

Computer Vision

Homogeneous Geometry

2D Transformations

• Projective (project)

$$\underline{H} = \begin{pmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{pmatrix} \quad 8 \text{DOF} \quad \square \rightarrow \triangle$$

invariants: concurrency, colinearity, tangency
cross ratio $\rightarrow \frac{\overline{AC} \cdot \overline{BD}}{\overline{BC} \cdot \overline{AD}} = \text{const.}$ (if A,B,C,D points on line)

• Affine (move+rotate+squish+skew)

$$\underline{H} = \begin{pmatrix} a_{11} & a_{12} & t_x \\ a_{21} & a_{22} & t_y \\ 0 & 0 & 1 \end{pmatrix} \quad 6 \text{DOF} \quad \square \rightarrow \parallelogram$$

invariants: parallelism, ratios of areas $\frac{A_1}{A_2} = \frac{A'_1}{A'_2}$,
ratios of lengths on parallel lines $\frac{\overline{AB}}{\overline{CD}} = \text{const.}$,
linear combinations $\lambda_1 \underline{x}_1 + \lambda_2 \underline{x}_2 = \underline{x}_3 \Leftrightarrow \underline{x}'_3 = \lambda_1 \underline{x}'_1 + \lambda_2 \underline{x}'_2$,
line at inf. L_∞

• Similarity (move+rotate+scale)

$$\underline{H} = \begin{pmatrix} sR_{11} & sR_{12} & t_x \\ sR_{21} & sR_{22} & t_y \\ 0 & 0 & 1 \end{pmatrix} \quad 4 \text{DOF} \quad \square \rightarrow \square$$

invariants: ratio of lengths, angles,
circular points I; J, $\underline{\underline{\Omega}}_\infty^*$

• Euclidean (move+rotate) (keeps lengths)

$$\underline{H} = \begin{pmatrix} r_{11} & r_{12} & t_x \\ r_{21} & r_{22} & t_y \\ 0 & 0 & 1 \end{pmatrix} \quad 3 \text{DOF} \quad \square \rightarrow \square$$

invariants: lengths, areas

3D Transformations

• Projective (project) (keeps straight lines)

$$\underline{H} = \begin{pmatrix} \underline{\underline{A}} & \underline{t} \\ \underline{\underline{V}}^\top & v \end{pmatrix} \quad 15 \text{DOF} \quad \square \rightarrow \triangle$$

invariants: intersections, tangency

• Affine (move+rotate+squish+skew) (keeps parallel)

$$\underline{H} = \begin{pmatrix} \underline{\underline{A}} & \underline{t} \\ \underline{\underline{O}} & 1 \end{pmatrix} \quad 12 \text{DOF} \quad \square \rightarrow \parallelogram$$

invariants: parallelism of planes, volume ratios,
centroids, plane at inf. $\underline{\underline{\Pi}}_\infty$

• Similarity (move+rotate+scale) (keeps angles)

$$\underline{H} = \begin{pmatrix} sR & \underline{t} \\ \underline{\underline{O}} & 1 \end{pmatrix} \quad 7 \text{DOF} \quad \square \rightarrow \square$$

invariants: angles, ratios of lengths
absolute conics $\underline{\underline{\Omega}}_\infty$; $\underline{\underline{\Omega}}_\infty^*$

• Euclidean (move+rotate) (keeps lengths)

$$\underline{H} = \begin{pmatrix} R & \underline{t} \\ \underline{\underline{O}} & 1 \end{pmatrix} \quad 6 \text{DOF} \quad \square \rightarrow \square$$

invariants: volume, ...

2D points (\underline{x})

- $\tilde{\underline{x}} = (\tilde{x}, \tilde{y}) \rightarrow \underline{x} = (x, y, w) \in \mathbb{P}^2 = \mathbb{R}^3 / (0, 0, 0)$
- $\underline{x} = (x, y, w) = w(\tilde{x}, \tilde{y}, 1) = w \tilde{\underline{x}} \quad (\tilde{x} = \frac{x}{w}, \tilde{y} = \frac{y}{w})$
- $\tilde{\underline{x}}$: inhomogeneous vector \underline{x} : homogenous vector
 $\tilde{\underline{x}}$: augmented vector

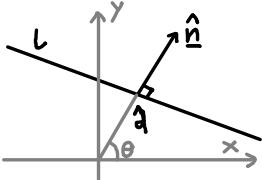
2 coord. $\underline{x}_1, \underline{x}_2$ describe the same point if:
 $x_1/w_1 = x_2/w_2$ and $y_1/w_1 = y_2/w_2$

point at inf.: $\underline{x} = (x, y, 0) \rightarrow$ can't be inhomog.

- transform a point: $\underline{x}' = \underline{H} \underline{x}$

2D lines (\underline{l})

$$ax + by + cw = 0 \rightarrow \underline{l} = (a, b, c)$$



normalized line: $\underline{l} = (\hat{n}_x, \hat{n}_y, \hat{d}) \quad (\|\hat{n}_x, \hat{n}_y\| = 1)$

$\hookrightarrow (\hat{n}_x, \hat{n}_y)$: normal vect., \hat{d} : dist from orig.

line at inf.: $\underline{l}_{\infty} = (0, 0, 1)$

\hookrightarrow includes all points at inf. $\underline{x} = (x, y, 0)$

polar coord. line: $\underline{l} = (\cos \theta, \sin \theta, d)$

- point is on line: $\underline{x}^T \cdot \underline{l} = 0$

- intersection of 2 lines: $\underline{x} = \underline{l}_1 \times \underline{l}_2$

- line from 2 points: $\underline{l} = \underline{x}_1 \times \underline{x}_2$

- transform a line: $\underline{l}' = (\underline{H}^{-1})^T \cdot \underline{l}$

- angle between lines: $\cos \theta = \frac{a_1 a_2 + b_1 b_2}{\sqrt{(a_1^2 + b_1^2)(a_2^2 + b_2^2)}}$

3D lines

\downarrow point \underline{X} on line \downarrow

- line as 2 points $\underline{X}_1, \underline{X}_2$: $\underline{X} = \lambda_1 \underline{X}_1 + \lambda_2 \underline{X}_2$
- line as point + dir. $\underline{X}_1, \underline{d}$: $\underline{X} = \underline{X}_1 + \lambda (\underline{d}, 0)$
- line as 2 planes $\underline{\Pi}_1, \underline{\Pi}_2$: $\begin{pmatrix} \underline{\Pi}_1^T \\ \underline{\Pi}_2^T \end{pmatrix} \cdot \underline{X} = \underline{0}$
- Plücker line: $\underline{\underline{L}} = \underline{X}_1 \underline{X}_2^T - \underline{X}_2 \underline{X}_1^T \rightarrow ?$

3D points (\underline{X})

- $\tilde{\underline{X}} = (\tilde{x}, \tilde{y}, \tilde{z}) \rightarrow \underline{X} = (x, y, z, w) \in \mathbb{P}^3 = \mathbb{R}^4 / (0, 0, 0, 0)$
- $\underline{X} = (x, y, z, w) = w(\tilde{x}, \tilde{y}, \tilde{z}, 1) = w \tilde{\underline{X}} \quad (\tilde{x} = \frac{x}{w}, \tilde{y} = \frac{y}{w}, \tilde{z} = \frac{z}{w})$
- $\tilde{\underline{X}}$: inhomogeneous vector \underline{X} : homogenous vector
 $\tilde{\underline{X}}$: augmented vector

2 coord. $\underline{X}_1, \underline{X}_2$ describe the same point if:

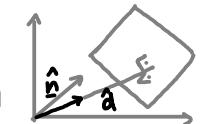
$\tilde{x}_1 = \tilde{x}_2, \tilde{y}_1 = \tilde{y}_2, \tilde{z}_1 = \tilde{z}_2$ (w_1, w_2 : scale)

point at inf.: $\underline{X} = (x, y, z, 0) \rightarrow$ can't be inhomog.

