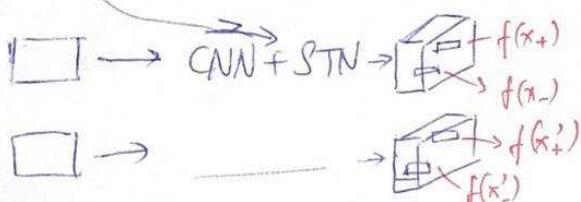
scale
rotation

+ Universal Correspondence Network : compute a patch descriptor

- Usage depends on input data

\rightarrow could be used for

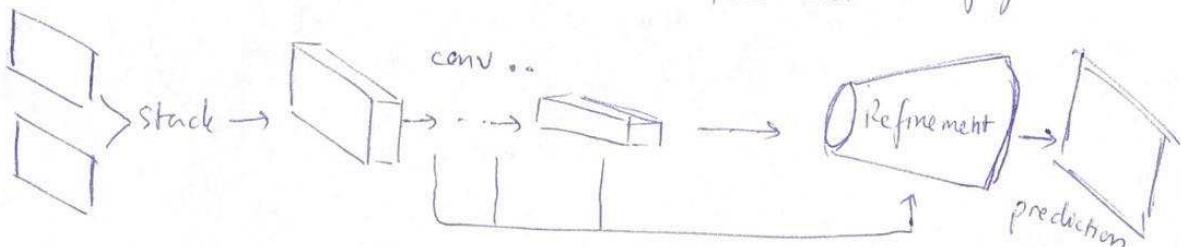
stereo matching
estimating optical flow
semantic matching



Optical Flow Estimation:

+ FlowNet S: Simple Design

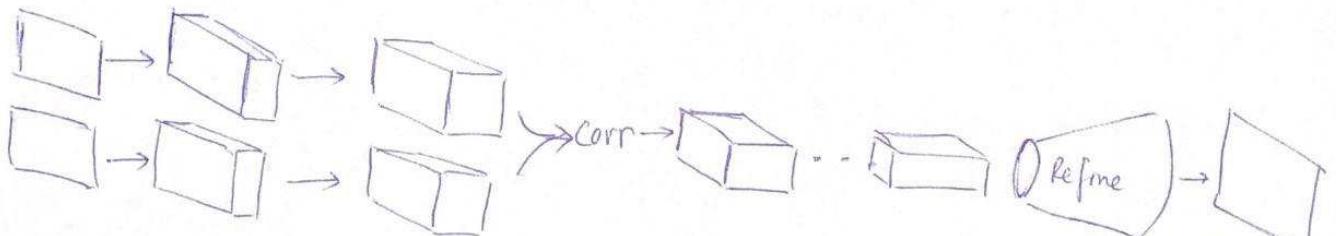
Simply Stack 2 images
The network figures out itself



close to:
pre-CNN
methods

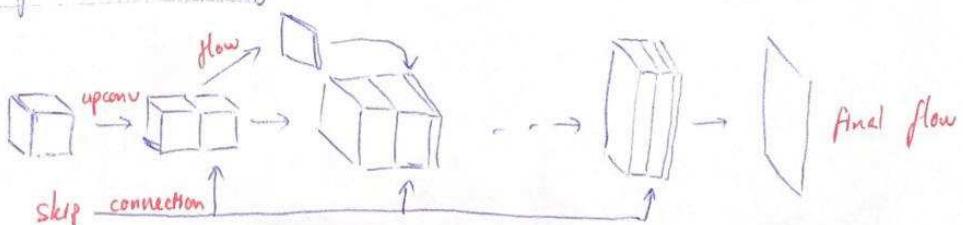
+ FlowNet C: Correlation Design

but much
faster



$$c(x_1, x_2) = \sum_{\delta \in [-k, k] \times [-k, k]} \langle f_1(x_1 + \delta), f_2(x_2 + \delta) \rangle$$

