

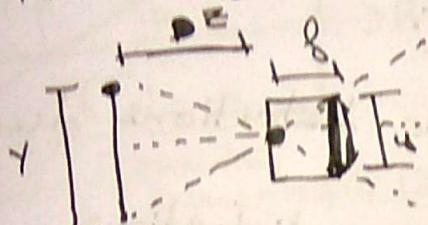
②

## Human Eye:

- Retina  $\rightarrow$  Fovea (High sensitivity)

## \* Pinhole Camera: (Everything is in focus)

- Focal length =  $f$ , Depth =  $z$



Inverted image formed.

$$\frac{y}{z} = \frac{u}{f}, \quad u = f - \frac{y}{z}$$

$\Rightarrow$  Issue is that a very small amount of light enters through the pinhole. (Long exposure)

- If pinhole size is large, we get blurry images since light from multiple regions are summed up in each point within the pinhole camera.

(Too small pinhole  $\rightarrow$  Diffraction issue)

## \* Camera with Lens: (Bidirectional inversion)

\* Camera with Lens: (Bidirectional inversion)

- Lens helps focus light in a range onto a single point. The aperture is large and yet we get a sharp image.

To get a sharp image.  $f$  = Focal length of the lens.

$$\frac{1}{f} = \frac{1}{d_o} + \frac{1}{d_i}$$

$d_o$  = Object to lens.  
 $d_i$  = Lens to screen.

$\Rightarrow$  Depth of Focus (DOF) = Region of depth of object which appears to be in focus keeping camera at the same point.

Depth of Field

Note: Lens-Aperture Camera: (SLR cameras)

- Lens is present along with a variable hole size pin hole behind it.

Note: Zooming  $\rightarrow$  f of lens changes or multiple lenses.  
Focusing  $\rightarrow$  Moving position of lens.

### \* Distortions due to lenses:

1) \* Pincushion & Barrel geometric distortions due to bad lenses.  
ex: Fish eye lenses etc have these distortions

### 2) Lens Flare

- When there are multiple lenses, some light gets reflected back and light bounces between lenses / inside lens if its a single lens.

### 3) Chromatic Aberration:

- The focal point of different light R, G, B is not the same.  
 $\Rightarrow$  White light falling on lens, R, G, B gets split and we can see the light separating out.

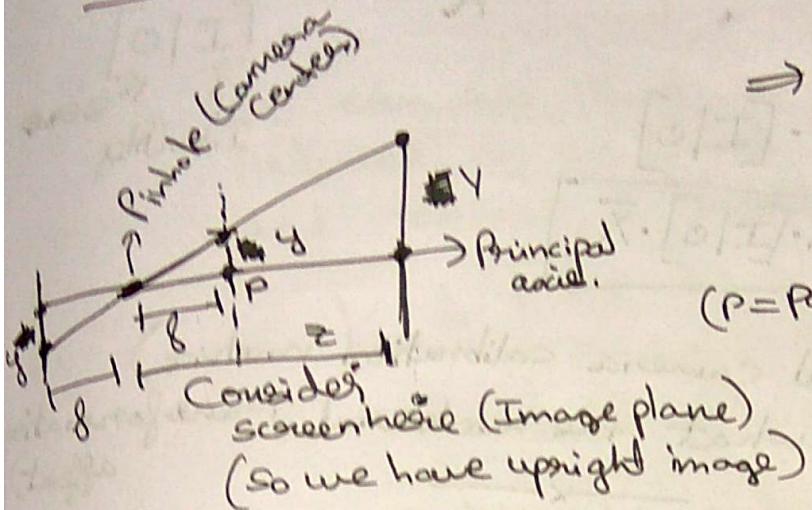
Note: CCD & CMOS sensor arrays for capturing light.

Note: Sampling:  $f[m, n] = f(x, y) \sum_{n, m} g(x - nM_1, y - mM_2)$   
\*  $M_1, M_2$  = Sampling distances  
 $\Rightarrow g$ 's all together called 2D impulse train.  
• we then perform quantization on  $f[m, n]$ .

Note: Test on thursday (12<sup>th</sup>), 2.1-2.4, 3.1, 3.2.1,  
~~3.2.2~~ (H&Z) 3.2.2

③

## Geometry & Camera Matrices:



⇒ Consider screen ahead of pinhole for mathematical simplicity  
(P = Principal point)

- In 2D,  $\bar{x} = \frac{fx}{z}$ ,  $\bar{y} = \frac{fy}{z}$

⇒  $(\bar{x}, \bar{y})$  are the coordinates on screen.  
 $(x, y)$  are world points.

- $\begin{bmatrix} \bar{x} \\ \bar{y} \end{bmatrix} = \begin{bmatrix} ? & ? \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix}$

We cannot form this camera matrix without  $x, y, z$  in it, so we move to a different formulation

### \* Homogeneous Coordinates:

- $\begin{bmatrix} x \\ y \\ w \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$

$\bar{x}$        $\xrightarrow{\text{Camera matrix}}$        $\bar{X}_c$

(Unit of  $f$  gives unit of  $x, y$ )  
(ex: Pixels)

- We interpret coordinates as  $(\frac{x}{w}, \frac{y}{w})$

⇒ Helps us model the nonlinearity with ease.

(Note:  $X, Y, Z$  are <sup>word point coords</sup> world coordinate axis of camera - (camera coordinates))

$$\Rightarrow P = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}_{3 \times 4} = \underbrace{\begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix}_{3 \times 3}}_K \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}_{3 \times 4}}_{\begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix}} \text{ Identity}$$

$$\therefore P = K \cdot [I|0]$$

$$\therefore \boxed{\bar{x} = K \cdot [I|0] \cdot \bar{x}_c}$$

\* ( $K$  = Internal camera calibration matrix)

( $[I|0]$  is what the external transformation affect)

Note: If pixel distances on screen are not same, change  $f$ 's in  $K$ . ( $f$ 's units are in pixels), scale each  $f$  based on the pixel pitch of the corresponding axis. (pixel distances) (Internal property) ( $x$  &  $y$ )

\* Internal Properties:

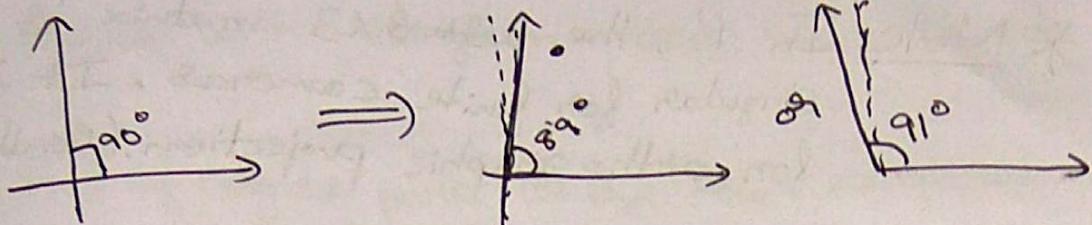
a) We however want our image center  $(0,0)$  to be at the bottom left corner & all coordinates  $\bar{x}$  are true.

→ We do this by making

$$K = \begin{bmatrix} fx_0 & x_0 \\ 0 & fy_0 & y_0 \\ 0 & 0 & 1 \end{bmatrix}$$

$(x_0, y_0) = (\text{Height}/2, \text{Width}/2)$  of screen  
(This is again internal to the camera)

b) Some cameras have issues where axes are not at a  $90^\circ$ . There is a skew. So to correct for the skew,  $x$  coordinate needs to be adjusted



- We can see that  $\alpha$  coordinate gets reduced, so we would have to modify and increase it based on skew ( $\gamma$  coord)

$$\therefore K = \begin{bmatrix} \alpha_x & 0 & \alpha_x \\ 0 & \alpha_y & \alpha_y \\ 0 & 0 & 1 \end{bmatrix}$$

(Assume  $\alpha$  axis is correct)  
(convention)

$K$  = Upper triangular matrix with 5 degrees of freedom.

### \* External Transformations:

- Translation, rotation etc is done by using another ~~matrix~~ matrix as required.

$$\Rightarrow \bar{x}_c = \begin{bmatrix} R & -RC \\ 0 & 1 \end{bmatrix} \bar{x}_w$$

To get camera coordinates, we use another transformation matrix.

$R$  =  $3 \times 3$  rotation matrix

$-RC$  =  $3 \times 1$  vector

(Translation & then rotation)  $\rightarrow$  Graphics notes for details.

Finally,

$$\Rightarrow \bar{x} = K[1|0]\bar{x}_c = K[R|-RC]\bar{x}_w$$

$$\Rightarrow \bar{x} = P\bar{x}_w = \begin{bmatrix} K|R| & -KRC \\ 0 & 1 \end{bmatrix} \bar{x}_w = [M|P_4] \bar{x}_w$$

Note:

A few representations where  $P_i = 4 \times 1$  column

$$\left[ P_1 \ P_2 \ P_3 \right]^T \quad \text{or} \quad \left\{ \begin{array}{l} M = [P_1 \ | \ P_2 \ | \ P_3] \\ P_i = \text{column vectors} \end{array} \right.$$

so it becomes  $3 \times 4$  matrix.

\* Note: In P, the left  $3 \times 3$  matrix is non singular for finite cameras. It is singular for orthographic projection. (Parallel projection)

④

## Camera Calibration: (Ch 2.1 to 2.3)

- Camera center does not have image on screen.
- \* Camera center = world point that gives  $(0,0,0)$  when projected.

$$P = [p_1 \ p_2 \ p_3 \ p_4]$$

\*  $\Rightarrow$  Columns  $p_1, p_2 \ \& \ p_3$  are the ~~screen~~ coordinates of X, Y, Z vanishing points in the world.

$$\text{ex: } [p_1 \ p_2 \ p_3 \ p_4] \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow X \text{ vanishing point}$$

$$\qquad\qquad\qquad \textcircled{p}_1 \rightarrow \text{Screen coordinates.}$$

$\Rightarrow$  Column  $p_4$  is the image of world origin.

$$\Rightarrow P = [p_1 \ p_2 \ p_3]^T,$$

\*  $\hookrightarrow p_3$  = Principal plane = Parallel plane to screen passing through camera origin. ( $z=0$  always)

$$(a, b, c, d) \Rightarrow (a, b, c, 0)$$

; Passes

through origin.

