



UNIVERSITÀ DEGLI STUDI
DI MILANO

Corso di
Visione Artificiale

Laurea Magistrale in Informatica (F94)

Docente:
Raffaella Lanzarotti

*Dipartimento di Informatica
Università degli Studi di Milano*

Where are we?

First part:

Image formation and Early vision

- Image formation
 - Geometric Camera Models
 - Color spaces
- Image Processing
 - Punctual and spatial processing
 - Feature Extraction
- Reconstruction
 - **Camera calibration**
 - Stereo Vision
 - Structure from Motion and RGB-d Cameras
 - Optical flow and Tracking

Second part:

Machine learning for CV

- Linear Neural Network
- Multi Layer Perceptron
- Convolutional Neural Networks
- Recurrent Neural Networks
- Transformers
- Variational Auto-Encoders
- Generative Adversarial Networks
- Graph Neural Networks
- Self-supervised learning
- Vision Language Models

1. IMAGE FORMATION

Camera Calibration

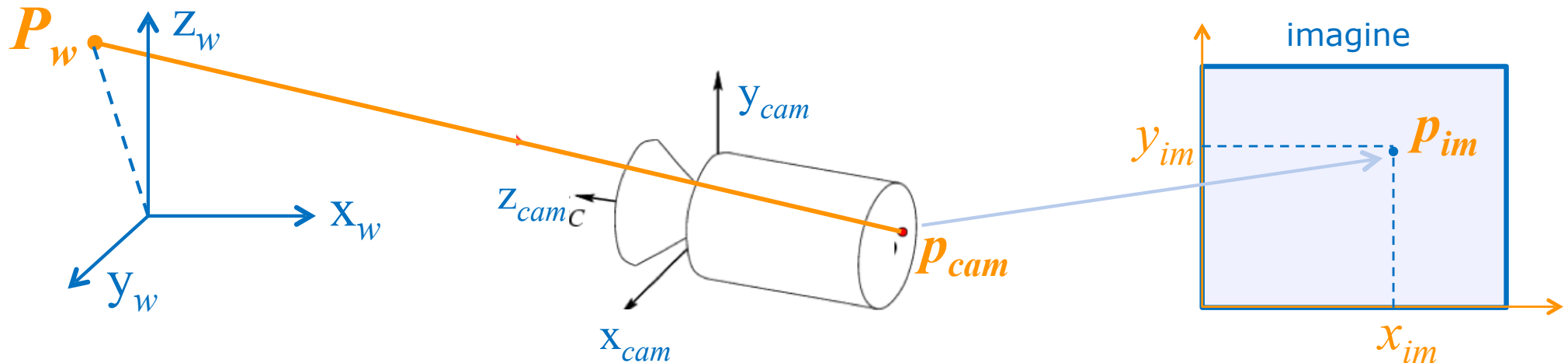
Chapter 1 – Forsyth Ponce

credits, F. Pedersini, S. Nayar

Camera calibration – definition

Camera Calibration:

Process to determine the geometric model of a camera



Projection Matrix: $\mathbf{M}(\xi) : \tilde{p}_{im} = \mathbf{M} \cdot \tilde{P}_w$

Calibration: determine \mathbf{M} (or the camera parameters ξ)

REMARK:

\tilde{x} used to indicate homogeneous coords

RECALL: Complete perspective projection camera model

$$\tilde{\mathbf{p}}_{IM} = \begin{matrix} M_{in} \\ \left[\begin{array}{cccc} f & 0 & x_c & 0 \\ 0 & f & y_c & 0 \\ 0 & 0 & 1 & 0 \end{array} \right] \end{matrix} \cdot \begin{matrix} M_{ext} \\ \left[\begin{array}{c|c} \mathbf{R} & \mathbf{T} \\ \hline \mathbf{0} & 1 \end{array} \right] \end{matrix} \cdot \tilde{\mathbf{P}}_w = \mathbf{M}(\xi) \cdot \tilde{\mathbf{P}}_w$$

- **Linear model in 11 parameters** ($\mathbf{M}_{3 \times 4}$, up to scale)
 - **only 9 params are independent:**

$$\xi = [\mathbf{R}, \mathbf{T}, f, \mathbf{C}] = [\varphi, \vartheta, \rho, t_x, t_y, t_z, f, x_c, y_c]$$