- transform a point: $\underline{X}' = \underline{H} \underline{X}$

3D planes ($\underline{\Pi}$)

$$ax + by + cz + dw = 0 \rightarrow \underline{\Pi} = (a, b, c, d)$$



normalized plane: $\underline{\Pi} = (\hat{n}_x, \hat{n}_y, \hat{n}_z, \hat{d}) \quad (\|\hat{n}_x, \hat{n}_y, \hat{n}_z\| = 1)$

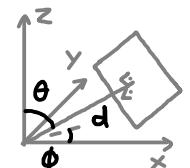
$\hookrightarrow \hat{n}$: normal vect., \hat{d} : dist from orig.

plane at inf.: $\underline{\Pi}_{\infty} = (0, 0, 0, 1)$

\hookrightarrow includes all points at inf. $\underline{X} = (x, y, z, 0)$

spherical coord. plane:

$$\underline{\Pi} = (\cos \theta \cos \phi, \sin \theta \cos \phi, \sin \theta, \hat{d})$$



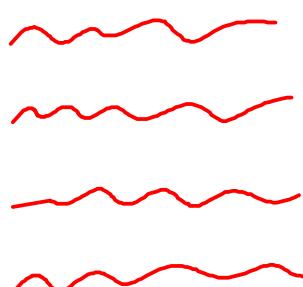
- is point on plane: $\underline{x}^T \cdot \underline{\Pi} = 0$

- plane from 3 points: solve $\begin{pmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \underline{x}_3^T \end{pmatrix} \cdot \underline{\Pi} = \underline{0}$ for $\underline{\Pi}$

- point is on 3 planes: solve $\begin{pmatrix} \underline{\Pi}_1^T \\ \underline{\Pi}_2^T \\ \underline{\Pi}_3^T \end{pmatrix} \cdot \underline{X} = \underline{0}$ for \underline{X}

- transform a plane: $\underline{\Pi}' = (\underline{H}^{-1})^T \cdot \underline{\Pi}$

\downarrow plane $\underline{\Pi}$ on line \downarrow



2D conics ($\underline{\underline{C}}$)

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

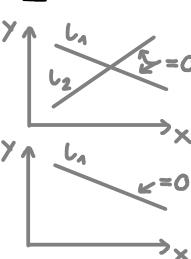
$$\underline{\underline{C}} = \begin{pmatrix} a & b/2 & d/2 \\ b/2 & c & e/2 \\ d/2 & e/2 & f \end{pmatrix}$$

conic: $\underline{\underline{C}}$ / dual-conic: $\underline{\underline{C}}^* = \underline{\underline{C}}^{-1}$

degenerate-conic:

$$\rightarrow \text{rank}(\underline{\underline{C}}) = 2 : \underline{\underline{C}} = \underline{L}_1 \underline{L}_2^T + \underline{L}_2 \underline{L}_1^T$$

$$\rightarrow \text{rank}(\underline{\underline{C}}) = 1 : \underline{\underline{C}} = \underline{L}_1 \underline{L}_1^T$$



• point is on conic: $\underline{x}^T \underline{\underline{C}} \underline{x} = 0$

• line is tang. on conic: $\underline{L}^T \underline{\underline{C}}^* \underline{L} = 0$

• tang @ point on conic: $\underline{L} = \underline{\underline{C}} \underline{x}$

• intersection conic-line: $(\underline{x}_1 + \lambda \underline{x}_2)^T \underline{\underline{C}} (\underline{x}_1 + \lambda \underline{x}_2) = 0$

$$\underline{x}_1^T \underline{\underline{C}} \underline{x}_1 + 2\lambda \underline{x}_1^T \underline{\underline{C}} \underline{x}_2 + \lambda^2 \underline{x}_2^T \underline{\underline{C}} \underline{x}_2 = 0$$

• transform a conic: $\underline{\underline{C}}' = (\underline{H}^{-1})^T \underline{\underline{C}} \underline{H}^{-1}$

$$\underline{\underline{C}}'^* = \underline{H} \underline{\underline{C}}^* \underline{H}^T$$

circular points

intersection of: { all circular conics ($b=0$)
 \downarrow all points at inf. ($\rightarrow \underline{L}_\infty$) ($w=0$) }

$$\underline{x}^2 + \underline{y}^2 = 0 \text{ and } w=0 \rightarrow \underline{I} = (1, i, 0), \underline{J} = (1, -i, 0)$$

resulting conic: $\underline{\underline{C}}_\infty^* = \underline{I} \underline{J}^T + \underline{J} \underline{I}^T = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$

\rightarrow contains all lines going through \underline{I} or \underline{J}

$\rightarrow \underline{L}_\infty$ is its nullvector = line common to \underline{I} and \underline{J}

• invariant for similarity transform and up:

$$\underline{I} = \underline{H}_s \underline{I} \quad | \quad \underline{J} = \underline{H}_s \underline{J} \quad | \quad \underline{\underline{C}}_\infty^* = \underline{H}_s \underline{\underline{C}}_\infty^* \underline{H}_s^T$$

• angle between 2 lines through transform H :

$$\cos \theta = \frac{\underline{L}_1^T \underline{\underline{C}}_\infty^* \underline{L}_2}{\sqrt{(\underline{L}_1^T \underline{\underline{C}}_\infty^* \underline{L}_1)(\underline{L}_2^T \underline{\underline{C}}_\infty^* \underline{L}_2)}} \quad (\text{even if } \underline{\underline{C}}_\infty^*, \underline{L}_1, \underline{L}_2 \text{ transf. by some } \underline{H}!)$$

3D quadrics ($\underline{\underline{Q}}$)

$\underline{\underline{Q}} \in \mathbb{R}^{4 \times 4}$ with 9 DOF

quadric: $\underline{\underline{Q}}$ / dual-quadric: $\underline{\underline{Q}}^* = \underline{\underline{Q}}^{-1}$

degenerate-quadric: $\text{rank}(\underline{\underline{Q}}) \neq 4$

• point is on quadric: $\underline{x}^T \underline{\underline{Q}} \underline{x} = 0$

• plane is tang. on quadric: $\underline{\pi}^T \underline{\underline{Q}}^* \underline{\pi} = 0$

• tang @ point on quadric: $\underline{\pi} = \underline{\underline{Q}} \underline{x}$

• transform a quadric: $\underline{\underline{Q}}' = (\underline{H}^{-1})^T \underline{\underline{Q}} \underline{H}^{-1}$

$$\underline{\underline{Q}}'^* = \underline{H} \underline{\underline{Q}}^* \underline{H}^T$$

absolute conic

intersection of: { all spherical quadrics
 \downarrow all points at inf. ($\rightarrow \underline{\pi}_\infty$) ($w=0$) }

conic on $\underline{\pi}_\infty$ with $x^2 + y^2 + z^2 = 0$

$\rightarrow \underline{\Omega}_\infty: (x, y, z) \underline{\pi}(x, y, z)^T = 0 \text{ and } w=0$
 $\rightarrow \underline{\Omega}_\infty^* = \begin{pmatrix} \underline{\pi} & 0 \\ 0 & 0 \end{pmatrix}$