$\Rightarrow (a, b, c) \Rightarrow$  Normal of plane/screen.

— x —

## Camera Calibration: (Ch 2.1, 2.2, 2.3)

- 1) 3D Reference object based calibration  
 $(\bar{x} = P\bar{X}_w)$  (Use known  $\bar{x}, \bar{X}_w$ ) (12 unknowns)
- 2) Precisely moving plane (Physical measurements)

\*3) Plane with unknown motion (Checkered pattern on plane)

4) Set of collinear points where the line passes through fixed points.

5) Self calibration (rigid static world & point correspondence)

→ Going from ① to ⑤ precision drops, but easy to perform.

— X —

\*1) 3D reference object based: (Linear equations)

$$x_i = P_{11}X_i + P_{12}Y_i + P_{13}Z_i + P_{14}$$

$$y_i = P_{21}X_i + P_{22}Y_i + P_{23}Z_i + P_{24}$$

$$z_i = P_{31}X_i + P_{32}Y_i + P_{33}Z_i + P_{34}$$

→ We know  $(x_i, y_i, z_i)$  &  $\left(\frac{x_i}{w_i}, \frac{y_i}{w_i}\right)$

for each pair of world/screen points.

$$\frac{x_i}{w_i} = u_i = \frac{P_1X_i + P_{14}}{P_3X_i + P_{34}}$$

$$P_1X_i - u_iP_{34} + P_{14} - u_iP_{34} = 0 \rightarrow 1$$

$$\frac{y_i}{w_i} = v_i = \frac{P_2X_i + P_{24}}{P_3X_i + P_{34}}$$

$$P_2X_i + P_{24} - v_iP_3X_i = 0 \rightarrow 2$$

Using this we get camera matrix. (Once we find approx P)

→ We now need to decompose this

We can define ①, ② for each point in a matrix where  $G = 2n \times 2$

$$G \cdot p = 0$$

(where  $|p|=12$ )

$G =$  rows for each point

2n rows for n points

From this get Rft.

$$P = K[R|t], K = \begin{bmatrix} \alpha & \beta & u_0 \\ 0 & \gamma & v_0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$P = [M | p_4], M = KR, p_4 = Kt$$

→ We know, M & p4, we need to find K

In rotation  $R^T \Rightarrow$  Opposite rotation,  $\therefore R^T R = I$

$$\text{Now, } MMT = KRR^TK^T = KKT. \text{ We can get K from this.}$$

• We then divide K to make bottom right element 1 since that scales images.

This scales images  
Always an issue (cannot be found)

## $\Rightarrow$ Refining $P$ : (Non linear optimizations)

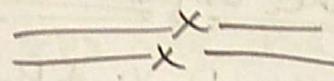
- \* The  $P$  we calculate is an approximation, so we minimize w.r.t.:

$$\min_P \sum_i \|m_i - \phi(P, N_i)\|^2$$

Original position  
on image ( $m_i, N_i$ )  
Points projected by  $P$

$\Rightarrow$  Perform optimizations like gradient descent.  
(Useful to remove noise etc. that might be present in the image) (sort of)

⑥

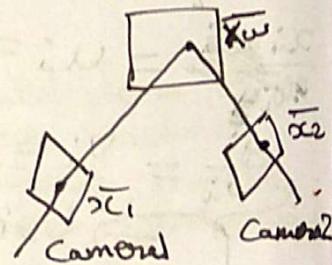


## 2-View Geometry: (2D World)

- \* Assume world is planar & the XY plane. we can transform world plane to XY if its not originally on the XY plane.

$$\therefore \text{Coordinates} = \begin{bmatrix} X \\ Y \\ 0 \\ 1 \end{bmatrix}$$

$$\bar{x} = [P_1 \ P_2 \ P_3 \ P_4] \begin{bmatrix} X \\ Y \\ 0 \\ 1 \end{bmatrix}$$



$$\bar{x} = [P_1 \ P_2 \ P_4] \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}$$

Non singular (Rank 3) matrix

$$\bar{x}_1 = H_{3 \times 3}' \cdot \bar{x}_w = H_1 \cdot \bar{x}_w$$

$$\bar{x}_2 = H_{3 \times 3}'' \cdot \bar{x}_w = H_2 \cdot \bar{x}_w$$

$$\bar{x}_1 = H_1 \cdot \bar{x}_w = H_1 \cdot H_2^{-1} \bar{x}_2$$

$$\therefore \bar{x}_1 = H_{12} \bar{x}_2 \quad (H_{12} = H_1 \cdot H_2^{-1})$$

$$\bar{x}_2 = H_{21} \bar{x}_1 \quad (H_{12} = H_2^{-1})$$

- If we know  $H$  and one of the images, we can compute the other image automatically (similar to interpolation on input for each output pixel)

\* ( $H = \text{Homography matrix}$ )

$\Rightarrow H$  is a  $3 \times 3$  matrix with 8 free variables.  
So to estimate it we need 4 point correspondences.

\* 2-View: (Same Camera Center) (3D World)

$$*\bar{x}_1 = K_1 \cdot R_1 [I | -C] \bar{X}_w \quad (K \text{ can change so zooming etc is allowed})$$

$$\bar{x}_2 = K_2 \cdot R_2 [I | -C] \bar{X}_w$$

- Assumption is that camera center does not change. We only rotate it

$\Rightarrow K_1 R_1$  &  $K_2 R_2$  are full rank matrices

$$\bar{x}_2 = K_2 \cdot R_2 \cdot (K_1 R_1)^{-1} \cdot \bar{x}_1$$

$$\bar{x}_2 = H_{21} \cdot \bar{x}_1$$

$3 \times 3$  Homography matrix.

$\Rightarrow$  Similar to previous case, here we assume world is planar but camera center can change. Here world can be 3D but camera center should not change

Note: Other homographies are also possible, parallel lines etc.

Note: Panorama — Stitch images with homography  
Use merging techniques etc

## \*Feature Matching | Description | Detection:

- Feature detection, description & matching.
  - SIFT, Harris corners, SURF etc

⑦ Harris Corner Detection: (Feature Detection) (DIP Notes)

- Edges → Illumination edges, surface color edges  
→ Depth edges (example in slides)  
→ Surface normal edges.
  - Depth & illumination edges depend on illumination of the scene & viewpoint.
  - DoG → Difference of Gaussians operator.

$$D(G(I)) \Rightarrow (D(G)I)$$

$\downarrow$   
 $D \circ G(I) \rightarrow$  Edge detection

$$\bullet E(x,y) = \sum_{u,v} (I(x+u,y+v) - I(x,y))^2$$

$$I(x+u, y+v) = I(x, y) + \frac{\partial I}{\partial x}(x, y) \cdot u + \frac{\partial I}{\partial y}(x, y) \cdot v$$

$$\Rightarrow I(x+u, y+v) - I(x, y) = \begin{bmatrix} I_{xx}^{(xy)} & I_{xy}^{(xy)} \\ I_{yx}^{(xy)} & I_{yy}^{(xy)} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

$$\Rightarrow E(X^4) = \sum_{\text{even}} [(Ex I_4) \begin{bmatrix} u \\ v \end{bmatrix}]^2$$

$$= \sum_{\text{all } u, v} [u \vee] \begin{bmatrix} I_x \\ I_y \end{bmatrix} [I_x \ I_y] \begin{bmatrix} u \\ v \end{bmatrix}$$

$$[uv] \begin{bmatrix} \Sigma Ix^2 & \Sigma Ixy \\ \Sigma Ixy & \Sigma Iy^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \sum_{(u,v) \in \text{new}} [uv] \begin{bmatrix} Ix^2 & Ixy \\ Ixy & Iy^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

(As we vary  $x, y$ )  $\rightarrow$  We want to identify minimum of  $E(x, y)$  & then maximize it for all  $(u, v)$  that is Maximize the direction  $(u)$  which makes  $E(u, v) \rightarrow \min$

i.e. Maximize the direction ( $u, y$ ) giving min change  $\rightarrow$  min eigenvalue

⇒ Eigen vectors & values for H give us the direction & amount of change in those two major directions.

$$\Rightarrow Hx_+ = \lambda_+ x_+ \quad , \quad Hx_- = \lambda_- x_-$$

Value      ↓      Max charge vector      Value      ↓      Min charge vector  
 $\lambda_+$       ↓      value =  $\lambda_-$

- We want min charge vector to be high. i.e. Both eigen values should be large. ( $\because \lambda_- < \lambda_+$ ) (easier to check)

$\Rightarrow \lambda_+$  = High for edges,  $\lambda_-$  is also high for corners.

- $f = \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2} = \frac{\text{Determinant}(H)}{\text{Trace}(H)} \rightarrow \text{Harris operator}$

If f is large then select point.

$\Rightarrow$  We then apply local thresholding to select single points in a region so that we pick only the local maxima.

(8)

— x —

### Feature Description:

#### \*1) Multiscale Oriented Patches: (MOPS)

- Take multiple windows of various sizes around feature points & store them as the descriptors.
- For orientation of windows, use the direction of maximum gradient as the base orientation.
- 8x8 patches at 5 scales (scale image)
- Bias gain normalization to remove mean & divide by variance  $\rightarrow$  illumination invariance.

Dimension of 8x8 =  $8 \times 8 \times k$   $\rightarrow$   $k$  features.

## 2) Scale Invariant Feature Transform (SIFT)

- Illumination Invariance — Difference based metrics  
(Difference between points etc)

- Scale Invariance →
  - Pyramids (resized images)
  - Difference of Gaussians  
(Subtract levels from gaussian pyramid without resize)  
↑  
Different blurs.