- **Extrinsic Parameters:** depend on the relative position camera-scene
 - **Rotation:** Euler angles: $\mathbf{R} = [\varphi, \theta, \rho]$
 - **Translation:** translation vector: $\mathbf{T} = [t_x, t_y, t_z]$
- **Intrinsic Parameters:** depend on the camera characteristics
 - **Focal length:** f
 - **Optical Centre position:** $\mathbf{C} = \langle x_c, y_c \rangle$

Camera calibration – Calibration pattern

Set of fiducial points: easy to locate, with high precision

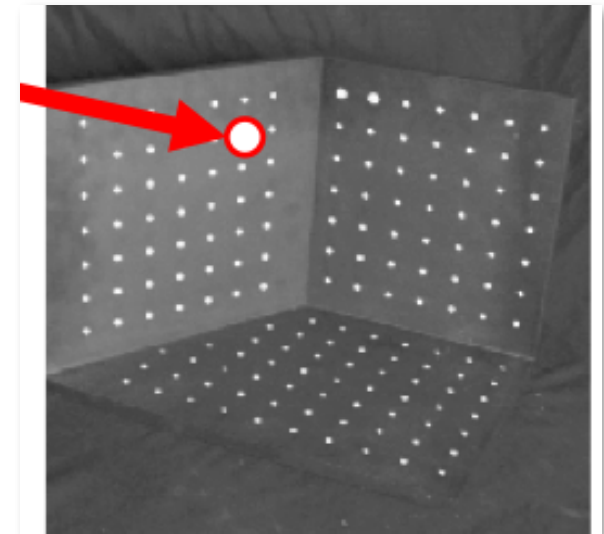
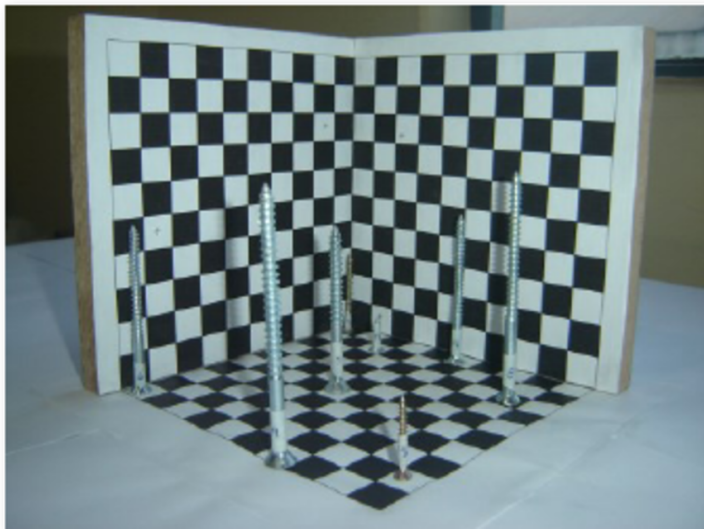
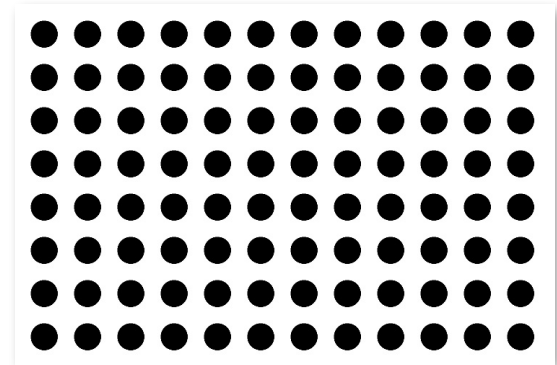
PATTERN:

- *spheres*
- *circles*
- *chessboard*

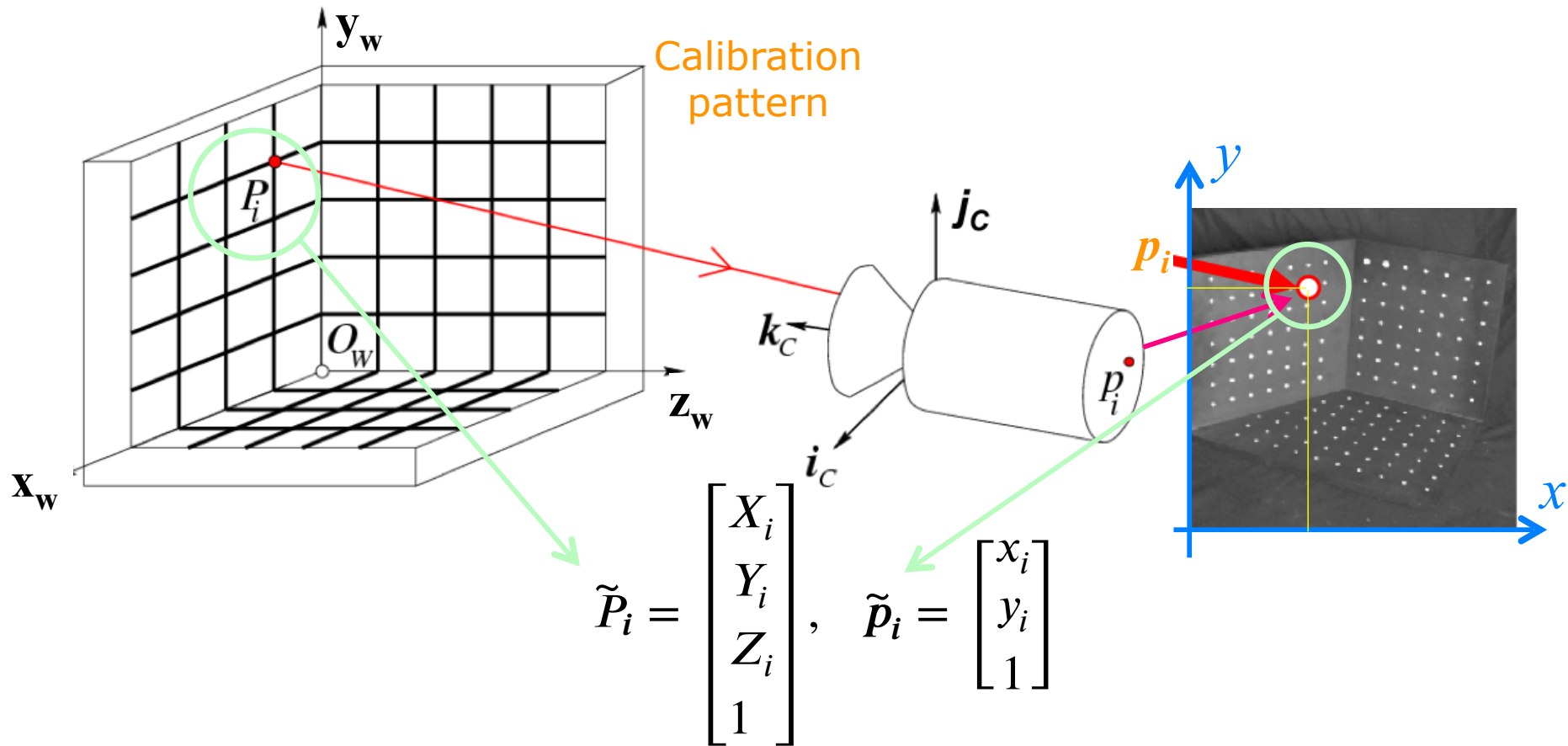
Fiducial Point:

- *sphere centre*
- *circle centre*
- *square vertices*

- Fiducial points over a non-degenerate 3D-space
 - If I use **planar patterns**, I need **at least 2 images**, on different planes



Camera calibration – Problem Definition



Problem definition:

- Given a set of N fiducial points P_i , of known 3D world position
- and given the corresponding image-coordinates p_i
- determine the camera model M (function of ξ) such that:

$$\tilde{p}_i = M \cdot \tilde{P}_i, \quad i = 1..N$$

Camera calibration - Linear approach

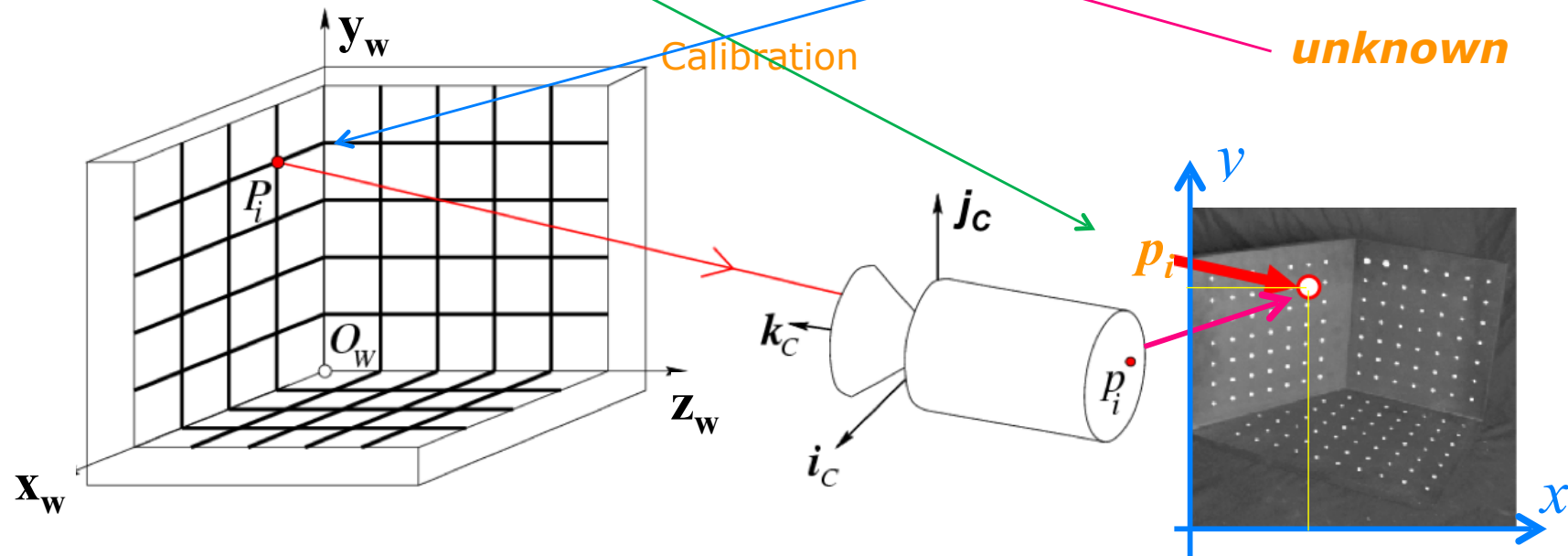
Determine the matrix \mathbf{M} [3 x 4] of the linear model $\tilde{p}_i = \mathbf{M} \cdot \tilde{\mathbf{P}}_i$ given:

- the coords-World : $P_i, i = 1..N$
- the coords-image: $p_i, i = 1..N$

For each $i=1..N$:

N equations:

$$\tilde{p}_i = \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} = \mathbf{M} \cdot \tilde{\mathbf{P}}_i = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \mathbf{m}_3 \end{bmatrix} \cdot \begin{bmatrix} X_i \\ Y_i \\ Z_i \\ 1 \end{bmatrix}, \quad i = 1 \dots N$$



Camera calibration - Linear approach

For a pair $\langle P_i, p_i \rangle$:

$$\tilde{\mathbf{p}}_i = \begin{bmatrix} \tilde{x}_i \\ \tilde{y}_i \\ \tilde{z}_i \end{bmatrix} = \mathbf{M} \cdot \tilde{\mathbf{P}}_w = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \mathbf{m}_3 \end{bmatrix} \cdot \tilde{\mathbf{P}}_i = \begin{bmatrix} \mathbf{m}_1 \cdot \tilde{\mathbf{P}}_i \\ \mathbf{m}_2 \cdot \tilde{\mathbf{P}}_i \\ \mathbf{m}_3 \cdot \tilde{\mathbf{P}}_i \end{bmatrix} \quad (\text{eq. 1})$$

- Remember Euclidean vs Homogeneous coords:

$$\mathbf{p}_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} = \begin{bmatrix} \frac{\tilde{x}_i}{\tilde{z}_i} \\ \frac{\tilde{y}_i}{\tilde{z}_i} \end{bmatrix} \quad (\text{eq. 2})$$

