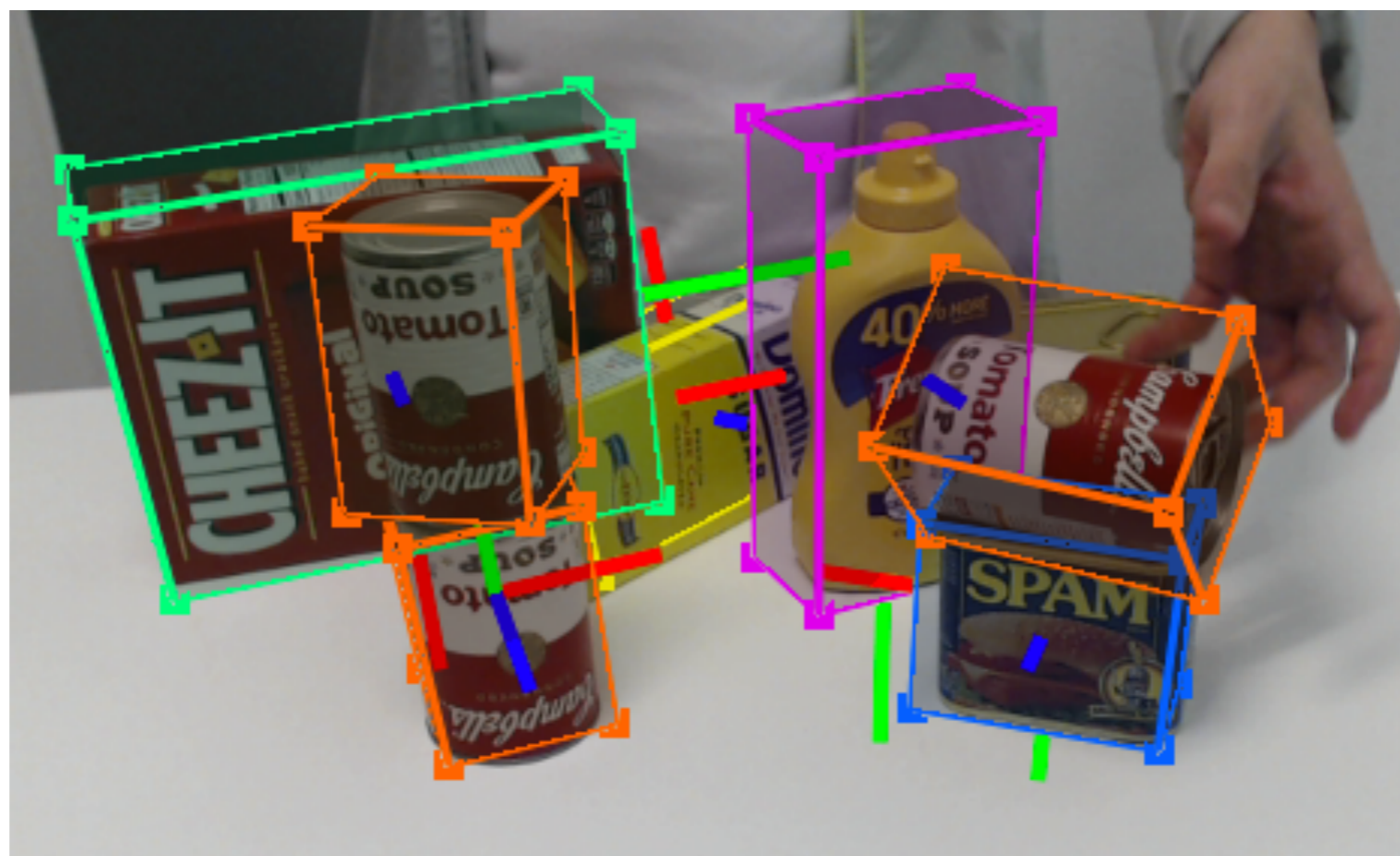


6.2 Pose Estimation

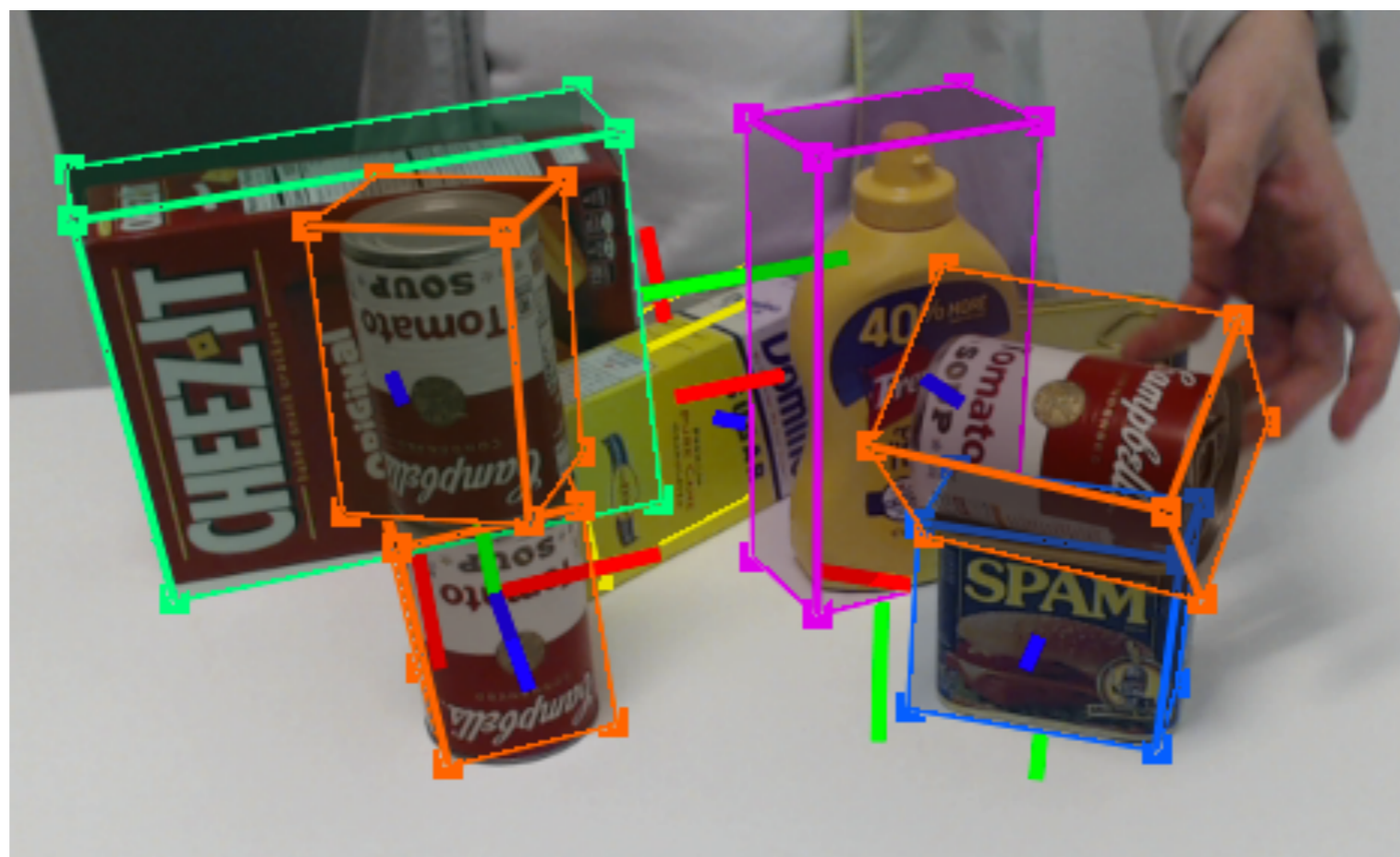
Problem Definition

- ❖ Given:
 - A 3D model of an object
 - An image of the object
- ❖ Estimate:
 - The 3D pose of the object relative to the camera