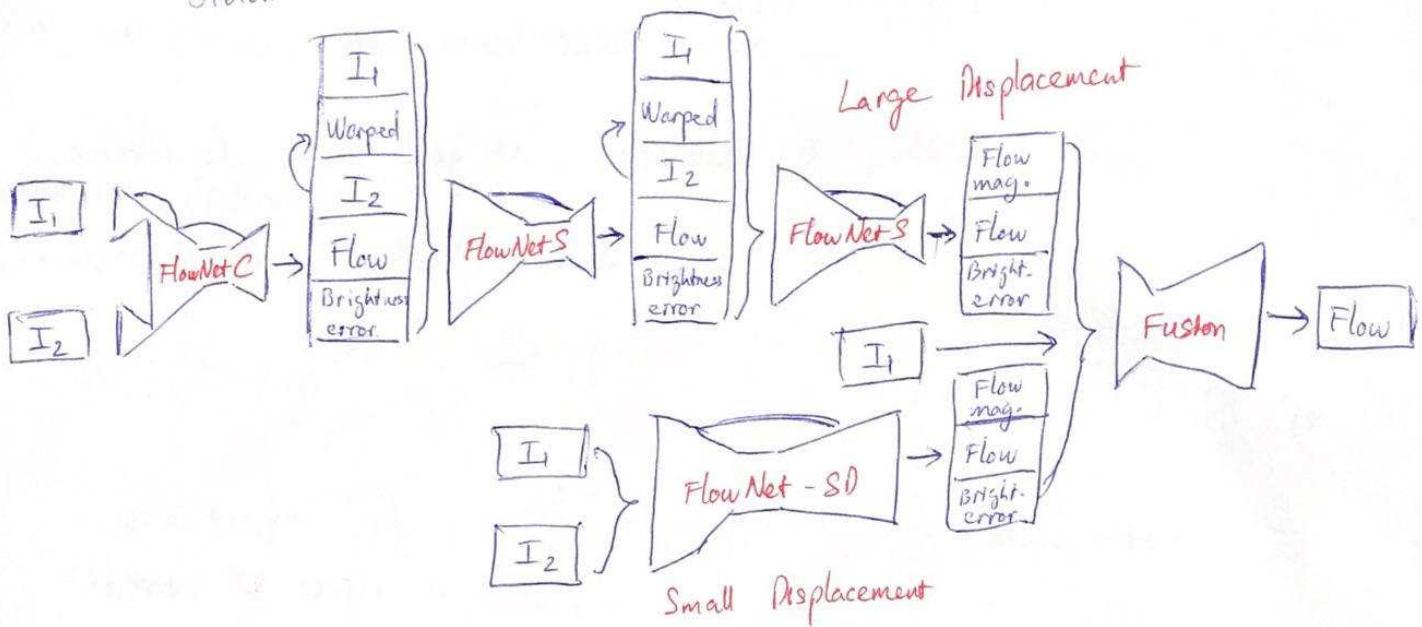
+ Refinement stage



+ Flow Net 2.0

improved design

- Stack FlowNetC and FlowNetS

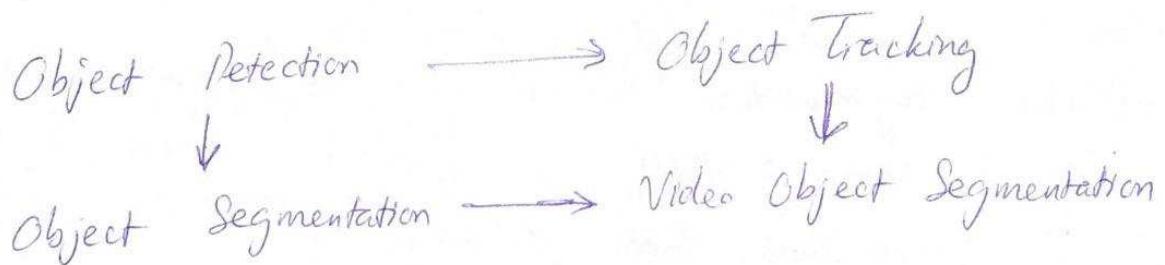


+ FlowNet - SD (FlowNet - Small Displacement)

- Replace 5×5 , 7×7 by 3×3 kernels
assume stride 1 instead of 2

→ FlowNet 2.0: similar accuracy as best pre-CNN methods
but much faster

L20, 3, Video Object Segmentation (VOS)



+ Task formulation:

- Semi-supervised VOS: given object mask in the first frame
→ pixel-accurate segmentation of entire video
- Unsupervised VOS (UVOS/ZVOS): no 1st-frame mask

+ Two Strategies:

- Mask propagation : given mask in 1st frame
try to learn an online classification for that
- Proposal Tracking -by- Detection : at each frame, do instance segmentation
then \rightarrow Temporal association problem

+ Basic elements

- First-frame fine-tuning: pre-training for objectness
then adaptation to object of interest
Ex: OSVOS - applicable, but slow
OSVOS-S with Semantic Prior
- Online Adaptation: not like Online AdaBoost, as we do it with CNN
 \Rightarrow extremely slow
Ex: OnAVOS: result improved greatly
but with great computational cost
- Mask Refinement: refine mask from previous frame
Ex: Mask Track
- Optical Flow Mask Propagation: warp segmentation mask from one frame to the next one
- Data Augmentation: because we could overfit to 1st-frame appearance
 \Rightarrow crop out object, \rightarrow augmentation (position, lighting, background!)
Ex: Lucid Data Dreaming
- Object Appearance Re-Identification: when object go in/out of view
 \Rightarrow need to re-identify objects
Ex: ReID-VOS
- Proposal Generation: Instance Segmentation Networks
Ex: Mask-RCNN