(Blurring is sort of like continuous scaling)

In DOGs, after certain amount of blur, we lose as much info as we lose on resize, so to save space we resize the image & then proceed from there.

- Rotation Invariance → Histogram of gradient directions for all pixels in window. Base reference direction is average direction of all gradients.

- ⇒ Create  $4 \times 4$  regions in window & get histograms of orientations in each region, the gradients are given weighted contributions.
- Gradients of pixels are slightly averaged, so that different region spanning edges are taken care of.

- ⇒ Feature detection → Max in multiple levels of difference of gaussians (DOGs)

- Feature size:  $16 \times 8 = 128$  dimensions.  
↑  
Regions Histogram with 8 elements.

## \* Feature Matching:

- Direct matching  $\rightarrow n^2$  comparisons of features  
(Too slow) (If  $n$  points are present)

$\Rightarrow$  Hashing of descriptors  $\rightarrow$  Matches are directly found

→ Hard to hash though since features are similar but not identical.

$\Rightarrow$  Nearest neighbours  $\rightarrow$  k-d trees etc-  
\* (Approximate algorithms that are fast)

- Outliers can be dropped using thresholds on match similarities.

## \* RANSAC:

- Select 4 feature pairs at random (from matches)
- Compute homography
- Compute inliers where  $SSD(p_i^1, H p_i^2) < \epsilon$
- Repeat this process and the one where we get maximum inliers, keep them and discard outlier matches.
- Finally compute  $H^1$  on all inliers. This is the final homography matrix to align images based on features -

## ⑨ Geometry (3D world + Camera center Shift allowed.)

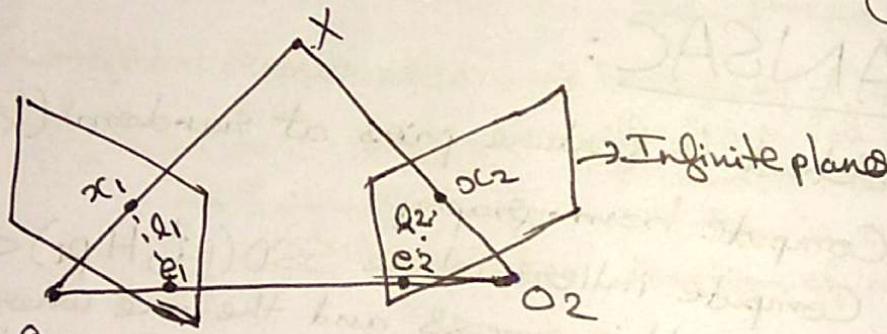
- Consider 2 cameras, you know the rotation & translation between these 2 cameras.

⇒ If a point  $X$  in 3D world appears at  $x_1$  on camera  $C_1$ , with origin  $O_1$  &  $X$  appears at  $x_2$  on camera  $C_2$  with origin  $O_2$ .

\*  $X$  lies on intersection of  $\overline{O_1x_1}$  &  $\overline{O_2x_2}$ .  
 (Essentially we are identifying the depth as well) → Reconstruction of 3D world from images. (Slight errors in  $x_1$  &  $x_2$  give huge errors)  
 (Noise in image etc)

### Epipolar Geometry:

\* \*



- $e_1$  &  $e_2$  are epipoles → They don't change if we change  $X$ .

⇒ All world points mapping to  $x_1$  in  $I_1$  (pre image of  $x_1$ ) map to a line  $l_2$  in  $I_2$ .  $l_2$  is essentially the line on the image plane of  $I_2$ .  $l_1$  &  $l_2$  are called epipolar lines.

- Image of  $O_1$  in  $I_2 = e_2$  & image of  $O_2$  in  $I_1 = e_1$  are epipoles.

Epipolar constraint: • The plane of  $O_1e_2$  = Epipolar plane

- \* •  $\lambda_1 x_1 = X \quad \} \text{since } x_1 \text{ & } X \text{ lie on same line}$   
 (Assume origin at  $O_1$ )

$$\Rightarrow \lambda_2 x_2 = RX + T, \text{ where } R, T \text{ are rotation \& translation of O}_2 \text{ wrt O}_1.$$

$$= R(\lambda_1 x_1) + T$$

- Additions are inconvenient, so we try get rid of it with  $\hat{T}$

$$\cdot \hat{T} \lambda_2 x_2 = \hat{T} R \lambda_1 x_1 + \hat{T}, \hat{T} = \begin{bmatrix} 0 & -T_z & T_y \\ T_z & 0 & -T_x \\ T_y & T_x & 0 \end{bmatrix}$$

$$\Rightarrow \lambda_2 x_2^T \hat{T} x_2 = \lambda_1 x_2^T \hat{T} R x_1$$

$$\cdot x_2^T \hat{T} R x_1 = 0$$

$$\therefore x_2^T \hat{T} x_2 = 0$$

we can find this as above

(Cross product,  $90^\circ$  to 3D plane  
(Dot product with  $90^\circ = 0$ ))

$$\therefore x_2^T E x_1 = 0 \text{ & } x_1^T E x_2 = 0$$

If we repeat the process with  $O_2$  as origin.

$$\cdot \text{Now } AXA = 0$$
 ~~$\cdot A^T A = 0$~~ 

$$\therefore A^T A = \hat{A} \hat{A} = 0$$

$$\cdot \text{We know } E = \hat{T} R \quad (\text{Since we know how our cameras are configured})$$

$$\Rightarrow \text{Now if we know } x_2, \text{ since we know } E$$

we can get  $x_2^T E = 1 \times 3$  vector and,  $x_1$  lies on this line.

$$0 = (x_2^T E) x_1 \Rightarrow \text{Epipolar line } l_1$$

$$0 = (x_1^T E) x_2 \Rightarrow \text{Epipolar line } l_2, E = \text{Essential matrix}$$

Similarly,  $0 = (x_1^T E) x_1$

Note: If internal parameters of  $C_1$  &  $C_2$  are different ( $K_1$  &  $K_2$ )

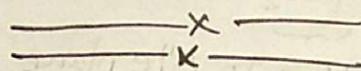
\* then,  $x_1 = K_1 X, x_2 = K_2 X$

we get,  $x_2^T K_2^{-T} \hat{T} R K_1^{-1} x_1 = 0$

$F \rightarrow$  Fundamental matrix

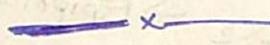
Note: The reason we move to epipolar geometry \* is because finding  $x_2$  corresponding to  $x_1$  we need not search the entire image but just on  $l_2$ .

- Note: • To find  $X$ , we need  $x_1$  &  $x_2$ . To get corresponding match of  $x_1$  using SIFT etc., we need to search entire image for feature matching but from epipolar constraint, given  $E$  &  $x_1$  we only need to search on  $\ell_2$  for our match point  $x_2$ . (Reduces search space).
- Once we have  $x_1$  &  $x_2$ , we can find  $X$  as intersection of  $O_1x_1$  &  $O_2x_2$ .



## ⑩ Multiview Reconstruction:

- Binocular stereo — Two cams with known intrinsic and extrinsic params
- Multiview stereo — Multiple known cameras
- Structure from motion — Multiple cameras & project with correspondences. Retrieve 3D coords & camera matrices  $C_i$ .

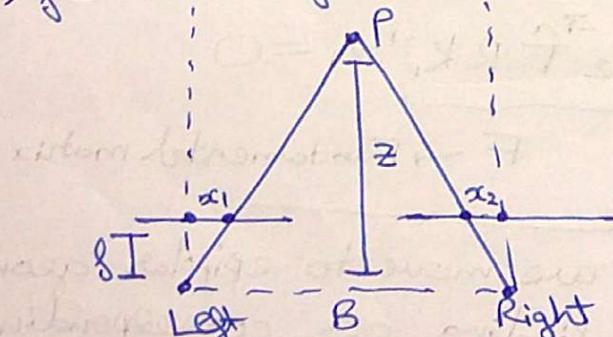


## \* Binocular Stereo:

- Assume both cameras as similar (internal params)

$$\Rightarrow f = \text{Focal length}, z = \text{Depth}$$

$$\bullet \frac{x_1 - x_2}{f} = \frac{B}{z} \quad \begin{matrix} \text{-ve direction} \\ \text{from similar triangles.} \end{matrix}$$



$$\Rightarrow z = \frac{Bf}{x_1 - x_2}$$

$$* \rightarrow \underbrace{x_1 - x_2} = \text{Stereo disparity} = \frac{Bd}{z} = \Delta$$

or Parallax. ( $\Delta$ )

↳  $A_2 \equiv \Gamma$ , parallel

Same world point in both images.

→ Larger baseline gives better estimates but hard to get matches etc. → Since we want to avoid  $\Delta$  zero ( $\Delta \approx 0$ )

— X —

Note: Matching patches (vectors)

$$* \cdot \text{Correlation} \rightarrow \frac{\mathbf{v}^T \mathbf{v}}{\sqrt{\mathbf{v}^T \mathbf{v}} \cdot \sqrt{\mathbf{v}^T \mathbf{v}}}$$

blw  $\mathbf{v} \otimes \mathbf{v}^T$

• Normalized correlation → Subtract mean from vectors & correlate.

— X —

Note: Constraints on matching:

- \* • Epipolar lines • Colour constancy
- Uniqueness — one point on one image only matches 1 other point. Not always true.
- \* • Ordering & Monotonicity: Points appearing in an order will remain same for small camera rotations or translations.

— X —

Note: In Binocular Stereo, if one of the cameras (Rectification) is rotated as well instead of just translation,   
 → \* \* image  $I_2$  convert to  $I'_2$  such that  $I'_2$  is

Camera calibration for stereo  
↓ Estimate  $H_1$  &  $H_2$

rotation of  $I_2$  (Affine transform) • Compute & apply the rotation homography & you can now perform binocular stereo.

→ (If translations along both axis b/w cameras perform 2 homographic transformations and make same parallel to line joining their centers)

Note: Dense pixel-pixel matching — DP Algo. (Slides)  
★ (Not just matches b/w certain feature points)  
Globally optimal solution. → Some issues exist as well. — X —

### Structured Lighting:

- Project patterned light onto the world & then captures images so that features & match detection becomes simpler.