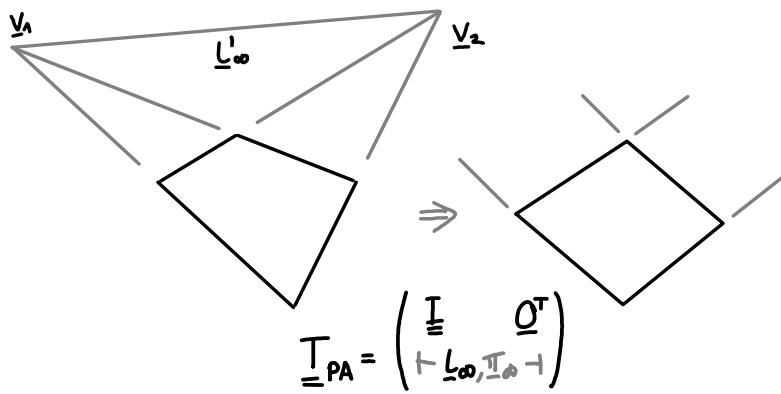
• invariant for similarity transform and up:

$$\underline{\Omega}_\infty^* = \underline{H}_s \underline{\Omega}_\infty^* \underline{H}_s^T$$

• angle between 2 lines through transform H :

$$\cos \theta = \frac{\underline{\pi}_1^T \underline{\Omega}_\infty^* \underline{\pi}_2}{\sqrt{(\underline{\pi}_1^T \underline{\Omega}_\infty^* \underline{\pi}_1)(\underline{\pi}_2^T \underline{\Omega}_\infty^* \underline{\pi}_2)}} \quad (\text{even if } \underline{\Omega}_\infty^*, \underline{\pi}_1, \underline{\pi}_2 \text{ transf. by some } \underline{H}!)$$

Transformation: Projective \rightarrow Affine



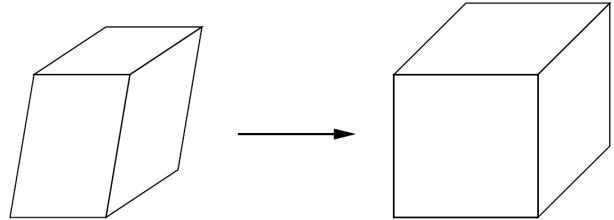
- 1) find vanishing points V_1, V_2, V_3 by intersecting known parallel lines
- 2) find line/plane $L_{\infty}' / \Pi_{\infty}'$ going through V_1, V_2, V_3
- 3) find transform that maps $L_{\infty}', \Pi_{\infty}'$ to $L_{\infty} = (0,0,1)$, $\Pi_{\infty} = (0,0,0,1)$
 $\hookrightarrow (T_{PA}^{-1})^T \cdot L_{\infty}, \Pi_{\infty} = (0,0,0,1)$

Transformation: Affine \rightarrow Similarity

get $C_{\infty}^*, Q_{\infty}^*$ by finding known angles/length-ratios constraints (3 for 2D, 6 for 3D) and solving.

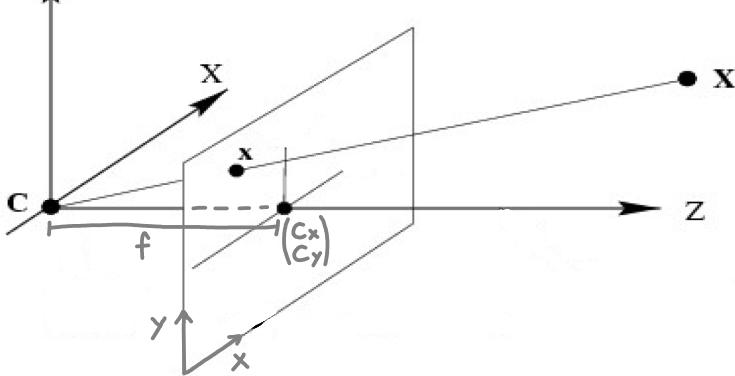
$$C_{\infty}^*, Q_{\infty}^* = \begin{pmatrix} A & \underline{\alpha} \\ \underline{\Omega}^T & 1 \end{pmatrix} \begin{pmatrix} I & O \\ \underline{\Omega}^T & 0 \end{pmatrix} \begin{pmatrix} A & \underline{\alpha} \\ \underline{\Omega}^T & 1 \end{pmatrix}^T = \begin{pmatrix} AA^T & \underline{\Omega} \\ \underline{\Omega}^T & 0 \end{pmatrix} \rightarrow T_{AM} = \begin{pmatrix} A^{-1} & O \\ \underline{\Omega}^T & 0 \end{pmatrix}$$

?!



Pinhole-camera model (3D \rightarrow 2D)

$$\begin{pmatrix} x \\ y \\ w \end{pmatrix} = \underbrace{\begin{pmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix}}_K \underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}}_{P_R} \underbrace{\begin{pmatrix} R & t \\ \underline{\Omega}^T & 1 \end{pmatrix}}_{T_{CAM}} \underbrace{\begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}}_P$$



f: image plane dist

p_x, p_y : width/height of pixel

c_x, c_y : coord (2D) of center pixel

α : skew angle of pixel

$$\hookrightarrow f_x = f/p_x, f_y = f/p_y, s = \tan \alpha \cdot f/p_y$$

K: calibration matrix of the camera

P_e : projection to plane $z=1$

T_{CAM} : move world in camera frame

$$\hookrightarrow \text{if } R, t \text{ move camera instead} \rightarrow T_{CAM} = \begin{pmatrix} R^T & -R^T t \\ \underline{\Omega}^T & 1 \end{pmatrix}$$

P: projection mat. $P = K \cdot P_e \cdot T_{CAM}$

pseudo inv.

project point: $x = P X$

project line: $L = P X_1 \times P X_2$ (X_1, X_2 : line)

project plane: ?

project quadric: $C^* = P Q_{\infty}^* P^T$
 \hookrightarrow degenerate!

Direct linear transform (DLT) (get camera proj. \underline{P} from pairs $\underline{x}_i, \underline{x}'_i$ (inhomog.!!))