Perspective 3-Point Problem

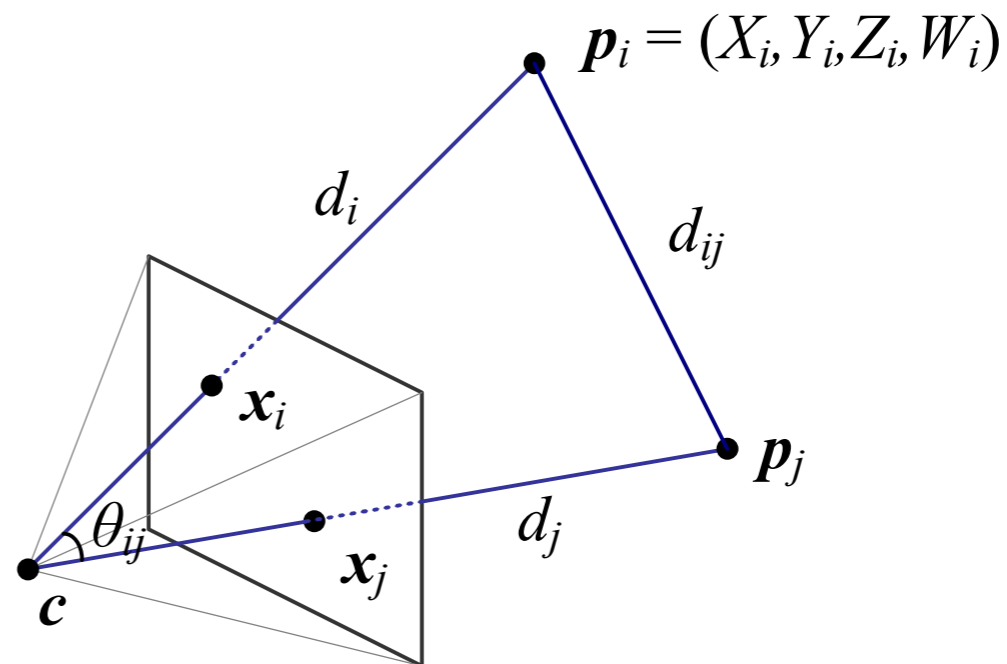
- ❖ How many degrees of freedom (parameters) are we estimating?
- ❖ How many point correspondences between 3D object and 2D image do we need?



Linear Algorithms

❖ 3 x 4 camera projection matrix P

$$\tilde{x}_s = K \left[R \mid t \right] \bar{p}_w = P \bar{p}_w$$



$$P = \begin{bmatrix} p_{00} & p_{01} & p_{02} & p_{03} \\ p_{10} & p_{11} & p_{12} & p_{13} \\ p_{20} & p_{21} & p_{22} & p_{23} \end{bmatrix}$$



$$x_i = \frac{p_{00}X_i + p_{01}Y_i + p_{02}Z_i + p_{03}}{p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}}$$

$$y_i = \frac{p_{10}X_i + p_{11}Y_i + p_{12}Z_i + p_{13}}{p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}}$$

Linear Algorithms

$$x_i = \frac{p_{00}X_i + p_{01}Y_i + p_{02}Z_i + p_{03}}{p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}}$$

$$y_i = \frac{p_{10}X_i + p_{11}Y_i + p_{12}Z_i + p_{13}}{p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}}$$

- ❖ As for estimation of 2D homographies, we can form a linear estimate of the parameters p_{ij} by multiplying through by the denominator, which yields

$$\begin{bmatrix} X_i & Y_i & Z_i & 1 & 0 & 0 & 0 & 0 & -x_iX_i & -x_iY_i & -x_iZ_i & -x_i \\ 0 & 0 & 0 & 0 & X_i & Y_i & Z_i & 1 & -y_iX_i & -y_iY_i & -y_iZ_i & -y_i \end{bmatrix} \begin{bmatrix} p_{00} \\ p_{01} \\ \vdots \\ p_{23} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