(Kinect uses this)

— X —

Note: Things trivial for humans are hard for computers → Grouping (Segmentation)

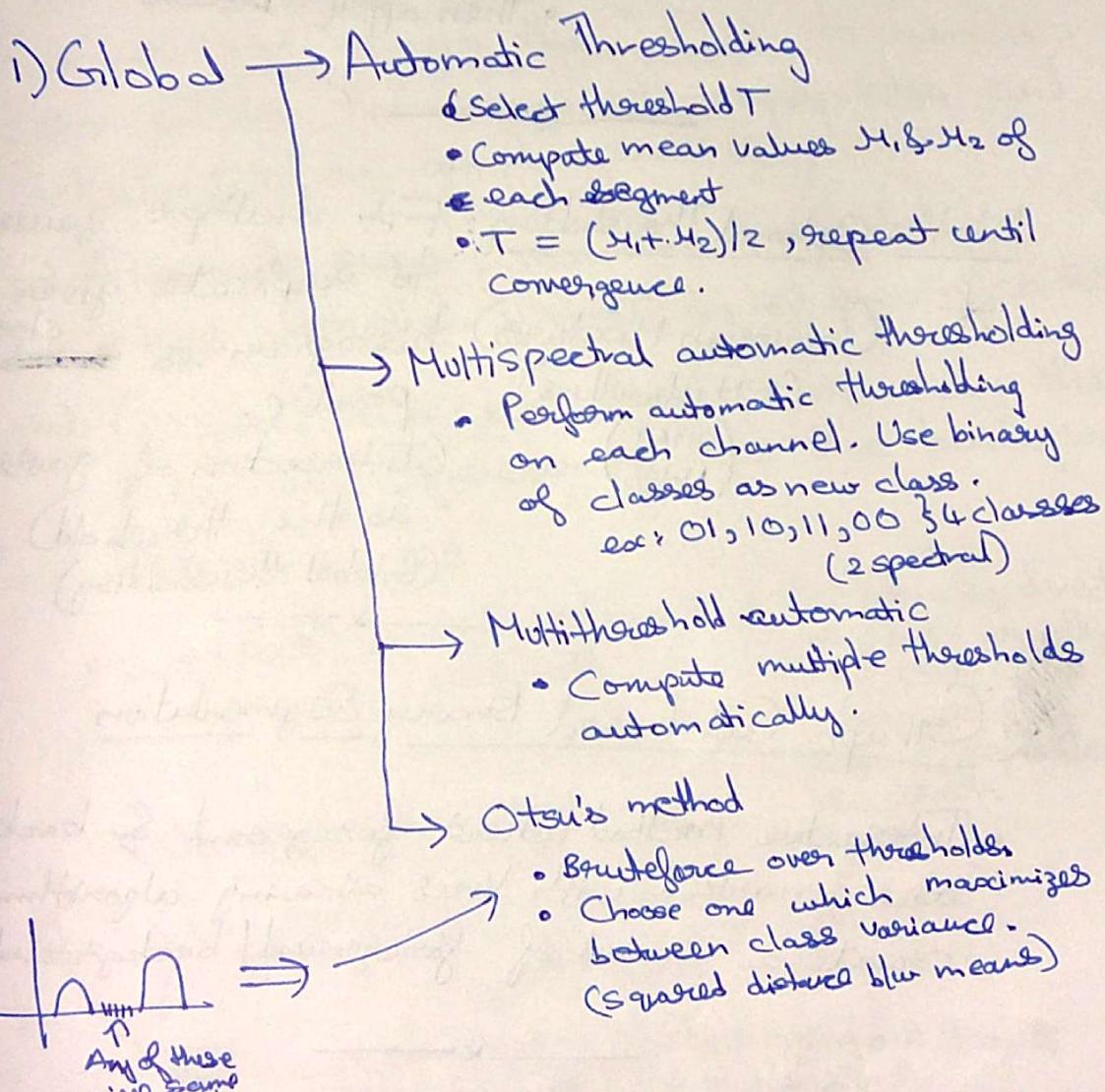
Things trivial for computers are hard for humans → Measuring more

# Image Segmentation:

- Classification based → Each pixel is classified
- Region based → Region growing & splitting
- Boundary based → Edges
- Motion based

## \* Thresholding: (DIP Notes)

- Global → Threshold selected via histogram.
- Adaptive (Local) →



~~Adaptive thresholding~~

## 2) Adaptive Thresholding

- Use mean in a local region.  
(Generally bad)
- \* → Mean C
  - Uniform region set all as background  
otherwise perform mean thresholding  
as above.
- \* → Chow-Kaneko
  - Apply mean filters multiple  
times to get background image  
and subtract from original -
  - Then apply Otsu etc

— X —

Note: Optimal Thresholding: Fit multiple gaussians  
to replicate given histogram as closely as  
(Gaussian Mixtures)  
(EM Algorithms)  
(SMAI)  
(Notes)

(Intersection of gaussians  
is the threshold)  
(Global thresholding)

— X —



## Graph Cut based Binary Segmentation

- Interactive method where foreground & background  
are marked with lines allowing algorithm to  
create a model of foreground / background.

— X —

~~AMRE~~

## Superpixels:

- If adjacent pixels are quite similar, group them together to form a superpixel.  
 $\Rightarrow$  Essentially these superpixels are highly grainy segmentations.

## \* Graph Cuts: (Normalized Cuts paper)

- Creating graphs directly where each pixel is a node makes them hard to handle, so divide the image into superpixels and then perform graph operations.  
 $\Rightarrow$  Adjacent superpixels are connected by edges with their weights based on similarities.  
 $\Rightarrow$  Perform min cut where we can then split the foreground and the background using this cut.

$$\bullet \text{Cut}(A, B) = \sum_{u \in A, v \in B} w(u, v) \quad \left. \begin{array}{l} \text{Performing direct} \\ \text{min cut usually} \\ \text{just puts one} \\ \text{superpixel in a} \\ \text{region & separates} \\ \text{the rest.} \end{array} \right\}$$

$\therefore$  We use normalized cuts.

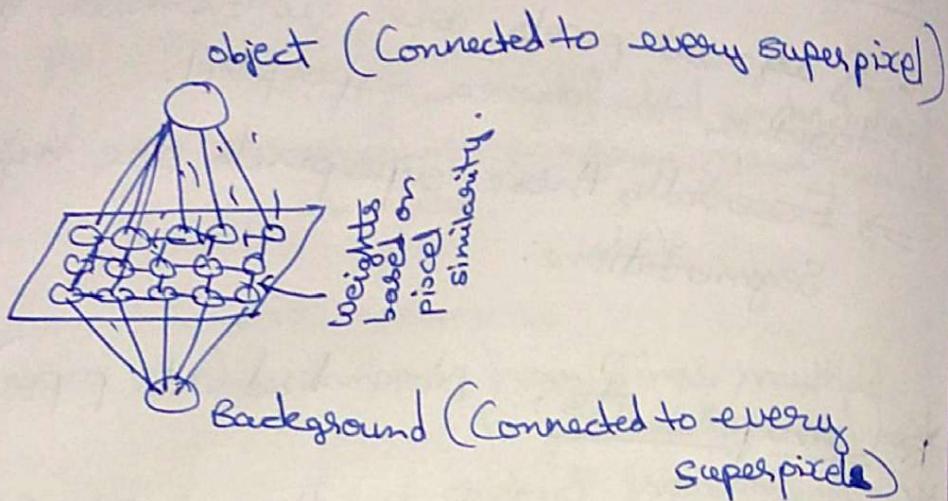
$$* N\text{Cut}(A, B) = \frac{\text{cut}(A, B)}{\text{assoc}(A, V)} + \frac{\text{cut}(A, B)}{\text{assoc}(B, V)}$$

Minimize this      Total Number of edges going out from all elements in A + edges in A.

(Read the paper)

## \* Binary Image Segmentation: (Graph Cut Based)

- Perform graph cut but add model information as well.



⇒ Ask user to specify 2 lines, one on the foreground and one on the background

- Strength of connection from Object to pixel & Background to pixel depends on the input user gave (essentially we are adding in class information)

⇒ Graph cut where one constraint is added that "Object" & "Background" should

$$E(x) = \sum_i c_{ii} x_i \quad \text{belong to different sides of cut}$$

$$+ \sum_{ij} c_{ij} x_i (1 - x_j)$$

Note: Energy function: (example in slides)

- Assign labels to vertices such that overall cost is minimized. → Energy function (Overall cost)

- We want to minimize the energy
- If vertices are connected, define additional cost as.  $\hat{f}(a,b) = f(a) + f(b)$

Some weight  
cost between edge ab.  
Weight of edge ab.

where  $\psi(a,a) = 0$ ,  $\psi(a,b) = 1$  (Same labels  
 $\uparrow$   
 $\emptyset$  = Pairwise potential  
 adjacent pixels  
 add no cost whereas  
 You can define any costs here.  
 different labels add cost)

$$*\therefore \text{Finally } Q(f) = \sum_a \underbrace{\emptyset_{a,f(a)}}_{\substack{\text{Cost of a certain labelling}}} + \underbrace{\sum_{(a,b)} \emptyset_{ab,f(a)f(b)}}_{\substack{\text{Unary potential} \\ (\text{Cost of labelling})}} + \underbrace{\sum_{(a,b)} \emptyset_{ab,f(a)f(b)}}_{\substack{\text{Pairwise potential} \\ (\text{Cost due to connected vertices})}}$$

$\Rightarrow \psi^* = \min Q(f; \emptyset) = Q(f^*; \emptyset)$   
 This is a huge problem:  $2^{\text{vertices}}$  assignments possible.

$\Rightarrow$  Cost of a source | sink cut = Sum of edges on cut going from source to sink.

These energy functions are minimized via mincuts.