- algebraic error: $(\underline{P} \in \mathbb{R}^{3 \times 4} \rightarrow 11 \text{DOF} \rightarrow 6+ \text{ points needed})$

$$\begin{aligned} x_i &= \frac{p_{00}X_i + p_{01}Y_i + p_{02}Z_i + p_{03}}{p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}} \Rightarrow p_{00}X_i + p_{01}Y_i + p_{02}Z_i + p_{03} - x_i(p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}) = 0 \\ y_i &= \frac{p_{10}X_i + p_{11}Y_i + p_{12}Z_i + p_{13}}{p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}} \Rightarrow p_{10}X_i + p_{11}Y_i + p_{12}Z_i + p_{13} - y_i(p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}) = 0 \end{aligned}$$

$\downarrow 6+$ of those pairs of equations

minimize $\|\underline{A}\underline{p}\|$: $\Leftarrow \underline{A} \cdot \left(\underbrace{\underline{p}_{00}, \underline{p}_{01}, \underline{p}_{02}, \underline{p}_{03}}_{\underline{P}_0^T}, \underbrace{\underline{p}_{10}, \underline{p}_{11}, \underline{p}_{12}, \underline{p}_{13}}_{\underline{P}_1^T}, \underbrace{\underline{p}_{20}, \underline{p}_{21}, \underline{p}_{22}, \underline{p}_{23}}_{\underline{P}_2^T} \right)^T = 0$

singular-value decomp. of \underline{A}
 $\underline{A} = \underline{U} \cdot \underline{S} \cdot \underline{V}^T \rightarrow \underline{p} = \text{rightmost col of } \underline{V}$
 (with smallest S entry)



\underline{P} found. now find $\underline{K}, \underline{R}, \underline{t}$ with $\underline{P} = \underline{K}(\underline{R} | \underline{t}) = (\underline{K} \underline{R} | \underline{K} \underline{t})$

- RQ decomposition of $M = (KR)^{-1} \rightarrow M = R^* \cdot Q^* \rightarrow K = R^*, R = Q^*$
- get camera position \underline{C} : solve $\underline{P} \underline{C} = 0$ (same way as solving $\underline{A} \underline{p} = 0$) $\rightarrow \underline{t} = -R \cdot \underline{C}$

- algebraic error + normalization:

normalize every coord $(x_i, y_i, X_i, Y_i, Z_i)$ so that $\frac{1}{N} \sum [\text{coord}] = 0$, $\frac{1}{N} \sum [\text{coord}]^2 = 1$

$\hookrightarrow \underline{x}_n = \underline{T} \cdot \underline{x}$, $\underline{X}_n = \underline{U} \cdot \underline{X}$ → get \underline{P}_n with above steps → $\underline{P} = \underline{T}^T \underline{P}_n \underline{U}$ → continue above steps.

- geometric error (gold standard):

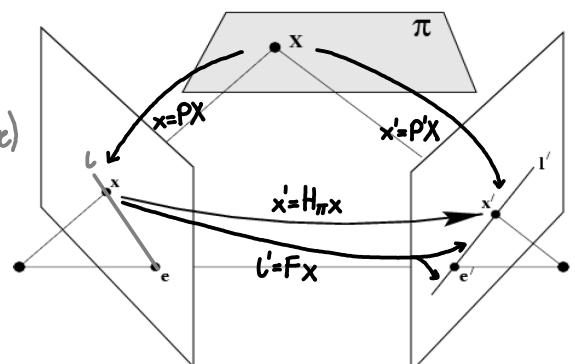
same as algebraic, but minimize $\sum \text{dist}(\underline{x}_i, \underline{P} \underline{X}_i)^2$

$$\underline{X}(\lambda) = \underline{P}^T \underline{X} + \lambda \underline{C}$$

Two-view geometry (find point ↔ line mapping of 2 cameras given matching pairs $\underline{x}_i, \underline{x}'_i$)

$$\begin{aligned} \underline{l}' = \underline{F} \underline{x} &\quad , \quad \underline{l} = \underline{F}' \underline{x}' \\ \underline{x}^T \underline{F} \underline{x} = 0 &\quad , \quad \underline{x}^T \underline{F}' \underline{x}' = 0 \end{aligned}$$

plücker matrix
 $\underline{F} = [\underline{e}']_x H_{\pi} = [\underline{e}']_x P' P^T$ (fundamental matrix)
 $\Rightarrow \text{DoF}$
 $\underline{F}' = \underline{F}^T$, $\underline{e}^T \underline{F} = 0$, $\underline{F} \underline{e} = 0$
 $e, e': \text{epipoles}$, $l, l': \text{epilines}$



- 8 point alg: (find $\underline{F}, \underline{P}'$ from 8 matches $\underline{x} \leftrightarrow \underline{x}'$, $\underline{P} = [\underline{I} | \underline{0}]$)

$$\underline{x}^T \underline{F} \underline{x} = 0 \rightarrow x' x f_{11} + x' y f_{12} + x' z f_{13} + y' x f_{21} + y' y f_{22} + y' z f_{23} + z' x f_{31} + z' y f_{32} + z' z f_{33} = 0$$

$$\hookrightarrow \left[\begin{matrix} x' x, x' y, x' z, y' x, y' y, y' z, z' x, z' y, z' z \end{matrix} \right] \cdot \left[\begin{matrix} f_{11}, f_{12}, f_{13}, f_{21}, f_{22}, f_{23}, f_{31}, f_{32}, f_{33} \end{matrix} \right]^T = 0$$

$$\underline{A} \cdot \underline{f} = 0$$

\hookrightarrow minimize $\|\underline{A} \underline{f}\|$: SVD of $A \rightarrow U \cdot S \cdot V^T \rightarrow$ rightmost col of V

$$\hookrightarrow \underline{P} = [\underline{I} | \underline{0}], \underline{P}' = [\underline{e}]_x \underline{F} + \underline{e}' \underline{V}^T | \lambda \underline{e}']$$

(\underline{V}, λ : any value!!)
 \hookrightarrow pick with e.g. known angles between walls

normalization:

$\underline{x}_n = \underline{T} \cdot \underline{x}$, $\underline{x}'_n = \underline{T}' \cdot \underline{x}'$ so that img bottom-left = (-1, -1), top-right = (1, 1)

\hookrightarrow get F_n like left.

$$\rightarrow F = T'^T F_n T$$

• 7 point alg: (1 point less + use $\det(F)=0$)

same as 8 point, but take 2 rightmost rows of $V \rightarrow F_1, F_2 \rightarrow$ find λ with $\det(F_1 + \lambda F_2)$ by: $\lambda = EW(F_2^T F_1) \rightarrow F = F_1 + \lambda F_2$
(alternative F from 8-point $\rightarrow USV^T = \text{svd}(F) \rightarrow$ set $S_{33}=0 \rightarrow$ recreate $F = USV^T$)

• 5 point alg: (3 point less + use $\det(F)=0$, $2FF^TF - \text{trace}(FF^T)F=0 \Leftrightarrow \lambda_1=\lambda_2$ and $\lambda_3=0$)

• known up: (constraint from knowing how camera up changed from P to P')

222

???

• 3D point \underline{X} constraint:

→ stack any 4 rows to \underline{X}

$$\text{matching 2D point: } \begin{pmatrix} w \cdot x \\ w \cdot y \\ w \end{pmatrix} = \begin{pmatrix} P_1 \\ P_2 \\ P_3 \end{pmatrix} \cdot \underline{X} \rightarrow \begin{pmatrix} P_3 \cdot x - P_1 \\ P_3 \cdot y - P_2 \end{pmatrix} \cdot \underline{X} = 0 \rightarrow \begin{pmatrix} -A_1 \\ -A_2 \\ -A_3 \\ -A_4 \end{pmatrix} \cdot \underline{X} = 0 \rightarrow \boxed{\text{if constraint is met: } \det(\underline{A}) = 0}$$

Factorization (reconstruct camera projections \underline{P}_i ($i=1..m$) and 3D points \underline{X}_j ($j=1..n$), given all 2D points \underline{x}_{ij})