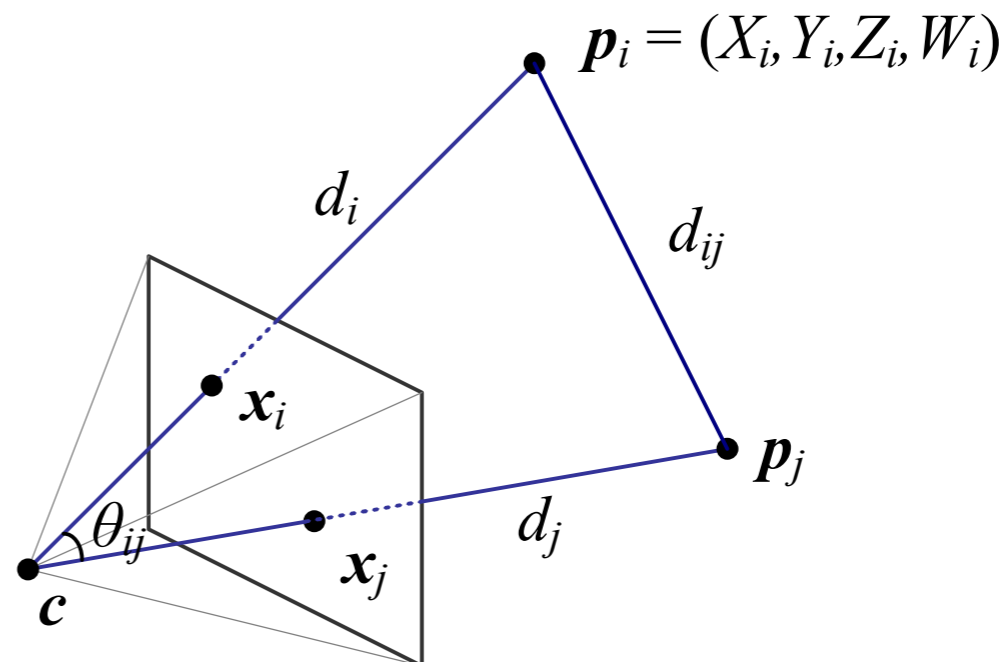
- ❖ How many pairs of matching points do we need?
- ❖ Again, this estimate does not minimize the squared deviation but can be used as an initial guess for an iterative solution.

Linear Algorithms

- ❖ 3 x 4 camera projection matrix P

$$\tilde{x}_s = K \left[R \mid t \right] \bar{p}_w = P \bar{p}_w$$

$$P = \begin{bmatrix} p_{00} & p_{01} & p_{02} & p_{03} \\ p_{10} & p_{11} & p_{12} & p_{13} \\ p_{20} & p_{21} & p_{22} & p_{23} \end{bmatrix}$$



- ❖ Once P has been estimated, its constituents K , R and t can be recovered.
- ❖ Recall that R is orthonormal and K is normally treated as upper triangular:

$$K = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

f_x and f_y : encode focal length and pixel spacing, which may be slightly different in x and y dimensions.

c_x and c_y : encode principal point (intersection of optic axis with sensor plane) - usually very close to centre of image

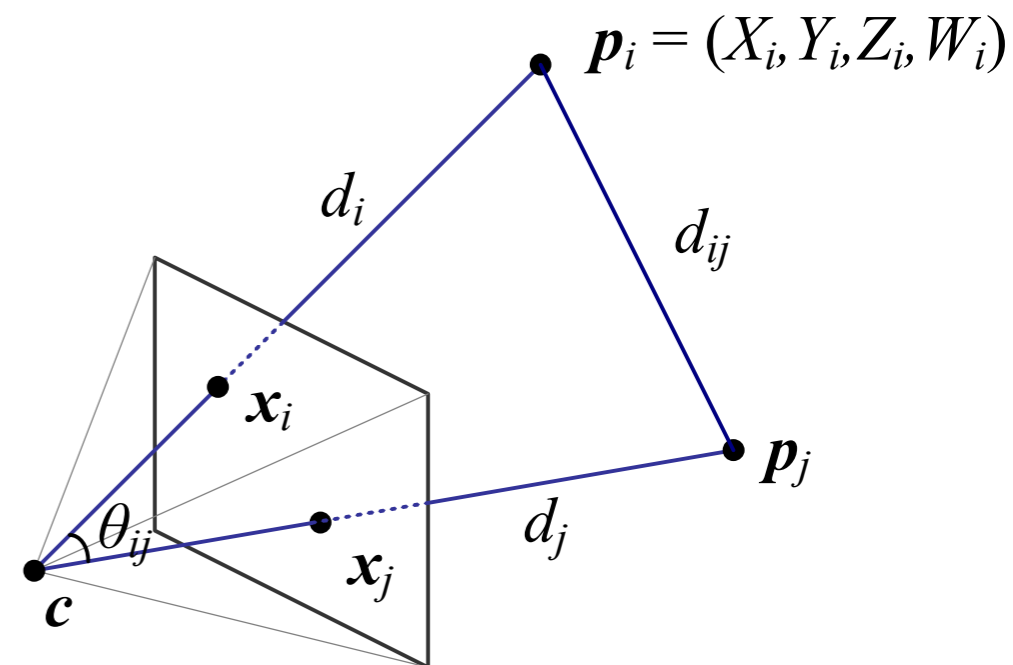
s : encodes possible skew between sensor axes (usually close to 0).