- Considering \mathbf{p}_i (in Euclidean coords) and combining (eq. 1) and (eq. 2) we obtain:

$$\mathbf{p}_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} = \begin{bmatrix} \frac{\tilde{x}_i}{\tilde{z}_i} \\ \frac{\tilde{y}_i}{\tilde{z}_i} \end{bmatrix} = \begin{bmatrix} \frac{\mathbf{m}_1 \tilde{\mathbf{P}}_i}{\mathbf{m}_3 \tilde{\mathbf{P}}_i} \\ \frac{\mathbf{m}_2 \tilde{\mathbf{P}}_i}{\mathbf{m}_3 \tilde{\mathbf{P}}_i} \end{bmatrix} \Rightarrow \begin{cases} \mathbf{m}_1 \tilde{\mathbf{P}}_i - x_i \mathbf{m}_3 \tilde{\mathbf{P}}_i = 0 \\ \mathbf{m}_2 \tilde{\mathbf{P}}_i - y_i \mathbf{m}_3 \tilde{\mathbf{P}}_i = 0 \end{cases},$$

Camera calibration - linear approach

For a pair $\langle P_i, p_i \rangle$:

$$p_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} = \begin{bmatrix} \tilde{x}_i \\ \tilde{z}_i \\ \tilde{y}_i \\ \tilde{z}_i \end{bmatrix} = \begin{bmatrix} \frac{\mathbf{m}_1 \tilde{\mathbf{P}}_i}{\mathbf{m}_3 \tilde{\mathbf{P}}_i} \\ \frac{\mathbf{m}_2 \tilde{\mathbf{P}}_i}{\mathbf{m}_3 \tilde{\mathbf{P}}_i} \end{bmatrix} \Rightarrow \begin{cases} \mathbf{m}_1 \tilde{\mathbf{P}}_i - x_i \mathbf{m}_3 \tilde{\mathbf{P}}_i = 0 \\ \mathbf{m}_2 \tilde{\mathbf{P}}_i - y_i \mathbf{m}_3 \tilde{\mathbf{P}}_i = 0 \end{cases}, \quad i = 1..N$$

2 equations, 12 unknown m_{ij}

- In matricial form:

$$\begin{bmatrix} P_{1x} & P_{1y} & P_{1z} & 1 & 0 & 0 & 0 & 0 & -x_1 P_{1x} & -x_1 P_{1y} & -x_1 P_{1z} & -x_1 \\ 0 & 0 & 0 & 0 & P_{1x} & P_{1y} & P_{1z} & 1 & -y_1 P_{1x} & -y_1 P_{1y} & -y_1 P_{1z} & -y_1 \end{bmatrix} \begin{bmatrix} m_{11} \\ m_{12} \\ m_{13} \\ m_{14} \\ m_{21} \\ m_{22} \\ m_{23} \\ m_{24} \\ m_{31} \\ m_{32} \\ m_{33} \\ m_{34} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Camera calibration - Linear approach

For N pair $\langle P_i, p_i \rangle$ ($i = 1..N$) :

- write the **2N** eq as a linear system in the 12 unknowns m_{ij}

$$\begin{bmatrix} P_{1x} & P_{1y} & P_{1z} & 1 & 0 & 0 & 0 & 0 & -x_1 P_{1x} & -x_1 P_{1y} & -x_1 P_{1z} & -x_1 \\ 0 & 0 & 0 & 0 & P_{1x} & P_{1y} & P_{1z} & 1 & -y_1 P_{1x} & -y_1 P_{1y} & -y_1 P_{1z} & -y_1 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ P_{Nx} & P_{Ny} & P_{Nz} & 1 & 0 & 0 & 0 & 0 & -x_N P_{Nx} & -x_N P_{Ny} & -x_N P_{Nz} & -x_N \\ 0 & 0 & 0 & 0 & P_{Nx} & P_{Ny} & P_{Nz} & 1 & -y_N P_{Nx} & -y_N P_{Ny} & -y_N P_{Nz} & -y_N \end{bmatrix} \cdot \begin{bmatrix} m_{11} \\ m_{12} \\ m_{13} \\ m_{14} \\ m_{21} \\ m_{22} \\ m_{23} \\ m_{24} \\ m_{31} \\ m_{32} \\ m_{33} \\ m_{34} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{P}_{[2N \times 12]} \cdot \mathbf{m}_{[12 \times 1]} = \mathbf{0}_{[2N \times 1]}$$