• orthographic factorization: (assume \underline{P}_i is an orthographic projection)

$$\begin{array}{l} \underline{x}_{ij} = \begin{pmatrix} \tilde{x}_{ij} \\ \tilde{y}_{ij} \end{pmatrix}, \quad \underline{P}_i = \begin{pmatrix} P_i^x \\ P_i^y \\ P_i^z \end{pmatrix}, \quad \underline{X}_j = \begin{pmatrix} X_j \\ Y_j \\ Z_j \\ 1 \end{pmatrix} \rightarrow \underline{x}_{ij} = \underline{P}_i \underline{X}_j \\ \downarrow \quad \downarrow \quad \downarrow \\ \underline{x} = \begin{pmatrix} \underline{x}_{11} & \dots & \underline{x}_{1n} \\ \vdots & \ddots & \vdots \\ \underline{x}_{mn} & \dots & \underline{x}_{nn} \end{pmatrix}, \quad \underline{P} = \begin{pmatrix} \underline{P}_1 \\ \vdots \\ \underline{P}_m \end{pmatrix}, \quad \underline{X} = (\underline{X}_1, \dots, \underline{X}_n) \rightarrow \underline{x} = \underline{P} \underline{X} \end{array} \quad \left| \begin{array}{l} \text{rank}(\underline{x}) = 3 \\ \text{svd: } \underline{x} = U \Sigma V^T \rightarrow \text{set all EW in } \Sigma \text{ to 0 except largest 3} \\ \text{reconstruct affine } \tilde{\underline{P}}, \tilde{\underline{X}}: \tilde{\underline{P}} = U, \tilde{\underline{X}} = \Sigma V^T \\ \text{solve for } \underline{C}: \tilde{\underline{P}}^x \underline{C} \tilde{\underline{P}}^x = 1, \tilde{\underline{P}}^y \underline{C} \tilde{\underline{P}}^y = 1, \tilde{\underline{P}}^z \underline{C} \tilde{\underline{P}}^z = 0 \\ \text{solve for } \underline{A}: \underline{C} = \underline{A}^{-1} \cdot \underline{A}^T \text{ (Cholesky factorization)} \\ \text{reconstruct metric } \underline{P}, \underline{X}: \underline{P} = \tilde{\underline{P}} \underline{A}^{-1}, \underline{X} = \underline{A} \tilde{\underline{X}} \end{array} \right. \quad \text{(force rank 3 constraint.)}$$

• perspective factorization:

$$\begin{array}{l} \underline{x}_{ij} = \begin{pmatrix} \tilde{x}_{ij} \\ \tilde{y}_{ij} \\ 1 \end{pmatrix}, \quad \underline{P}_i = \begin{pmatrix} P_i^x \\ P_i^y \\ P_i^z \end{pmatrix}, \quad \underline{X}_j = \begin{pmatrix} X_j \\ Y_j \\ Z_j \\ 1 \end{pmatrix} \rightarrow \lambda_{ij} \underline{x}_{ij} = \underline{P}_i \underline{X}_j \\ \downarrow \quad \downarrow \quad \downarrow \\ \underline{x} = \begin{pmatrix} \lambda_{11} \underline{x}_{11} & \dots & \lambda_{1n} \underline{x}_{1n} \\ \vdots & \ddots & \vdots \\ \lambda_{mn} \underline{x}_{mn} & \dots & \lambda_{nn} \underline{x}_{nn} \end{pmatrix}, \quad \underline{P} = \begin{pmatrix} \underline{P}_1 \\ \vdots \\ \underline{P}_m \end{pmatrix}, \quad \underline{X} = (\underline{X}_1, \dots, \underline{X}_n) \rightarrow \underline{x} = \underline{P} \underline{X} \end{array} \quad \left| \begin{array}{l} (\lambda_{ij}: \text{factor from w-coord } w_{ij} = \underline{P}_i \underline{X}_j, \text{ UNKNOWN}) \\ 1) \text{ assume } \lambda_{ij} = 1 \forall (i,j) \\ 2) \text{ calculate } \underline{P}_i, \underline{X}_j \text{ like orthographic (force rank 4 constraint.)} \\ 3) \text{ if 5th EW of } \Sigma \text{ was small } \rightarrow \text{stop.} \\ 4) \text{ estimate new } \lambda_{ij} \text{ from } \lambda_{ij} \underline{x}_{ij} = \underline{P}_i \underline{X}_j \\ 5) \text{ goto 2.} \end{array} \right.$$

• refine results of above methods: minimize reprojection error $\min_{\underline{P}, \underline{X}} \sum_{ij} \| \underline{x}_{ij} - \underline{P}_i \underline{X}_j \|_2$ (bundle adjustment)

feature definition

local feature = properties of a small patch of an image used to describe an object

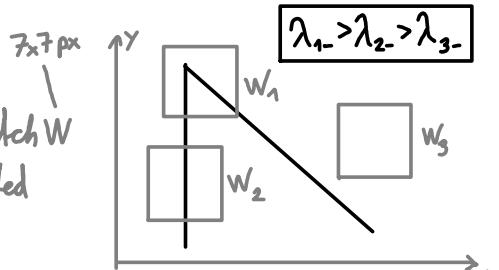
- has to be: invariant to changes in viewpoint, illumination, background, occlusion, ...
- has to be: unique in comparison to other features.

detection (finds feature positions in an image)

• standard corner detector:

$$E(u,v) = \sum_{(x,y) \in W} (I(x+u, y+v) - I(x, y))^2 \\ \approx \sum_{\underline{H}} [u \ v] \begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

$I(x,y)$: intensity of image
 $(x,y) \in W$: all pixels in a small patch W
 (u,v) : direction that patch W is shifted
 $(I_x(x,y)) = \frac{\partial I(x,y)}{\partial x}$



- $EV(H)_+ = \lambda_+$: (u,v) with largest increase in E
- $EV(H)_- = \lambda_-$: (u,v) with smallest increase in E
- $EW(H)_+ = \lambda_+$: amount of increase in λ_+
- $EW(H)_- = \lambda_-$: amount of increase in λ_-

$$(I_x(x,y) = \frac{I(x+1,y) - I(x-1,y)}{2}; I_y(x,y) = \frac{I(x,y+1) - I(x,y-1)}{2}; \text{ points from } \lambda_-(x,y): \text{local maxima} \rightarrow \text{threshold})$$

• Harris corner detector:

instead of $\lambda = EW(H)_-$ do: $\lambda \approx \det(\underline{H}) - k \cdot \underbrace{(h_{11} + h_{22})^2}_{\text{trace}(H)} = \lambda_+ \lambda_- - k(\lambda_+ + \lambda_-)$ ($k=0,04 \dots 0,06$)

• difference of gaussian (DoG) corner detector: ??

• MSER blob detector:

slowly raise threshold intensity → when a region grows the slowest it's a blob (bounded by strong edges)

• LoG blob detector: smooth image with gaussian filter → do laplacian → mark on center of blobs

• ANMS adaptive non-maxima suppression: ??

description (gets values describing the patch around a detection)

• transformation invariants:

- similarity (3 Points): $\frac{\|A-B\|}{\|A-C\|} = \text{const.}$
- affinity (4 Point): $\frac{|A-B \times A-C|}{|A-B \times A-D|} = \text{const.}$
- projectivity (5 Point): $\frac{|A-B \times A-E|}{|A-C \times A-E|} / \frac{|B-D \times B-E|}{|C-D \times C-E|} = \text{const.}$

• photometric invariants:

approx of color transf.: $\begin{pmatrix} r' \\ g' \\ b' \end{pmatrix} \approx \begin{pmatrix} S_R & 0 & 0 \\ 0 & S_g & 0 \\ 0 & 0 & S_b \end{pmatrix} \cdot \begin{pmatrix} r \\ g \\ b \end{pmatrix} + \begin{pmatrix} 0_r \\ 0_g \\ 0_b \end{pmatrix}$

if patch transformation matches, then color normalized patches (mean=0, var=1) are invariant.