Linear Algorithms

$$\tilde{\mathbf{x}}_s = \mathbf{K} \left[\mathbf{R} \mid \mathbf{t} \right] \bar{\mathbf{p}}_w = \mathbf{P} \bar{\mathbf{p}}_w$$

$$\mathbf{P} = \begin{bmatrix} p_{00} & p_{01} & p_{02} & p_{03} \\ p_{10} & p_{11} & p_{12} & p_{13} \\ p_{20} & p_{21} & p_{22} & p_{23} \end{bmatrix}$$

$$\mathbf{K} = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$$



- ❖ Thus \mathbf{K} and \mathbf{R} can be recovered from the first 3 columns of \mathbf{P} using QR decomposition.

Complexity: $O(MN^2 + N^3)$ for an $M \times N$ matrix (3×4 in our case).

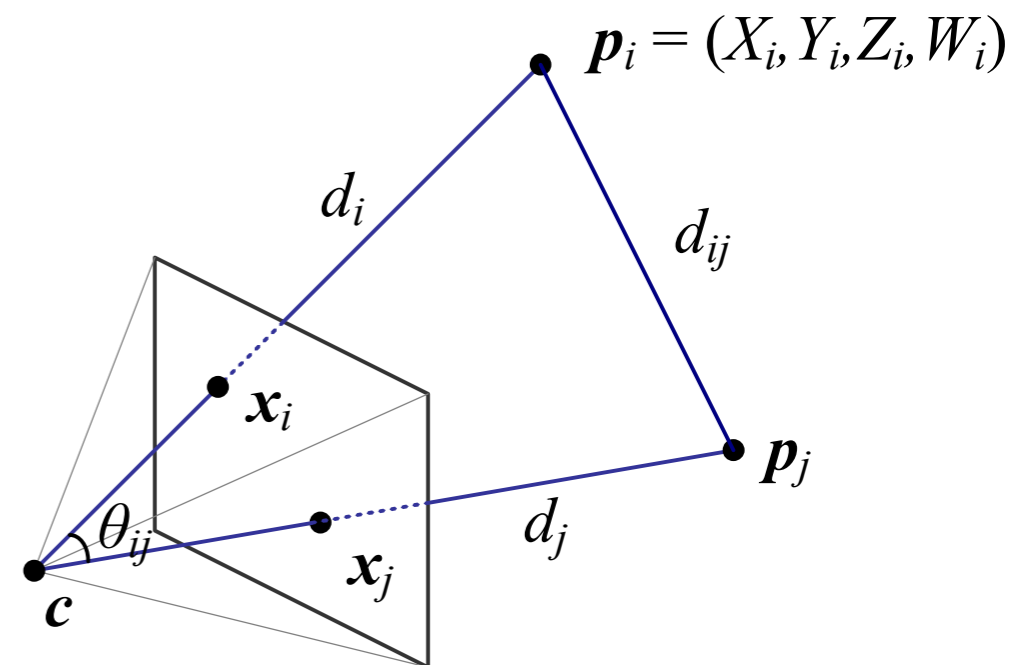
MATLAB function qr(A)