\*Note: Max flow algorithms:

- Start with flow 0
- Find path with the flow from source to sink and push max into it as possible and add to flow.
- Subtract ~~the~~ capacities this pushed cap and add to back edges.
- Repeat.

$\Rightarrow$  Mincut = Edges whose capacity is 0 after max flow. Value of mincut = max flow obtained.

— x —

Note: Papers on energy minimization. Fast approx etc.

— x —

Note: Maximum Apriori P?

## ① Grab Cut : (Segmentation)

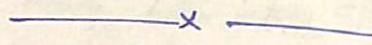
- \* • Specify region of image that does not contain the background. (To create a background/foreground model).

⇒ Uses graph cut for segmentation.

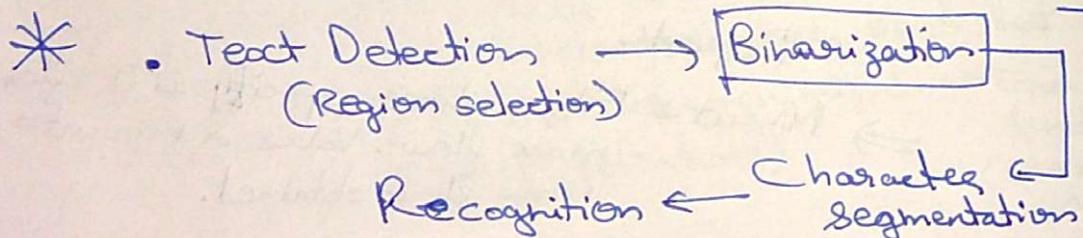
⇒ Foreground / Background models based on mixture of gaussians.

- Iterative graph cut procedure where we start with region outside rectangle as background & inside rectangle as foreground. → Perform graph cut & then based on assignment, retrain foreground / background models & repeat the process till convergence.

⇒ Foreground / Background border is separated via matting (smooth transition)  
(Like a sigmoid)



## ② MRFs for Binarizing Natural Scene Text

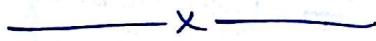


⇒ Pairwise & individual energy terms.

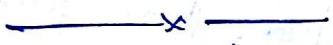
⇒ Also added an ~~edgeless~~ <sup>edges</sup> term to the energy → Generally on segment boundary there are edges.

⇒ Heuristic to create foreground / background models → Perform edge detection & if edge near top, above it is background.  
 Two parallel nearby lines → area b/w them belong to the foreground.

→ Perform standard graphcut.

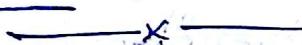


Note: OBJ CUT paper. Read GrabCut & MRF\* (2) papers.



→ Graph Based methods

## Multilabel MRF (Semantic Segmentation)



More than 2 labels makes graph cut NP hard.

↪ We use  $\alpha$ -expansion approximation algorithms

↪ Approximation quality =  $2 \times$  Best.  
 (Potts model is an example) (Slides)

### \* Metric Functions:

- $V(\alpha, \beta) = 0 \Leftrightarrow \alpha = \beta$  } Semimetric functions.

$$V(\alpha, \beta) = V(\beta, \alpha)$$

$$V(\alpha, \beta) \leq V(\alpha, \gamma) + V(\gamma, \beta)$$

$$\Rightarrow E(f) = \sum_{\{p, q\} \in N} V_{p,q}(f_p, f_q) + \sum_{p \in E} D_p(f_p)$$

Energy function generalization

Note: Move based algorithms: Have a local search neighbourhood and move to best solution in that region.

Keep moving as long as you are improving on your solution.

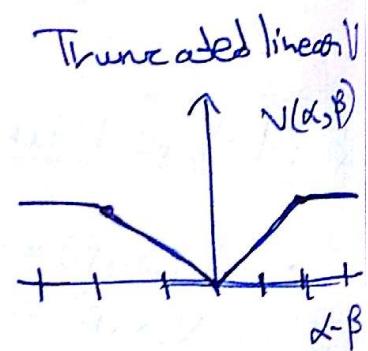
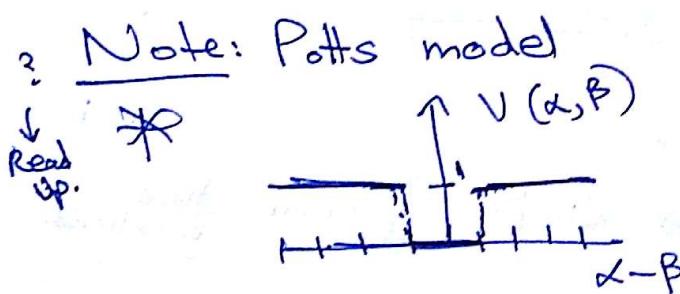
## \* $\alpha$ -Expansion Algorithm: (Semantic Segmentation)

- 1). Start with random label assignments to pixels.
- \* Move over each label  $\alpha$  performing binary classification ( $\alpha$  or Rest) and then check if score increases. (We only move pixels not having  $\alpha$  to having labels of  $\alpha$ , hence the name  $\alpha$ -Expansion).
- If your energy function does not reduce any further, then stop. (Expansion does not change assignment anymore) (Example in slides)

(Local algorithm (might get stuck at local optima at times))

- At each step we perform graph cut since its only 2 classes now ( $\alpha$  &  $\sim\alpha$ )

Graph Cut



- These models are used in energy functions which are used on graphs.

- This linearity helps counter structured noise in images while classifying etc.

(Create graph cut  
& now check if the  
energy function is  
minimized after cut)

Graph Cut

## \*2) Ink Bleed: (Removal)

\* • Ink written on other side of paper bleeds through

\*  $\Rightarrow$  Front (Recto) & Back (Verso) sides of the paper.

- Flip the Back and align the two images.

- We assume we have 3 labels - Foreground, Background and Bleed.

$\Rightarrow$  Dual Layer MRF: First layer has vertices from front image  
Second layer has vertices from second image.

- We define weights between each node in these layers (corresponding) such that the cost function is defined such that "InkBleed-InkBleed", "InkBleed-Background", "Background-InkBleed" etc get a large value since we don't want them (We know it's not possible)

$\Rightarrow$  Unary assignments:

$$\text{Foreground } E(l=F) = \frac{S_I + S_B}{2 \sum (S_F + S_I + S_B)} \quad \begin{matrix} \text{From model!} \\ (\text{Similarities}) \\ (\text{SUM}) \\ (\text{outputs}) \end{matrix}$$

$$\text{Background } E(l=B) = \frac{S_F + S_I}{2 \sum (S_F + S_I + S_B)}$$

$$\text{InkBleed } E(l=I) = \frac{S_F + S_B}{2 \sum (S_F + S_I + S_B)}$$

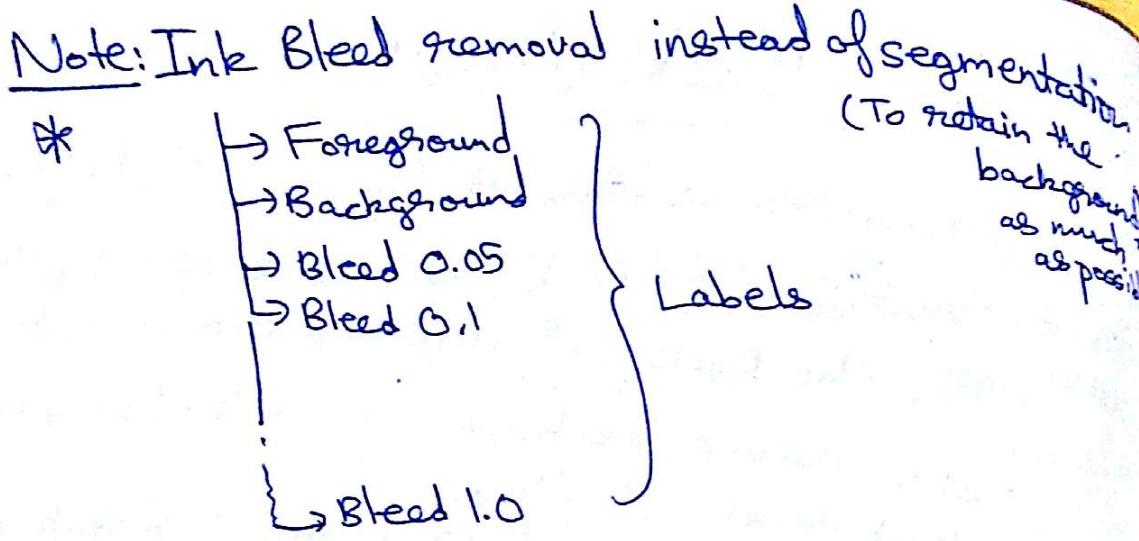
$$\Rightarrow \text{Pairwise: } E_S = \sum_{p \in VEN} V_1(l_p, l_q) + \sum_{p, p' \in EM} V_2(l_p, l_{p'})$$

Matrices define these costs

Intralayer  
(~~different~~ class assignments get a large cost)

Between layers  
(we give preference to a few pairs)

- Then use  $\alpha$ -expansion to solve the problem.

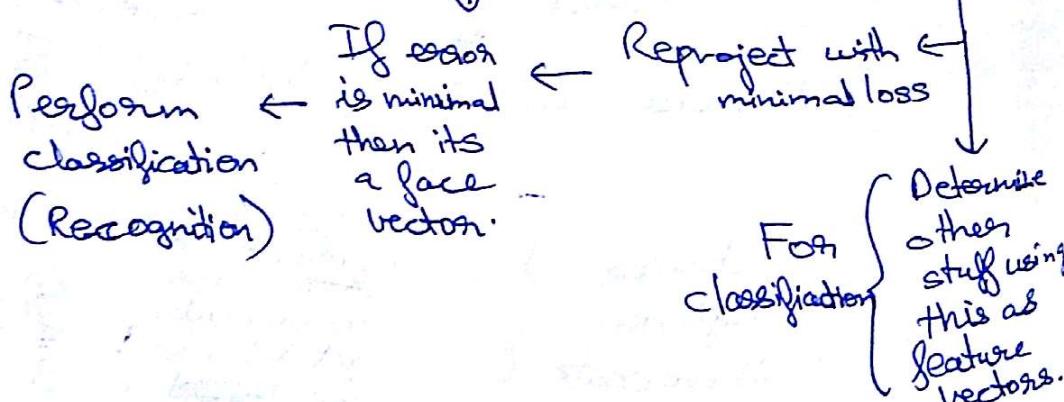


- $\text{Ink}_p = I^{(\alpha_p + \beta_p)} \cdot \text{Background}_p$  } Light absorption physics
  - Estimate  $\alpha_p$  &  $\beta_p$  for all pixels
  - Ink Bleed.
- Then perform standard  $\alpha$ -expansion after defining the energy functions.