(if patch is assumed small, similarity holds for affine and even projective transforms!)

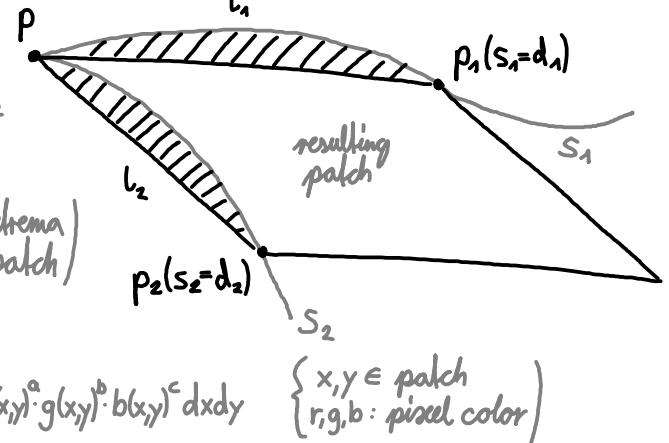
- normalized cross correlation (NCC) (directly compare pixels of 2 patches!)

$$S(u, v) = \frac{\sum_i (u_i - \bar{u})(v_i - \bar{v})}{\sqrt{\sum_i (u_i - \bar{u})^2} \sqrt{\sum_i (v_i - \bar{v})^2}}$$

u_i, v_i : intensities of 2 patches being compared.
 i : usually $1 \dots 49 = 7 \times 7 \text{ px}$ around detected edge.
 \bar{u}, \bar{v} : avg. intensities of 2 patches
 \rightarrow high \rightarrow good match

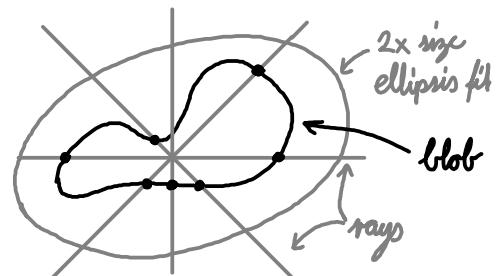
- parallelogram next to edge corner (invariant to: affine transf. + illumination)

- 1) Harris corner detection: for feature position p
- 2) Canny edge detection: find 2 edges of corner s_1, s_2
- 3) find 2 points on 2 edges with same area $l_1 = l_2$
- 4) $l = \int_0^d |dp(s)/ds \times (p(0)-p(s))| ds$ (pick d_1, d_2 to get local extrema)
 $(\text{of avg color intensity in patch})$



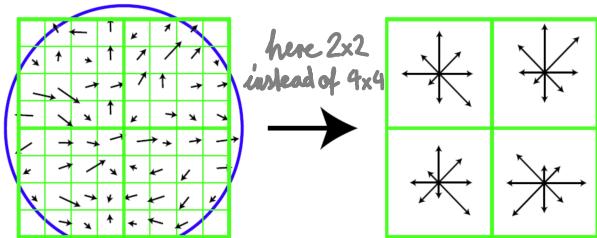
- intensity extrema (blobs) + ellipse patch (invariant to: affine transf. + illumination)

- 1) find blobs / intensity extrema
- 2) find points of max intensity changes along rays
- 3) fit ellipse to points and double its size
- 4) compute invariant of elliptic patch (e.g. D_{02} from above)



- SIFT (scale-invariant feature transform) (invariant to similarity + illum.)

- 1) use LoG method to find blobs
- 2) get square patch around blob and compute gradients
- 3) subdivide (square) patch of blob in 4×4 zones
- 4) bin gradient in each zone into 8 principal directions ($4 \times 4 \times 8$ bins)
- 5) normalize + rotate, so strongest dir is up. \uparrow (weight with dist. from center)



- SURF: faster version of SIFT

- HoG (histogram of oriented gradients): like step 4 of SIFT

matching

smallest sum-of-squared-differences (SSD),

- one-way nearest neighbor: for every descriptor of image 1, find best match in image 2
- mutual nearest neighbor: do one-way and one-way with swapped img $1 \leftrightarrow 2$. intersect matches.
- one-way + ratio test: do one-way and find also 2nd best. if $1^{\text{st}} \text{ best} / 2^{\text{nd}} \text{ best} < \sim 0.5 \rightarrow$ discard.

optical flow (find speed of every pixel in a video)

$I(x(t), y(t), t)$: intensity of a pixel that moves ("fixed" to an object) in a video
 $v = \frac{dx}{dt}, u = \frac{dy}{dt}$: momentan speed of a pixel

- $\frac{dI}{dt} = \frac{\partial I}{\partial x} \cdot u + \frac{\partial I}{\partial y} \cdot v + \frac{\partial I}{\partial t} \approx 0$ (assume that moving pixel does not change intensity)
- $(\frac{du}{dx})^2 + (\frac{dv}{dx})^2 + (\frac{du}{dy})^2 + (\frac{dv}{dy})^2 \approx 0$ (assume u, v map is "smooth")

find $u(x, y), v(x, y)$ that minimizes: $\iint (u_x^2 + u_y^2 + v_x^2 + v_y^2) + \lambda (I_x u + I_y v + I_t)^2$

↓ (calculation of variations → Horn-Schunck algorithm)

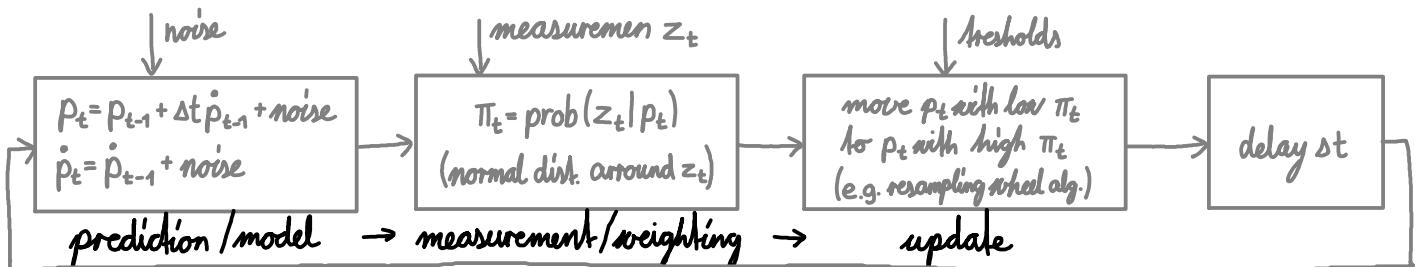
$$\frac{\partial u}{\partial i} = \Delta u - \lambda (I_x u + I_y v + I_t) \cdot I_x ; \quad \frac{\partial v}{\partial i} = \Delta v - \lambda (I_x u + I_y v + I_t) \cdot I_y$$

$$\left\{ \begin{array}{l} \frac{\partial I}{\partial x} = \frac{I(x+1, y, t) - I(x-1, y, t)}{2} \\ \frac{\partial I}{\partial y} = \frac{I(x, y+1, t) - I(x, y-1, t)}{2} \\ \frac{\partial I}{\partial t} = \frac{I(x, y, t+1) - I(x, y, t-1)}{2} \\ \frac{\partial u, v}{\partial x, y} : \text{same as above} \\ \lambda : \text{some const.} \end{array} \right.$$

$$\left\{ \begin{array}{l} \Delta[\cdot] = \frac{\partial[\cdot]^2}{\partial x^2} + \frac{\partial[\cdot]^2}{\partial y^2} \\ i : \text{iteration} \end{array} \right.$$

condensation filter (object tracker for video (similar to Kalman/Bayesian filter))

p_t : many position estimates at time t (mimicking probability distribution) ; π_t : prob. that p_t is true



model fitting

• hough transform:

1) discretize model parameter space

2) count up cells that explain a datapoint

3) find peaks in cell value

! problem if param. space unbounded!