Linear Algorithms

$$\tilde{\mathbf{x}}_s = \mathbf{K} \left[\mathbf{R} \mid \mathbf{t} \right] \bar{\mathbf{p}}_w = \mathbf{P} \bar{\mathbf{p}}_w$$

$$\mathbf{P} = \begin{bmatrix} p_{00} & p_{01} & p_{02} & p_{03} \\ p_{10} & p_{11} & p_{12} & p_{13} \\ p_{20} & p_{21} & p_{22} & p_{23} \end{bmatrix}$$

$$\mathbf{K} = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$$



- ❖ Given a calibrated camera (\mathbf{K} known), \mathbf{R} and \mathbf{t} can be recovered with as few as 3 matched points
- ❖ Basic idea: visual angle between any pair of 2D points \mathbf{x}_i and \mathbf{x}_j in the image must be the same as the visual angle between their corresponding 3D points \mathbf{p}_i and \mathbf{p}_j .

Linear Algorithms

- ❖ Basic idea: visual angle between any pair of 2D points \mathbf{x}_i and \mathbf{x}_j in the image must be the same as the visual angle between their corresponding 3D points \mathbf{p}_i and \mathbf{p}_j .

Let $\hat{\mathbf{x}}_i$ represent the unit vector pointing to image point \mathbf{x}_i from the camera centre \mathbf{c} :

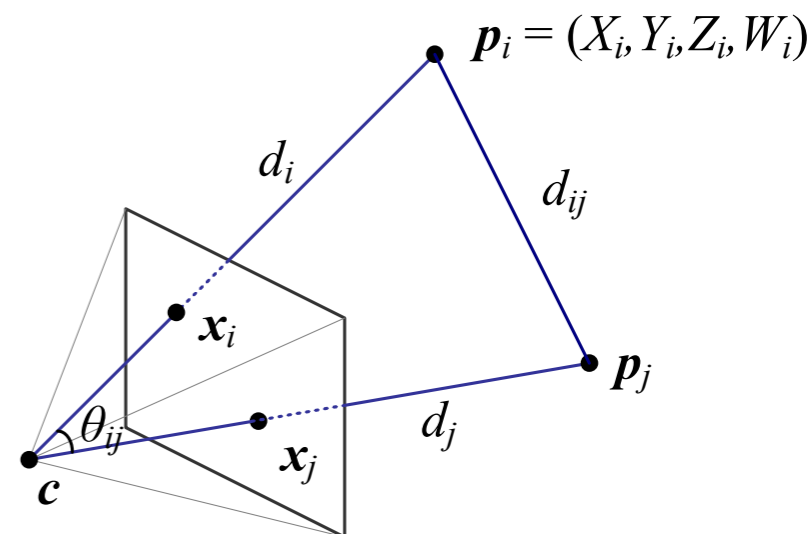
$$\hat{\mathbf{x}}_i = \mathcal{N}(\mathbf{K}^{-1}\mathbf{x}_i) = \mathbf{K}^{-1}\mathbf{x}_i / \|\mathbf{K}^{-1}\mathbf{x}_i\|$$

the unknowns are the distances d_i from the camera origin \mathbf{c} to the 3D points \mathbf{p}_i , where

$$\mathbf{p}_i = d_i \hat{\mathbf{x}}_i + \mathbf{c}$$

The cosine law for triangle $\Delta(\mathbf{c}, \mathbf{p}_i, \mathbf{p}_j)$ gives us