Note: Depth/Disparity estimation, Normal/Shape estimation, Image restoration/enhancement.

### \* Face Detection with Eigenfaces:

- Face Image  $\rightarrow$  Vector ( $n \times m$ )  $\rightarrow$  Perform PCA using training data



## \* Face Detection: (Viola and Jones)

- \* Try every window by check if its a face or not. This is too slow since there are too many windows.

→ False positive rate should be very low and detection rate should be high.

### i) Image Features:

- Use binary masks convolving with the image where feature values are the convolution sum. (bright - dark)

→ Some features give the true values at certain parts of the image which we can use to figure out parts of the face.

- To compute these binary convolutions, we need rectangular sums of the image, this is done very quickly using the integral image.

### ii) Window rejection: (Speedup) (Cascading)

Perform at multiple scales  
(Image pyramid)  
for scale invariance  
(Windows used)

- Prioritize binary masks so that if some mask fails ~~at a window~~ at a window then we can directly reject else we try more masks and check if the features are valid or not (to confirm faces)

### iii) Boosting: (Adaboost)

- Weak learners → features that perform just a little better than random.  
Each filter that is thresholded → Adjust the threshold so that no face is missed (spurious are allowed)

- We then combine multiple of these using AdaBoost,  $\rightarrow$  to create stronger classifier
- $\Rightarrow$  Since our thresholds are adjusted such that even a single reject means it's not a face, we have fast window rejection here.

Note: Cascading is using a series of classifiers, accept = move to next classifier, reject = stop. As we move forward classifiers are made more complex using many weak learners at once using AdaBoost.

Note: Training is very slow, to create classifiers using AdaBoost requires multiple iterations. Testing is quite fast though.

(Boosting illustration in slides)  
(Cascade structure in slides)

# Pedestrian Detection

## (HOG + SVM)

↳ Generative / Discriminative models

↳ Global / Part based descriptors

↓  
Handles  
low res  
images  
well.

Each part of the human body  
↳ Handles occlusions etc well

• Gradient Histograms (SIFT, GLOH, HOG etc)

→ Gradients:  $\begin{bmatrix} -1 & 1 \end{bmatrix}$  } One sided      } Then magnitude &  
                   $\begin{bmatrix} -1 & 0 & 1 \end{bmatrix}$  } Two sided      } angles

• Training data → Crop people → ~4000 dimensional HOG features

Binary classifier ←  
(SVM)

i) Detection → Use scale space of images and  
↓ combine adjacent window detections  
64x128 pixel windows      (sliding window)

ii) Descriptor → Compute gradients. (Two sided derivative)

\*      • 8x8 pixels treated as cells and  
compute histogram for each cell giving  
8x16 cells.

Different kinds of  
L2 normalization  
was possible.  
Normalization was good.

• Use ~~a~~ 2x2 cells as a block  
and normalize the histogram. These  
blocks have a 50% overlap as  
well. This gives us 7x15 blocks  
In total:

(For each pixel, gradient = max gradient over all channels)

• Gradient histogram uses 9 bins.

\*      • Trilinear interpolation of gradients

- \*  $\rightarrow$  Gradient magnitude linearly interpolated between nearest 2 bins based on angle.
- $\rightarrow$  Bilinear interpolation of pixel gradient in the closest 4 cells.  
(Overall, the magnitude is split twice,  $\frac{m}{x_1-x_2}$ , both these interpolations together ~~are~~ is called trilinear interpolation)

Note: Gradient image is calculated only once initially to speed up the algorithm.

\*  $\Rightarrow$  Total feature shape =  $7 \times 15 \times 4 \times 9 \underset{\text{(dimensions)}}{\underbrace{\times \times}} \approx 4k$ .  
 $\underset{\text{Blocks}}{\underbrace{\times \times}}$  Cells in a block       $\underset{\text{Bins in a cell.}}{\underbrace{\times \times}}$

Note: Support ~~Vector~~ Vector Machine Detectors, Dynamic Pedestrian detection (Viola & Jones extension), similar feature detections using binary filters HOG however beat them all. Extension of HOG pedestrian  $\rightarrow$  2d Global Detector paper.

### \* SVMs:

$$\cdot R(\alpha) \leq R_{\text{train}}(\alpha) + \sqrt{f(h)/N}$$

$\uparrow$   
Risk during testing assuming

data distribution is the same.

- $h = \text{VC dimensions}$ ,  $f(h) = \text{Monotonically increasing function}$ .
- ↳ We want  $h$  to be small to get better test results.

$$\Rightarrow h \propto \frac{1}{\text{Margin} \rightarrow \text{of the SVM}}$$

- This margin is relative to the diameter of hypothesis over all the points.

$\Rightarrow$  Relative margin:  $\frac{h}{D}$

$$h \leq \min \left\{ d, \left\lceil \frac{D^2}{P^2} \right\rceil \right\} + 1 \quad \begin{cases} \text{VC Dimension} \\ \text{Dimensionality of samples.} \end{cases}$$

$\Rightarrow$  This minimization statement disproved the curse of dimensionality, we just need  $h$  to be small for low errors.

We formulate an optimization problem to maximize the relative margin. (SMAI Notes)

Only support vectors influence the boundary ( $w$ )

- i) Allow slack variables for noisy data
  - ii) Kernels for non linear boundaries
- $\rightarrow K(x,y) = (1+x \cdot y)^n$  ? Huge dimension expansion
- (Note: Mercer's Theorem  $\rightarrow$  Rules for kernel functions)



## \* Visual Bag of Words:

- Histogram of words → Document summary  
(Use dictionary)

↳ Spatial info is lost  
↳ Syntax/Grammar is lost

- Bag of features models for images  
    ↓  
    visual words

### i) Texture Recognition:

- Repetition of textures (units of textures)
- Use bag of words → Histogram of all possible textures (dictionary)  
(Comparing histograms with distance measures)

### ii) Bag of Words Model:

- Histogram of words (Tag clouds etc)

I)  $\Rightarrow$  Feature extraction: Compute visual features from all images to learn a visual vocabulary  
↳ Quantize to get a finite vocab.

#### \* 1) Patch retrieval:

- i) ~~Sampling~~ Regular grid of a specific size.
- ii) Interest point detection & patches around them. (Ellipses around them etc)

#### 2) Features:

- Normalize patch & extract SIFT features

## II) Learning visual vocab:

- \* If SIFT features of ~~are~~ some patches are too close to each other, we should consider them as only one word in our vocab.
- Apply a clustering algorithm which would give you a vocab of best (most dissimilar) features. (Cluster centers)  
(K-Means clustering)  
→ This vocab is our "codebook".

## III) Histogram creation: (for test data)

- Given an input image, extract features in a similar manner and then hard cluster it into one of the elements of our codebook (nearest neighbor)
- Now use the histogram generated for other tasks.

Note: Scale issues in images → Pyramids or divide image into regions & get individual histograms.

Note: Dictionary size? Computation efficiency → Vocab tree? paper?  
Image retrieval using histograms. Generative or discriminative models while training - Sparse visual histograms usually → can be used for speeding up matching etc.

### \* Uses: (Visual histogram (bag) of words)

- i) → Image classification: Discriminative methods, standard classification techniques
- ii) Topic discovery: Learn about the space from which data is drawn

iii) Action recognition: dictionaries specific to task at hand.

iv) Image search:

- Copyright issues → discover these
- Retrieve similar images etc.

Note: Pascal VOC } original object detection & localization challenge

Note: Video Google → Features from videos, similar (paper) videos, video search etc.

## Feature Learning with DL:

- Perception learning
- Multi layer perception  $\rightarrow$  First autonomous vehicles
- CNNs - Animal neuron receptive fields gave idea.
- Dropout
- $\Rightarrow \frac{(N-F)}{S} + 1$  { Conv layer output size  
 $(N = N + Pad)$
- Max Pooling  $\rightarrow$  Useful (Preserve maximum activation)  
~~AlexNet, ZFNet, VGGNet, GoogleNet, ResNet (52 layers)~~  
~~(Inception modules)~~

## Image Indexing & Retrieval:

- Content based image retrieval (CBIR)
  - Attribute based
  - Object based
- CBIR  $\rightarrow$  Database population phase  
Retrieval phase

Note: "Things"  $\rightarrow$  Foreground important objects

"Stuff"  $\rightarrow$  Background objects

### 1) Feature selection for Database population:

- Initially, color histograms, color moments, texture descriptors are used, for retrieval, use a distance metric for similarity.
- Practical issues with huge amounts of data & specific subdomains with specific matchers - faces, monuments etc.

- Text indexing & retrieval is well explored, converting images into text documents can help perform this well.

### i) Virtual Textual Representation:

- Get segments of image that seem interesting
- Transform segments into text.
  - ↳ Then perform quantization to bin these segments. (Each cell has a representative string, so each segment now has a text string)

### 2) Inverted Index:

- For every word in an extended dictionary (all names, etc) have an index which points each word to all the documents and where it occurs. (Crawling to build the index)
- Based on this inverted index, given a query, we can try find the most relevant documents.

⇒ Documents can be represented as word histograms.

- Inverted indexes still however give you a lot of results and then page ranking needs to be done to get the best result. Page ranking generally depends on number of other sites linking to the page.