! param. space cell size has to be "right"!

(guess...)

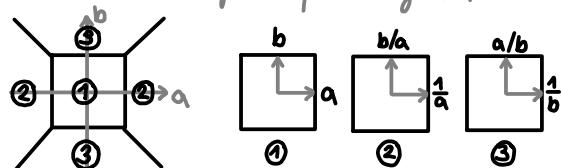
examples with line fitting:

model is $y = ax + b$, data space (x, y) , parameter space (a, b)

all cells (a, b) that explain a point get counted up +1

the cell (a, b) with the highest count is a good model

- use $x \cos \theta + y \sin \theta - g$ instead (still: bounding in g ?)
- use distorted subspaces for large (a, b)



• K-means line fitting (for when there are multiple lines in the data)

1) generate randomly placed lines (guess line count)

2) assign each point to closest line

3) refit lines to assigned points

least-square?

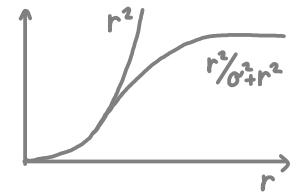
4) go to step 2 many times

5) go to step 1 many times

6) pick the iteration with lowest error.

robust estimation (alternative to least-squares, that's better against outliers)

instead of minimizing $\sum r^2$ (r : dist. datapoint to model), minimize $\sum \frac{r^2}{\sigma^2 + r^2}$



RANSAC:

n : min. # of points needed to determine model
 k : # of algorithm iterations
 t : threshold dist. to determine inliers of model

while less than k iterations

- take n randomly sampled points from data
- fit model to n points
- find all other datapoints within t of model
- refit model to all these points
- pick model with smallest fitting error

$$k: k = \log(1-p) / \log(1-(1-e)^n) \quad (p \approx 0.99)$$

p : probability 1+ sample has all inliers
 e : approximate proportion outlier/total

(recalc. k each iteration if e unknown.)
 $e = \text{lowest outlier/total up to that moment}$

t : empirical / so that X% data included

improvements:

- minimize median residual instead of maximizing inlier count (LMedS)
- sample 'good' points more frequently
- if quasi-degenerate data: do RANSAC reducing n until model starts getting bad
 \rightarrow sort out explainable data \rightarrow use remaining DoF to fit remaining data

segmentation (find regions of an image that belong together)

k-means:

- 1) give features to each pixel (1D: intensity, 3D: color, 5D: r,g,b,x,y,...)
- 2) randomly initialize k cluster centers c_1, \dots, c_k in the feature space
- 3) for each pixel p , put it in the closest cluster c_i with $\min \|p - c_i\|^2$
- 4) move each cluster to be at the mean of all contained pixels.
- 5) repeat steps 3,4 until clusters stop moving.

cons:

- how to choose k ?
- can be "local minimum"
- detects spherical clusters only
- sensitive to outliers
- sensitive to initial cluster locations

mixture of gaussians, EM:

similar to k-means, but clusters are defined by $\theta_b = (\alpha_b, \mu_b, V_b)$

$$R(x, \theta_b) = \alpha_b / \sqrt{(2\pi)^d |V_b|} \cdot \exp\left(-\frac{1}{2}(x - \mu_b)^T V_b^{-1} (x - \mu_b)\right)$$

$$P(x|\theta_b) = R(x, \theta_b) / \sum_{i=1}^k R(x, \theta_i)$$

(μ_b : mean vector, V_b : covariance matrix, $P(x|\theta_b)$: prob. that x is in θ_b)

assign every pixel to the θ_b with highest $P(x|\theta_b)$

given the assigned pixels of θ_b , change μ_b, V_b to maximize $P(x|\theta_b)$

cons:

same as k-mean, but better with outliers.

$$\begin{cases} \alpha_b^{new} = \frac{1}{N} \sum_{i=1}^N P(x_i|\theta_b) \\ \mu_b^{new} = \sum_{i=1}^N x_i \cdot P(x_i|\theta_b) / \sum_{i=1}^N P(x_i|\theta_b) \\ V_b^{new} = \sum_{i=1}^N (x_i - \mu_b^{new})(x_i - \mu_b^{new})^T P(x_i|\theta_b) / \sum_{i=1}^N P(x_i|\theta_b) \end{cases}$$

mean-shift

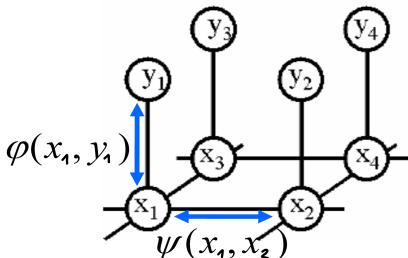
- 1) define feature space like in k-mean
- 2) place windows (e.g.: circle, sphere) on every pixel in the feature space
- 3) for every window: move it to the mean of all contained pixels
- 4) do step 3 until all windows stop moving
- 5) pixels with windows that end up in the same place are in one segment.

cons:

- window size selection difficult
- computationally expensive

Graph-cuts:

- 1) pick k key-features x
- 2) for every pixel/feature y calculate the "dist." to every key-feature $x \rightarrow \varphi(x, y)$
- 3) for every combination of 2 key-features define a "distance" $\psi(x_i, x_j) \rightarrow \psi(x_1, x_2)$
- 4) assign every pixel to a key-feature while minimizing $E = \sum \varphi + \sum \psi$
 ↪ use graph-cuts to find division (cuts) of regions with same key-feature



cons:

- how to pick key-features x

k-nearest neighbor (KNN) (!needs training dataset !)

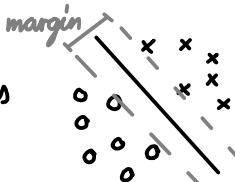
for every pixel find the nearest k points of the training dataset in featurespace, and assign it to a segment based on the k points.

cons:

- hard to choose k
- needs large training dataset

support vector machine (SVM):

find hyperplane (line in 2D case) that divides the featurespace and has maximal margin



cons:

- only 2 clusters
- not always solvable

random forest

- 1) build many random binary decision trees (=forest)

(feature-space description of pixel is input to the root. each node has a feature dimension and a threshold to decide if to go right or left. the leafs correspond to segments that the tree chooses.)

cons:

- lots of parameters to choose
- needs large training dataset
- training the forest is slow

- 2) train the trees by changing node parameters / killing, growing branches to make them predict a testing dataset. **how?**
- 3) use the average output of the trees to predict the segment of new pixels.