$$f_{ij}(d_i, d_j) = d_i^2 + d_j^2 - 2d_i d_j c_{ij} - d_{ij}^2 = 0,$$



where

$$c_{ij} = \cos \theta_{ij} = \hat{\mathbf{x}}_i \cdot \hat{\mathbf{x}}_j$$

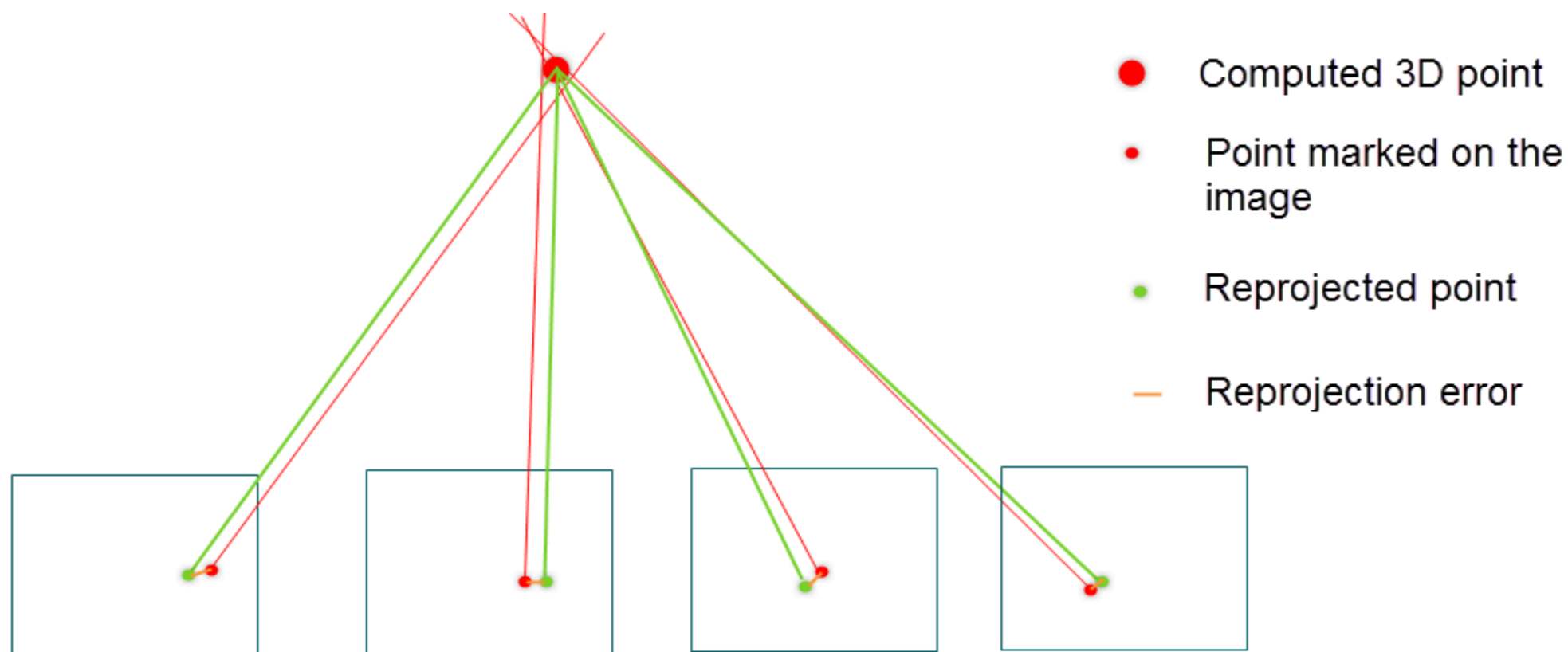
and

$$d_{ij}^2 = \|\mathbf{p}_i - \mathbf{p}_j\|^2.$$

Thus any triplet of constraints $f_{ij}(d_i, d_j), f_{ik}(d_i, d_k), f_{jk}(d_j, d_k)$ generates 3 equations in 3 unknowns.

Iterative Algorithms

- ❖ These minimal linear one-shot algorithms have limitations:
 - ⦿ Noisy (few points)
 - ⦿ Do not directly minimize error
- ❖ Given these limitations, they are most useful as a means to generate an initial guess that can then be refined iteratively to minimize the **reprojection error**.
- ❖ **Definition: Reprojection error**
 - ⦿ The deviation in the image between 2D image points \mathbf{x}_i and their corresponding 3D points \mathbf{p}_i , projected to the image.



Iterative Algorithms

- ❖ Let f now represent projection to the image:

$$x_i = f(p_i; R, t, K)$$

- ❖ We now iteratively minimize a measure of the linearized reprojection error

$$E_{\text{NLP}} = \sum_i \rho \left(\frac{\partial f}{\partial R} \Delta R + \frac{\partial f}{\partial t} \Delta t + \frac{\partial f}{\partial K} \Delta K - r_i \right),$$

where $r_i = \hat{x}_i - \tilde{x}_i$ is the current residual vector (2D error in predicted position)

and

← Sign reversed in textbook.

\hat{x}_i is the 2D image point.

\tilde{x}_i is the current estimate of the projection of 3D point p_i to the image.