$\rightarrow \mathbf{P} \cdot \mathbf{m} = \mathbf{0}$
Homogeneous linear system, in 11 unknowns (12, up to a scale factor) and 2N equations resolvable for $N = 6$ (at least 6 fiducial points)

Scale of Projection Matrix

- REMEMBER:
 - Projection Matrix M acts on homogeneous coords, i.e.:

$$\begin{bmatrix} u \\ v \\ w \end{bmatrix} \equiv k \cdot \begin{bmatrix} u \\ v \\ w \end{bmatrix} \quad (k \neq 0 \text{ is any constant})$$

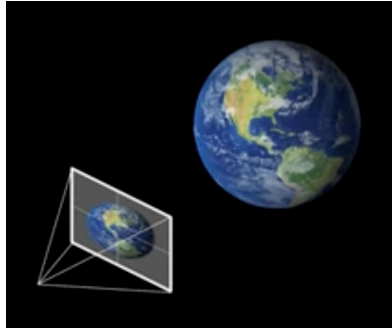
- That is:

$$\tilde{P} \cdot M = \tilde{P} \cdot (k \cdot M)$$

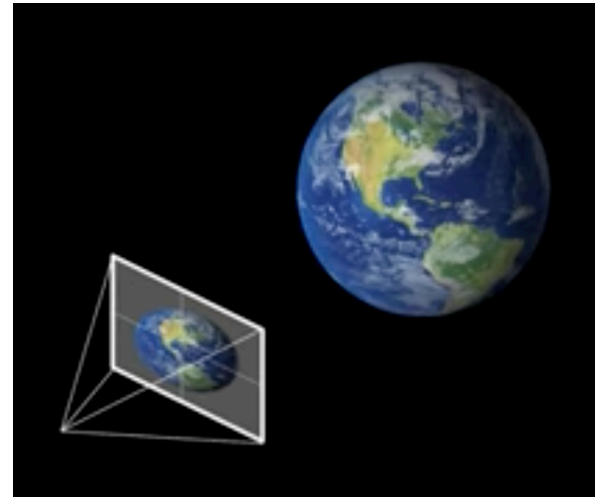
- so that: the projection matrices M and $(k \cdot M)$ produce the same homogeneous pixel coordinates

Projection matrix M is defined up to a scale factor

Scale Projection Matrix



Scale k_1



Scale k_2

Scaling projection matrix, implies simultaneously scaling the world and camera, which does not change the image

we can set **projection matrix** to some **arbitrary scale**

Least squares solution for m

- Option 1: set scale so that $m_{34} = 1$
- Option 2: set scale so that $\|m\|^2 = 1$
- We want: $\tilde{P} \cdot m = 0$, and $\|m\|^2 = 1$
- Formulated with the **constrained least squares problem**:

$$\min_m \|\tilde{P}m\|^2, \quad s.t. \quad \|m\|^2 = 1,$$

$$\min_m \|m^T \tilde{P}^T \tilde{P} m\|, \quad s.t. \quad \|m^T m\| = 1$$

- Let's define the **Loss function** $L(m, \lambda)$:

$$L(m, \lambda) = m^T \tilde{P}^T \tilde{P} m - \lambda(m^T m - 1)$$

- We want to **minimize** L wrt m

Constrained Least squares solution

- Let's take the derivatives of $L(m, \lambda)$ wrt m and set it to 0:

$$2\tilde{P}^T \tilde{P}m - 2\lambda m = 0$$

- equivalent to solve the **eigenvalue problem**:

$$\tilde{P}^T \tilde{P}m = \lambda m$$

- **Eigenvector** m corresponding to the smallest eigenvalue λ of the matrix $\tilde{P}^T \tilde{P}$ minimizes the loss function $L(m, \lambda)$
- or equivalently, m is the **singular vector** corresponding to the minimum singular value (not null) using the Singular Value Decomposition (**SVD**) of \tilde{P}



Camera calibration - Linear approach

Given the vector m we have to:

- rearrange it to form the projection matrix M (4×4)
- it remains to determine the intrinsic and extrinsic matrices:

$$M = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} = \underbrace{\begin{bmatrix} f & 0 & c_x & 0 \\ 0 & f & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}}_{M_{int}} \cdot \underbrace{\begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{M_{ext}}$$