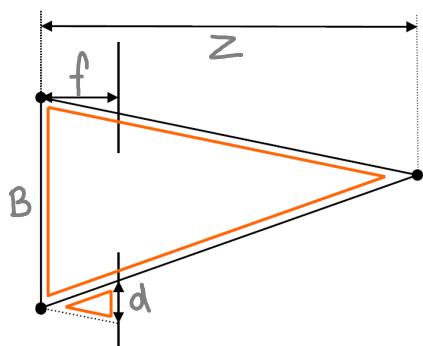
non-linear least-squares (find solution to eq. that minimizes error)

$X = f(P)$ (X : solution, f : some non-linear function, P : variable for solving the equation)

• Newton method: $P_0 = \text{guess.} \rightarrow P_{i+1} = P_i + \Delta, \Delta = (J^T J)^{-1} J^T (X - f(P_i)) , J = \frac{\partial f}{\partial P}$ $\leftarrow (f(P_0 + \Delta) \approx f(P_0) + J \cdot \Delta)$

• Levenberg-Marquardt method: same except $\Delta = (J^T J + \lambda \text{diag}(J^T J))^{-1} J^T (X - f(P_i))$ $\lambda_0 = 10^3 \rightarrow \text{how to change } \lambda \text{ ??}$

Stereo matching (find point correspondences in 2 images to compute depth + 3D model)



(cameras are identical and point in same direction)

$$F = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & -1 & 0 \end{pmatrix} \quad \begin{array}{l} \text{(fundamental matrix)} \\ \text{(defining B line as x-axis)} \end{array} \rightarrow \text{epipolar lines = horizontal}$$

for every feature in camera 1 search for a match on the corresponding epipolar line and get the offset d . (see below for methods)

$$\rightarrow \text{depth } Z = B \cdot f / d$$

- **feature match constraint:** 1 feature should have 1 or less matches. see above chapter for details.
- **ordering constraint:** order of matched features along epipolar lines should not be allowed to change. (only allow to skip features, not reorder!)
- **disparity constraint:** only realistic depths + depth gradients need to be considered.
 - foreshortening distortion of feature windows can be compensated after a rough depth map estimate
 - feature matches along epipolar lines can be skipped due to occlusion.
 - feature window matching improves by considering interpolated windows (between pixel)
 - window on background pixel containing foreground ends up in foreground!
 - by shifting window off-center, then taking best offset (not always, since depth accuracy loss)
 - **compact windows, rod-shaped windows, weight window based on similarity + proximity**
 - match each pixel by time varying illumination of the 2 cameras (if available)
 - depth map with semi-global optimization: **graph cuts, belief propagation, dynamic programming**
 - image pair rectification can make epipolar lines of any 2 img become horizontal = scan lines = fast.
 - a third camera is useful for verification...
 - if baseline dist. B can be chosen: too small \rightarrow depth error, too large \rightarrow much occlusion (use multiple helps)
 - plane-sweep stereo matching: assume depth \rightarrow render images on top each other \rightarrow find where result is sharp \rightarrow assign the assumed depth to that region.

Structure from motion (pipeline to get dense 3D point cloud from images)

- do feature detection+matching (with segmentation) to get sparse 2D point matches
- do factorization to get camera projections P + sparse 3D points between all image pairs
- discard outlier cameras + improve P matrices by enforcing:
 - 3D points of same feature in different views match + loops of relative projections have to return to initial camera pos.
- do stereo matching with each image pair along its epipolar lines to get dense 3D points
- somehow combine all stereo match results into one consistent dense 3D point cloud

specific object recognition (find pose+position of an object in an image if it exists)

(need decent slides to be uploaded!)

• model based (use a 3D mesh or feat. desc. model)

1) get features (best if invariant) from image

2a) try to position the model such that similar features appear

2b) or try to describe object with one feature descriptor and use that as model/comparison

↪ compact model + good with clutter

• hybrid techniques (combination of model based + image based)

- use features (best if invariant) of a few training images in different poses as model.

- detect features in image and determine match to model by counting feature matches

• image based (use images of object in different poses)

directly compare image pixels to all possible pose images to determine pose (position determined a priori with clear background + bounding box)

↪ fast evaluation, easy to produce model

Object class recognition (recognize class of object in image, e.g. cow, car,...)

• bag of words:

Training: 1) extract features (e.g. SIFT) from training images and assign each to its object class

2) segment feature space (e.g. mean shift) and save cluster centers = words

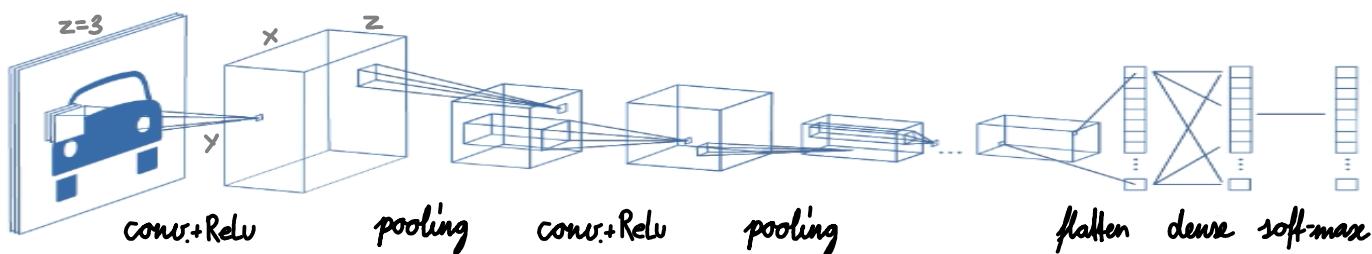
Usage: 1) get features from target image and assign them the class of a near word (e.g. k-nearest-neighbor)

2) let all feature classes vote on what is the most probable object class

! spatial information is not regarded! :    all match. "solutions":

- include pos. in feat. desc.
- phrases (co-occurring words)
- evaluate config. afterwards
- vector from words to obj center

• convolutional neural network (CNN): does feature extraction + image segmentation in one.



-(conv.) convolution layer: apply many conv. kernels to each previous channel pixel (grows in z)
kernel parameters are trained by backpropagation

-(ReLU) rectified layer: apply non-linear func to everything (so that CNN can capture non-linearities)
e.g. sigmoid, $x = x \text{ if } x > 0; 0 \text{ else, ...}$

- pooling: reduce x,y size by some factor by taking e.g. only the max value. (+- compensate for z growth)

- flatten: flatten $x \times y \times z$ cube to $1 \times (x \cdot y \cdot z)$ (necessary for dense layers)

- dense layer: multiply input with a dense matrix (linear combination of all inputs to output)

The dense matrix is trained by backpropagation like kernels

- soft-max: normalize input to sum=1 (so the output can be interpreted as % probability)

→ if trained correctly CNN outputs probability of object class (e.g. 70% cat, 30% dog)

↪ needs MANY training samples (can use pretrained kernels of different CNN to speed up training/require less samples!)

- **boosting classifier**: use many classifiers that are better than chance to build a new one that's better than all other
 - ↳ how? see Ada-Boost algorithm / Viola-Jones face detector (**not in these notes!**)

Object detection (like classification, but also determine location in image)

- **sliding window**: move detection window around in image and do class recognition at every move (SLOW!)

- **R-CNN / fast R-CNN**:

- 1) use selective search to get ~2000 bounding box candidates

do image segmentation with many techniques

→ generate box candidates from segments.

- 2) run each bounding box candidate through a CNN like for object class recognition.

(optionally add bounding box regression where CNN is set up to propose corrections to bbox)

fast R-CNN : run most of the CNN on the whole image and only then cut out bbox candidates and run rest of CNN on these sub- "images"