+ PReMVOS combines all of the previous VOS principles

- Proposal generation : Mask R-CNN (ResNet101 backbone)
- Refinement : Mask Refinement
- Merging : First-frame fine tuning + Data Augmentation
Mask propagation + Optical flow
Re-ID score

+ Combining Mask Propagation & Re-ID

- Ex RGMP (Region Guided Mask Propagation)

+ Using Recurrent NN : but with great cost for memory

- Ex: STM (Space Time Memory Networks)
Seq2Seq

⊗ VOS to MOTS (Multi Object Tracking & segmentation)

- VOS is restricted:
 - 1st-frame mask given
 - short video with objects present in almost all frames
 - few objects to track

⇒ MOTS deals with these short-comings

CV2 - Formulas

~~YB~~
Ng Kuan Due

1) Background Modelling:

$$+ I(t) - B \rightarrow \text{abs} \rightarrow \text{threshold } \lambda \rightarrow M(t)$$

$$+ I(t) - B(t-1) \rightarrow \text{abs} \rightarrow T(\lambda) \rightarrow M(t)$$

$$+ D(N) = \|I(t) - I(t+N)\|$$

$$+ B(t) = \alpha, I(t) + (1-\alpha) B(t-1)$$

$$I(t) - B(t-1) \rightarrow \text{abs} \rightarrow T(\lambda) \rightarrow M(t)$$

$$+ \hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n, \quad \hat{\sigma}^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \hat{\mu})^2$$

$$\hat{\mu}^{(t+1)} = \hat{\mu}^{(t)} + \frac{1}{N} x^{(t+1)} - \frac{1}{N} x^{(t+1-T)}$$

$$(\hat{\sigma}^2)^{(t+1)} = (\hat{\sigma}^2)^{(t)} + \frac{1}{N-1} [x^{(t+1)} - \hat{\mu}^{(t+1)}]^2 - \frac{1}{N-1} [x^{(t+1-T)} - \hat{\mu}^{(t+1-T)}]^2$$

$$\hat{\mu}^{(t+1)} = (1-\alpha) \hat{\mu}^{(t)} + \alpha, x^{(t+1)}$$

$$(\hat{\sigma}^2)^{(t+1)} = (1-\alpha)(\hat{\sigma}^2)^{(t)} + \alpha [x^{(t+1)} - \hat{\mu}^{(t+1)}]^2$$

$$+ p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k), \quad \Sigma_k = \sigma_k^2 I$$

$$B = \underset{b}{\operatorname{argmin}} \left(\sum_{k=1}^b \frac{\pi_k}{\sigma_k} > \epsilon \right)$$

$$[I_x \quad I_y] \cdot \begin{bmatrix} u \\ v \end{bmatrix} = -I_t$$

$$+ f_x \cdot u + f_y \cdot v + f_t = 0 \Rightarrow \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_y I_x & \sum I_y I_y \end{bmatrix} \cdot \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$$+ \Delta p = H^{-1} \cdot \sum_x \left(\nabla_I \frac{\partial w}{\partial p} \right)^T [T(x) - I(w(x, p))]$$

$$H = \sum_x \left(\nabla_I \frac{\partial w}{\partial p} \right)^T \left(\nabla_I \frac{\partial w}{\partial p} \right)$$

2) Bayesian Filtering

$$\rightarrow x_t \sim \mathcal{N}(D_t x_{t-1}, \Sigma_{dt}), \quad y_t \sim \mathcal{N}(M_t x_t, \Sigma_{mt})$$

$$\begin{aligned} \textcircled{*} \quad \bar{x}_t &= D_t x_{t-1}^+ \\ \bar{\Sigma}_t &= D_t \Sigma_{t-1}^+ D_t^\top + \Sigma_{dt} \end{aligned} \quad \left| \quad \begin{aligned} K_t &= \bar{\Sigma}_t M_t^\top (M_t \bar{\Sigma}_t M_t^\top + \Sigma_{mt})^{-1} \\ x_t^+ &= \bar{x}_t + K_t (y_t - M_t \bar{x}_t) \\ \bar{\Sigma}_t^+ &= (I - K_t M_t) \cdot \bar{\Sigma}_t^- \end{aligned} \right.$$