Weighted word histograms.

For other types of files such as images, their representation changes as noted in the next section.

### \* 3) Document Representations:

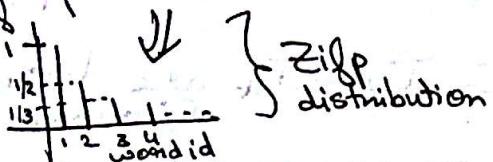
i) Binary weights:  $A_{ij} = 1$  if document  $j$  has word  $i$  or else its 0.

ii) Raw term frequency

\* iii) TFIDF:

It was noted that ratio of these frequencies in decreasing order had a specific pattern:

- Get frequency of each word over all documents you have
  - Get frequency of ~~each~~ words in relevant documents.  $\rightarrow$  Term frequency
- $\Rightarrow$  Now if TF is high & IDF is low then we consider the word as important.



• word frequency

$$\propto 1/i^k$$

$$\propto 1/i^k$$
</div

## \*Optical Flows & Tracking

- Feature tracking - sparse correlation
- Optical flow - dense ~~correlation~~ correlation (every pixel's motion is calculated)

⇒ At time  $t$ , pixel at  $(x, y)$  moves to  $(x+u, y+v)$  at time  $t+1$ .  $(u, v) = \text{Shift of pixel at } (x, y)$

$$\Rightarrow I(x, y, t) = I(x+u, y+v, t+1)$$

- Using Taylor approx,

$$I(x+u, y+v, t+1) = I(x, y, t) + I_x \cdot u + I_y \cdot v + I_t$$

( $I_x, I_y$  &  $I_t$  are derivatives of image).

$$\therefore \text{Difference} = I_x \cdot u + I_y \cdot v + I_t = 0$$

$$\Rightarrow \boxed{\nabla I \cdot [u \ v]^T + I_t = 0} \rightarrow ①$$

- This cannot be used since each pixel has only 1 equation but 2 variables  $u, v$  to solve.

⇒ We add in another constraint, if we take a vector parallel to edge (perpendicular to gradient),  $\nabla I \cdot [u' \ v']^T = 1$  then ① should be satisfied by this as well  $\Rightarrow (u+u', v+v')$  satisfies it.

Slides ↪ (The Aperture Problem demonstrates this)  
(Motion perpendicular to gradient is not visible)

Motion along direction of edge is not noticeable, so if actual direction  $= (u, v)$  for a pixel, then if  $(u', v') = \text{direction along edge}$  then

$(u+u', v+v')$  direction is also a valid solution so it would satisfy equations ①.

Note: This can also be observed in audio.

i) Spatial coherence constraint: (Lucas-Kanade)

- To solve 1 equation 2 variable issue, assume adjacent pixels move in the same direction.
- use  $5 \times 5$  window to get 25 equations for 2 variables. (Overconstrained)

$$\begin{bmatrix} I_{xx}(p_1) & I_{xy}(p_1) \\ I_{yx}(p_1) & I_{yy}(p_1) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_x(p_1) \\ I_y(p_1) \\ \vdots \\ I_x(p_{25}) \\ I_y(p_{25}) \end{bmatrix}$$

Get  $(u, v)$  using pseudo inverse.

$$Ad = b \Rightarrow (A^T A)d = A^T b$$

$$d = (A^T A)^{-1} A^T b$$

$$A^T A = \begin{bmatrix} \sum I_{xx} I_{xx} & \sum I_{xx} I_{yy} \\ \sum I_{xy} I_{yx} & \sum I_{yy} I_{yy} \end{bmatrix}$$

$$A^T b = \begin{bmatrix} \sum I_{xc} I_{xt} \\ \sum I_{yc} I_{yt} \end{bmatrix}$$

Essentially  $\det$  should not be 0,

- If plain image or just vertical/horizontal gradients would make  $\det 0$

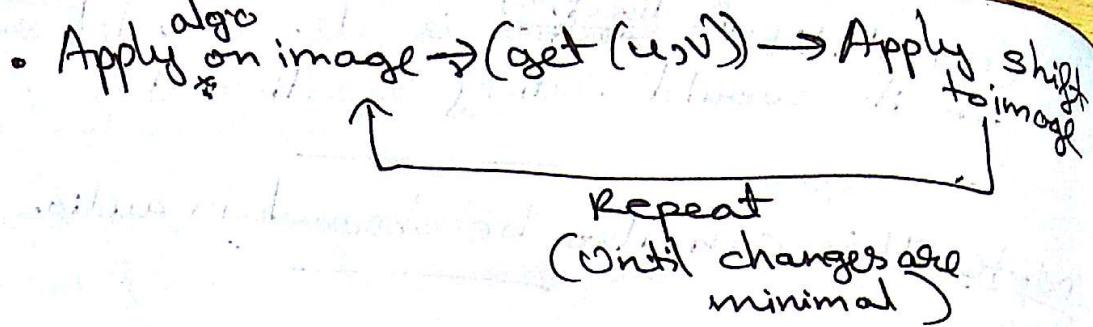
(Eigen vectors of  $A^T A$  should both be large)

(Similar to Harris conditions)

This would ensure our predictions of  $(u, v)$  are correct and not ambiguous.

ii) Lucas Kanade - Iterative (Algo in slides)

- Large movements are hard to approximate directly with linear approximation, so we apply the same procedure multiple times to get the total motion (Iterative updates)



- iii) Use scale space pyramids to approximate larger changes well.
- ↳ On spatial coherence constraint.  
(LK Algo)
  - (Can be iterative as well)  
at each level

### \* KLT - Tracking:

- Similar to flow algorithms
  - ↳ But also uses second moment matrix for better approximation.
- Find good points to track — Harris corners.
- Iterate & coarse-to-fine search for large motion.
- Small windows — sensitive to noise, miss long motion without pyramid.
  - ↳  $15 \times 15, 31 \times 31$  are common.  
(weighted windows used)

### Motion & Perceptual Organization

- ↳ Grouping of objects based on motion.  
(Even motion with multiple points helps us group better with a lot more confidence)

Note: Applications of optical flow,  
Segmentation based on motion cues, learning and tracking, dynamic models, activity recognition and improving video quality, MPEG-4 video compression.

\* Motion Field :

- Ideal optical flow (Ground truth sort of)
  - Compute actual motion on 3D objects & project onto the image.  $\rightarrow$  Motion field.
  - Optical flow we compute from just images is apparent motion. (Approximate)

$\Rightarrow$  We want optical flow to be equal to motion field.

Optical flow is based on illumination, Motion field is based on actual motion of objects

Note: Further approaches have added a lot more on top of this. Dense Flow  $\rightarrow$  Read.

~~first~~ early in the month of June (1st) and a  
(Bratt) week later all gathered ripe.  
All are oil (60%) if one month old  
green and soft at 20% oil (one year)  
(older) ready to eat if go with one  
(early) soft to tender after freezing.

## Projective Geometry: (Recap)

$\bar{l} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$ ,  $\bar{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$   $\rightarrow$  Homogeneous coords.  $\rightarrow$  Ray in 3D through origin.

$\bar{l}^T \bar{x} = 0$  { Equation of a line  
 (All points on line if  $\bar{l} = \text{const}$ )  
 $\rightarrow$  Set of lines through  $\bar{x}$   
 if  $\bar{x} = \text{const}$ .

Points at infinity  $\begin{bmatrix} x \\ y \\ 0 \end{bmatrix}$  can be seen as  $(\infty, \infty)$  but  
 in the direction  $(x, y)$ .  $\rightarrow$  Ideal point in direction  $(x, y)$

$\Rightarrow [x \ y \ 0]^T$  are all points at infinity, they  
 \* together form the line at infinity  $[0 \ 0 \ k]^T$

- Lines  $\begin{bmatrix} a \\ b \\ c \end{bmatrix}$  on projection plane can be thought of as a plane through origin intersecting the  $z=k$  plane. (In 3D)
- \*  $(a_0x + b_0y + c_0z = 0)$  if  $(x_0, y_0, z_0)$  lie on line and  $(a_0, b_0, c_0)$  is  $\perp$  to the line, so if we think of it as a plane  $(a_0, b_0, c_0)$  represent the normal of the plane)

$\Rightarrow$  The line at infinity  $(0, 0, k)$  is the plane that is  $\perp$  to our projection plane (since they intersect at  $\infty$ ).

\* Line joining 2 points:

•  $\bar{l}^T p = 0, \bar{l}^T q = 0$

$$\bar{l} = [(y_2 - y_1) \quad -(x_2 - x_1) \quad (x_2 y_1 - x_1 y_2)]^T$$

(In normal 2D lines)

⇒ Thinking of it in terms of plane through origin in 3D → each  $p, q$  are rays from origin and line is the plane containing these 2 rays so normal of plane is  $\perp$  to both rays  $p \& q$ .

\*  $\therefore l = \overline{p} \times \overline{q}$  {Cross product.  
(we will get the same result we had before)}

$$\text{ex: } \overline{p} = (5, 2, 1) \\ \overline{q} = (3, 2, 1)$$

$$l = \begin{bmatrix} i & j & k \\ 5 & 2 & 1 \\ 3 & 2 & 1 \end{bmatrix}$$

$$= i(0) - j(2) + k(4) \\ = -2\hat{j} + 4\hat{k}$$

$$l = \begin{bmatrix} 0 \\ -2 \\ 4 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \\ 2 \end{bmatrix}$$

(Scaling is fine)

ex: Ideal point of line  $[0 \ 1 \ -2]^T$  is  $[1 \ 0 \ 0]^T$ .  
This is same as  $[0 \ 1 \ k]^T$  for any  $k$ .

ex: Line joining  $[3 \ 4 \ 0]^T$

Point of intersection of 2 lines:

$$\bullet l^T x = m^T x = 0, x = l \times m$$

(Duality → similar to line through 2 points)

\* Note: Ideal point of  $l = [a \ b \ c]^T$  is  $[b \ -a \ 0]^T$

## \* Conics:

$$ax^2 + bxy + cy^2 + dx + ey + f = 0$$

$$\bullet ax^2 + bxy + cy^2 + dxw + eyw + fw^2 = 0$$

If  $w=1$ , we get the original equation of the conic.