- it holds that: A

$$\begin{bmatrix} m_{11} & m_{12} & m_{13} \\ m_{21} & m_{22} & m_{23} \\ m_{31} & m_{32} & m_{33} \end{bmatrix} = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} = KR$$

- Given K : upper triangular matrix, R is an orthonormal matrix, it is possible to decouple K and R from their product using the QR factorization method from linear algebra

Camera calibration - Linear approach

- It remains to determine the translation vector \mathbf{t} of the extrinsic matrix
- Given:

$$M = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} = \underbrace{\begin{bmatrix} f & 0 & c_x & 0 \\ 0 & f & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}}_{M_{int}} \cdot \underbrace{\begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{M_{ext}}$$

- it holds that:

$$\begin{bmatrix} m_{14} \\ m_{24} \\ m_{34} \end{bmatrix} = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix} = K\mathbf{t}$$

- and inverting we have:

$$\mathbf{t} = K^{-1} \cdot \begin{bmatrix} m_{14} \\ m_{24} \\ m_{34} \end{bmatrix}$$

Camera Calibration: Zhang algorithm (2000)

"A_flexible_new_technique_for_camera_calibration" in ***IEEE Transaction on Pattern Analysis and Machine Intelligence***, vol. 22, no. 11, pp. 1330-1334, 2000.

- Planar pattern in at least 2 views (chessboard)



- **Hypothesis:**
 - *The pattern dimensions are known*
 - *Trick: the scene reference system is joint to the chessboard (different extrinsic parameters for each photo, while shared intrinsic parameters)*

Camera calibration - Non Linear approach

- To be used in presence of radial distortion

$$\underbrace{\tilde{\mathbf{p}}_i = \begin{bmatrix} \tilde{x}_i \\ \tilde{y}_i \\ \tilde{z}_i \end{bmatrix} = \mathbf{M}(\xi) \cdot \tilde{\mathbf{P}}_i}_{\text{linear}} \rightarrow \underbrace{\tilde{\mathbf{p}}_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} = f(\xi, \mathbf{P}_i)}_{\text{non linear}}$$

Algorithm:

1. Solve the linear model

→ first estimate of the linear parameters: ξ_0

2. Non linear optimization starting from ξ_0 to minimize the reprojection error E

→ Final model: $\hat{\xi}$

$$E = \sum_i \left\| \mathbf{p}_i^{MEAS} - \mathbf{p}_i^{EST} \right\|^2 = \sum_i \left\| \mathbf{p}_i^{MEAS} - f(\xi, \mathbf{P}_i) \right\|^2 \rightarrow \hat{\xi} \quad t.c. \quad \hat{\xi} = \underset{\xi}{\operatorname{argmin}}(E)$$

Camera calibration - Non Linear approach

More specifically:

$$E = \sum_i \left\| \mathbf{p}_i^{MEAS} - \mathbf{p}_i^{EST} \right\|^2 = \sum_i \left\| \mathbf{p}_i^{MEAS} - f(\xi, \mathbf{P}_i) \right\|^2 \rightarrow \hat{\xi} \quad t.c. \quad \hat{\xi} = \underset{\xi}{\operatorname{argmin}}(E)$$

Given the initial model: $\xi_0 = \{\mathbf{R}, \mathbf{t}, f, x_C, y_C, k_D\}$ and N fiducial points P_i :
 $\xi = \xi_0$

$$p_i^{EST} = f(\xi, P_i) :$$

$$P_i \mapsto P_{cam,i} = \mathbf{R} \cdot P_i + \mathbf{t}, \quad i = 1..N$$

$$P_{cam,i} \mapsto p_{im,i} = \begin{bmatrix} x_C \\ y_C \end{bmatrix} + \frac{f}{z_{cam,i}} \begin{bmatrix} x_{cam,i} \\ y_{cam,i} \end{bmatrix}, \quad i = 1..N$$

$$p_{im,i} \mapsto p_i^{EST} = p_{im,i} (1 + k_{D1}r^2 + k_{D2}r^4), \quad i = 1..N$$

$$\text{error: } E(\xi) = \sum_{i=1}^N \left\| p_i^{MEAS} - p_i^{EST} \right\|^2 \rightarrow \text{new estimate } \xi$$

Lab time

- syntLinearCalibration
- RealCalibration