$$\rightarrow G_t = \frac{\partial g(x)}{\partial x} \Big|_{x=x_{t-1}^+}, \quad H_t = \frac{\partial h(x)}{\partial x} \Big|_{x=\bar{x}_t}$$

$$\rightarrow fme \left[\{x_t^i, w_t^i\}_{i=1}^N \right] = SIS \left[\{x_{t-1}^i, w_{t-1}^i\}_{i=1}^N, y_t \right]$$

$$\eta = 0$$

for $i = 1:N$

$$\begin{cases} x_t^i \approx q(x_t | x_{t-1}^i, y_t) \\ w_t^i = w_{t-1}^i \cdot p(y_t | x_t^i) \\ n = \eta + w_t^i \end{cases}$$

end

$$w_t = w_t / n$$

$$\rightarrow N_{eff} = \frac{1}{\sum_{i=1}^n (w_i^i)^2}$$

 uniform $N_{eff} = N$
degenerated $N_{eff} = 1$

\rightarrow Arulampalam's Resampling Algorithm

$$\left[\{x_k^{j*}, w_k^j, i^{j*}\}_{j=1}^{N_s} \right] = \text{Resample} \left[\{x_k^i, w_k^i\}_{i=1}^{N_s} \right]$$

3, Multi-Object Tracking

$$+) \quad \left\{ x_1^{(k)}, \dots, x_{N_k}^{(k)} \right\} \quad \text{and} \quad \left\{ y_1^{(k)}, \dots, y_{M_k}^{(k)} \right\}$$

$$z_l^{(k)} = j \quad \text{if} \quad y_j^{(k)} \leftarrow x_l$$

+ Nearest Neighbor filter:

$$z_l^{(k)} = \arg \min \left(x_{p,l}^{(k)} - y_j^{(k)} \right)^T \left(x_{p,l}^{(k)} - y_j^{(k)} \right)$$

$$\text{better}_j = \arg \min_j v_{j,l}^{(k)} \cdot (\sum_{p,l}^{(k)})^{-1} \cdot v_{j,l}^{(k)}$$

with innovation $v_{j,l}^{(k)} = y_j^{(k)} - x_{p,l}^{(k)}$

→ Track-splitting filter

$$L(\theta_{k,l}) = \prod_{j=1}^{N_k} p(z_{ij,l} | z_{(j-1),l}, \theta_{k,l})$$

$$\lambda_c(k) = \lambda_c(k-1) + v_{i_k, l}^{(k)\top} (\sum_l)^{-1} v_{c_k, l}^{(k)}$$

$$\Rightarrow \text{MHT} = p(\Omega_j^{(k)} | Y^{(k)}) = \eta \cdot p(Y^{(k)} | Z_j^{(k)}, \Omega_{pj}^{(k-1)}) \cdot p(Z_j^{(k)} | \Omega_{pj}^{(k-1)}) \cdot p(\Omega_{pj}^{(k-1)})$$

$$\Rightarrow \text{LAP: } \sum_{i=1}^N \sum_{j=1}^M w_{ij} z_{ij}, \quad \sum_{j=1}^N z_{ij} = 1, \quad \sum_{i=1}^M z_{ij} = 1$$

+), Network flow optimization

$$\begin{array}{l} \text{flow} \\ \text{conservation} \\ \hline \text{edges} \\ \text{capacities} \end{array} \quad f_{in,i} + \sum_j f_{j,i} = f_i = f_{out,i} + \sum_j f_{i,j} \quad \forall i$$

$$f_i \leq 1$$

4) Visual Odometry:

- Cameras:
 - Monocular cameras $I_{0:t} = \{I_0, \dots, I_t\}$
 - Stereo cameras $I_{0:t}^l = \{I_0^l, \dots, I_t^l\}$
 - Active RGB-D camera $I_{0:t}^r = \{I_0^r, \dots, I_t^r\}$
 - Active RGB-D camera $I_{0:t} = \{I_0, \dots, I_t\}$
- Direct Methods $E(\xi) = \int_{u \in \Omega} |I_1(u) - I_2(w(u, \xi))| du$
- Indirect $E(\xi) = \sum_i |y_{1,i} - w(y_{2,i}, \xi)|$

~~+) 2D to 2D~~ $E(T_t^{t+1}, X) = \sum_{i=1}^N \|\bar{y}_{t,i} - \pi(\bar{x}_i)\|_2^2 + \|\bar{y}_{t+1,i} - \pi(T_t^{t+1} \cdot \bar{x}_t)\|_2^2$

Error 2D to 3D $E(T_t) = \sum_{i=1}^N \|y_{t,i} - \pi(T_t, \bar{x}_i)\|_2^2$

3D to 3D $E(T_t^{t+1}) = \sum_{i=1}^N \|\bar{x}_{t+1,i} - T_t^{t+1} \cdot \bar{x}_{t,i}\|_2^2$