(Using homogeneous coordinates)

$$\begin{bmatrix} x & y & w \end{bmatrix} \begin{bmatrix} a & b/2 & d/2 \\ b/2 & c & e/2 \\ d/2 & e/2 & f \end{bmatrix} \begin{bmatrix} x \\ y \\ w \end{bmatrix} = 0$$

$$x^T C x = 0$$

↑ Symmetric matrix

(Degenerate cases where  $C$  reduces conic to a line)

## \* Tangents:

$$\Rightarrow x^T (Cx) = 0$$

$$x^T l = 0$$

$\boxed{l = Cx}$  Gives us a line through  $x$ . (So it touches the conic)

Tangent to conic at  $x$ .

Proof: (Suppose  $l$  intersects  $C$  at another point  $y$ ,

$$y^T C y = 0, y^T l = 0 \Rightarrow y^T C x = 0$$

$\therefore l = \text{Tangent}$

So the line  $l$  only intersects the conic at a single point.

Only possible if  $x$  &  $y$  are the same point.

$\Rightarrow$  If  $x, y$  lie on a line axis also lies on the line  $\Rightarrow$  So it would satisfy  $C$ , so they would all lie on  $C$

## Dual Conic:

- ~~l~~ Set of all lines tangential to the conic  $C$ .
- $l^T C^* l = 0 \quad \& \quad l^T C^{-1} l = 0 \quad \} \text{Set of lines tangential to } C$
- $\Rightarrow C^* = \text{Dual of } C = C^{-1}$

- Translation  $\xrightarrow{(2) \rightarrow 00F}$  Euclidean  $\xrightarrow{(3)}$  Similarity  $\xrightarrow{(4)}$  Affine  $\xrightarrow{(6)}$  Projective  $\xrightarrow{(8)}$
- $\uparrow$

$\Rightarrow$  Transformations.

\* Projective Transform: (Collinearity is preserved)

- $x' = H \in \mathbb{P}^3$  Transformed points
- $l' = H^{-T} l$  Transformed line
- $C' = H^{-T} C H^{-1}, C^* = H C^* H^T$  Transformed conics & duals.

Note! Isometric transforms are,

$$\star \begin{bmatrix} \delta \cos\theta & -\delta \sin\theta & a \\ \delta \sin\theta & \delta \cos\theta & b \\ 0 & 0 & 1 \end{bmatrix} \quad \delta = \pm 1$$

(If  $\delta = \text{Any non zero scalar} \rightarrow$  we have scaling as well).

(similarity transform)

Preserve angles, parallelism, ratio of lengths & circular points.

Note: Metric structure is the geometric reconstruction of an object that is defined upto an unknown similarity transformation. (True scale, rotation etc are not known just from image)

(Reconstructions are metric structures)

Note: Affine transforms map points at  $\infty$  to other points at  $\infty$ .

$$\begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \text{Affine transforms.}$$

Projective Transforms decomposition: (No need to know how)

- $\begin{bmatrix} A & \mathbf{t} \\ \mathbf{v}^T & 1 \end{bmatrix} = H_p \cdot H_A \cdot H_s = \begin{bmatrix} I & 0 \\ \mathbf{v}^T & 1 \end{bmatrix} \begin{bmatrix} K & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} S_R & 0 \\ 0 & 1 \end{bmatrix}$

$K = \text{Upper diag with } \det = 1$

- Preserves cross ratios & collinearity.
- 8 degrees of freedom (4 point matches needed)
- Finite points can map to ideal points & vice-versa (train tracks)

\* Cross Ratios on a line:

$$\text{Cross}(x_1, x_2, x_3, x_4) = \frac{|x_1 x_2| |x_3 x_4|}{|x_1 x_3| |x_2 x_4|} \quad \text{where } |x_i| = \det \begin{pmatrix} x_1 & y \\ x_2 & y \end{pmatrix}$$

(For 4 concurrent, coplanar lines, cross ratio of 4 points measured as intersections of any line with those lines is constant)

\* Circular points: (Preserved till similarity transforms)

- Intersection of circle with line at infinity.
- $$\Rightarrow ax^2 + ay^2 + dxw + eyw + fw^2 = 0$$
- $\Rightarrow$  Intersection with line at  $\infty \Rightarrow w=0$

$$\Rightarrow ax^2 + ay^2 = 0 \quad \Rightarrow x^2 + y^2 = 0$$

$\Rightarrow$  Simple solutions are  $(1, i) \Rightarrow \begin{bmatrix} 1 \\ i \\ 0 \end{bmatrix}$  and  $(1, -i) \Rightarrow \begin{bmatrix} 1 \\ -i \\ 0 \end{bmatrix}$

Note: Table comparing properties of all transforms in slides.

### Finding Line at Infinity: (From images)

Using tracks (Gibson) & table in Ravi's notes

length of shadow =  $\frac{1}{2} \times \text{object height} \times \tan(\theta)$   
where  $\theta = \text{angle between line of sight and tangent to shadow cast by object}$

length of the track shadow based on this

length from sun to shadow will be aligned with angle of the sun

length of track shadow of shadow traced out

length shadow on surface

(calculated)  $\theta = 45^\circ$  (calculated from diagram)  
length of track (approx 140)  $\theta = 45^\circ$

length of track  $\theta = 45^\circ$  (calculated from diagram)  
length of track  $\theta = 45^\circ$  (calculated from diagram)  
length of track  $\theta = 45^\circ$  (calculated from diagram)

length of base wall and front wall of track  
length of shadow of base wall of track

length of shadow of front wall of track

length of shadow of base wall of track  
length of shadow of front wall of track

length of shadow of base wall of track  
length of shadow of front wall of track

length of shadow of base wall of track  
length of shadow of front wall of track

length of shadow of base wall of track  
length of shadow of front wall of track

length of shadow of base wall of track  
length of shadow of front wall of track

## Graph Based Methods: (in CV)

- Irregularly sampled data domains are generally modelled as graphs.
- BSP (Bulk Simulated Processing)  $\rightarrow$  Local graph updates!

$\Rightarrow$  Spectral Graph Theory : Function defined on vertices. Smooth this function etc using the structure of the graph.

- MRF based methods also build a graph.
- Graphs are non euclidean so we can't compute derivatives etc.

$\Rightarrow$  We want our edges to be as local as possible over the euclidean manifold.

- \* • Graph Laplacian:  $L = D - W$  (Embedding)
- \*  $\hookrightarrow -\exp(-\|v_i - v_j\|^2/2) = w_{ij}$  } Set weights of edges as such.

$$L = U \Lambda U^T \quad \left. \begin{array}{l} \text{Ignore zero eigenvalues} \\ \text{and its corresponding eigen vectors} \end{array} \right\}$$

- The eigen vectors are then used to project graph nodes onto the euclidean space spanned by the d eigen vectors.  
(Embedding space) (Graph Laplacian Embedding)
- $\rightarrow$  The distance in this embedding space depicts the average connectivity in the original graph.

### \* Heat (Diffusion) Kernel Embedding:

$$H(t) = e^{-tL} = U e^{-t\frac{\lambda}{2}} U^T \quad \left. \begin{array}{l} (L = \text{Graph Laplacian}) \\ (\text{Diagonal matrix}) \end{array} \right\}$$

$\uparrow$   
Kernel Matrix.  
 $(n \times n \rightarrow n \text{ nodes in graph})$

$$= U e^{-t\frac{\lambda}{2}} (e^{-\frac{t\lambda}{2}})^T U^T$$

$$= X \cdot X^T$$

- The  $i^{th}$  column of the matrix at a time & show how the heat source spreads if rest are treated as the sinks.

(Similar to how Adjacency matrix shows one edge connectivity,  $A^2 \Rightarrow$  Two edge connectivity b/w nodes and so on...)

### \* Dense 3D Matching: (Paper) (Using Heat Diffusion)

- \*  $\Rightarrow$  Topology (Structure) changes in the graph give us different graph embeddings so compositions are hard.
- $\hookrightarrow$  Heat diffusion at small scales however works, so we can use this.

- Given certain reference points on a graph, we can use them to compute a heat diffusion based descriptors for every other point on the graph.

(Similar to barycentric point descriptors in 2D)

- $\Rightarrow$  So using these point matches, we find a dense matching between all nodes using the descriptors. We then apply EM for smoothing.

### \* Spectral Correspondence for Image Matching: (Paper) (CVPR)

- Take 2 images  $\rightarrow$  compute SIFT per pixel.

- Take 2 images  $\rightarrow$  compute SIFT per pixel.

$\hookrightarrow$  Graph b/w pixels that are similar

$\hookrightarrow$  Also connect pixels from image 1 to image 2 based on similarity.

$$\Rightarrow L = D - A = U \Lambda U^T$$

- Each eigen vector is of shape  $(n_1 \times n_2) \times 1$

$\hookrightarrow$  Split into  $(n_1 \times 1)$   $n_2 \times 1$

Use MGER  $\hookrightarrow$  Compute similarity between these matching can also be done using SIFT here and those eigen vectors.

## \* Image Annotation: (Paper)

- \* • Uses Hypergraph Laplacian

Edges can connect multiple nodes.

(we can represent a hyperedge using multiple dyadic edges)

$$\mathbf{T} = [h_{11}, \dots, h_{1p}]$$

Laplacian is defined here

- Some Laplacian is defined here.
- Each image in dataset is a node here.
  - ↳ Feature space (VGG features) used to create a k-neighbour graph.

- Hypergraph heat diffusion to diffuse labels from training image nodes to test image nodes.

(Initially test nodes are connected to all with zeros as label values)

→ Useful since more frequent labels we use small diffusion size and large diffusion size for less frequent labels

↳ This was an issue not addressed by earlier methods.