+) Eight point algorithm:

$$\bar{y}_i \cdot \bar{y}'_i = 0 = a_i E_s \rightarrow A \cdot E_s = 0$$

$$\rightarrow A = U_A \Sigma_A V_A^\top \rightarrow E_s : \text{last column of } V_A$$

$$E = U \cdot \text{diag}(6_1, 6_2, 6_3) \cdot V^\top$$

$$\rightarrow E = U \cdot \text{diag}(1, 1, 0) \cdot V^\top$$

$$\begin{cases} R = U \cdot R_z^\top (\pm \frac{\pi}{2}) \cdot V^\top \\ T = U \cdot R_z (\pm \frac{\pi}{2}) \text{diag}(1, 1, 0) \cdot U^\top \end{cases}$$

+) Relative scale recovery

$$r_{ij} = \frac{\|\bar{x}_{t+2,t-i,i} - \bar{x}_{t+2,t-1,j}\|_2}{\|\bar{x}_{t-1,t,i} - \bar{x}_{t-1,t,j}\|_2}$$

$$\rightarrow \text{DLT: } \tilde{y}_i = [x'_i \ y'_i \ w'_i]^T, \quad \tilde{x}_i = [x_i \ y_i \ z_i \ w_i]^T$$

$$\Leftrightarrow 2 \text{ constraints: } \begin{bmatrix} 0 & -w'_i \tilde{x}'_i & y'_i \tilde{x}'_i \\ w'_i \tilde{x}'_i & 0 & -x'_i \tilde{x}'_i \end{bmatrix} \cdot \begin{bmatrix} P_1^T \\ P_2^T \\ P_3^T \end{bmatrix} = 0$$

$$\text{from } \tilde{y}_i \times (P \tilde{x}_i) = 0$$

\rightarrow Closed-form solution of 3D \rightarrow 3D

$$\mu_{t+1} = \frac{1}{N} \sum_{i=1}^N x_{t+1,i}, \quad \mu_t = \frac{1}{N} \sum_{i=1}^N x_{t,i}$$

$$\text{Set } A = \sum_{i=1}^N (x_{t+1} - \mu_{t+1})(x_t - \mu_t)^T, \quad A = U \cdot \Sigma \cdot V^T$$

$$\rightarrow R_{t+1}^+ = V \cdot U^T \quad \& \quad t_{t+1}^+ = \mu_t - R_{t+1}^+ \mu_{t+1}$$

\rightarrow Direct Methods

$$\text{photometric error } I_i(y) = I_2(\pi(T(\xi) \cdot Z_1(y), \bar{y}))$$

$$\text{Geometric error } [T(\xi) Z_1(y) \bar{y}]_2 = Z_2(\pi(T(\xi) \cdot Z_1(y), \bar{y}))$$

$$E(\xi) = \sum_{y \in \Omega} \frac{r(y, \xi)^2}{\sigma_I^2}, \quad r(y, \xi) = I_i(y) - I_2(\pi(T(\xi) \cdot Z_1(y), \bar{y}))$$

minimize the residual

$$= r(\xi)^T \cdot W \cdot r(\xi)$$

$$\begin{aligned} \rightarrow \text{Gauss-Newton: } s_{i+1} &= \xi_i - [\nabla_{\xi}^2 E(\xi_i)]^{-1} \nabla_{\xi} E(\xi_i) \\ &= \xi_i - H_i^{-1} \cdot J_i^T W r(\xi_i) \end{aligned}$$

$$\rightarrow \text{Levenberg-Maquardt: } s_{i+1} = \xi_i - (H_i + \lambda I)^{-1} J_i^T W r(\xi_i)$$

$$\rightarrow \text{Lie algebra: } \xi = \begin{bmatrix} \omega \\ v \end{bmatrix} \in \mathbb{R}^6, \quad \hat{\xi} = \begin{bmatrix} \hat{\omega} & v \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{4 \times 4}$$

$$T = \exp(\hat{\xi}) = \begin{bmatrix} \exp(\hat{\omega}) & Av \\ 0 & 1 \end{bmatrix}$$

$$T = \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 4}$$

$$\hat{\xi} = \log(T) = \begin{bmatrix} \log R & A^{-1}t \\ 0 & 0 \end{bmatrix}$$

+)Contrastive Loss:

$$\checkmark \|f(x) - f(x_+)\| \rightarrow 0$$

$$\|f(x) - f(x_-)\| \geq m$$

+)Triplet loss:

$$\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2$$