

Lecture 8: Dense Reconstruction

→ Why do we need Dense Reconstruction?

→ Accurate 3D models → cultural heritage

→ Reconstruction of houses, buildings, famous sites.

→ Stereo Reconstruction →

→ Input → a left and right pair of images

→ Output → a disparity or depth map.

$$\text{disparity} \propto \frac{1}{\text{depth}}$$

* In stereo vision, the two images, captured by the two cameras separated by a distance, can be used to get the 3D location (x, y, z) of the points of images in real world i.e. the depth.

disparity → is the difference in image location of the same 3D point when projected under perspective to two different cameras.

* Any point in the scene that is visible in both cameras will be projected to a pair of image points in the two images, called a conjugate pair. The displacement b/w the positions of the two points is called the "disparity".

depth → (the actual location of 3D point) can be calculated by using the disparity of the corresponding point

$$\text{depth} = \frac{\text{baseline} * \text{focal length}}{\text{disparity}}$$

distance b/w 2 cameras

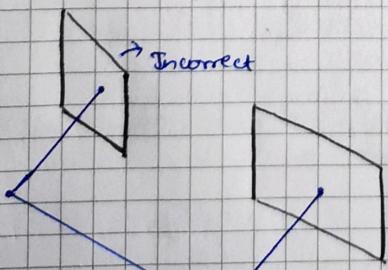
→ How can we recover the depth?

We can recover the depth by (matching + triangulation)

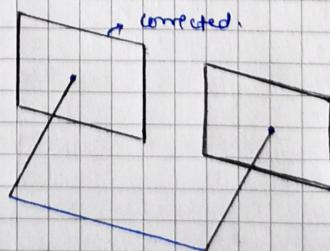
→ Stereo reconstruction → Image Rectification.

→ for stereo cameras, we assume that the cameras are parallel.

When we make a stereo camera, there may be many physical inaccuracies during manufacturing leading to the camera not being parallel.



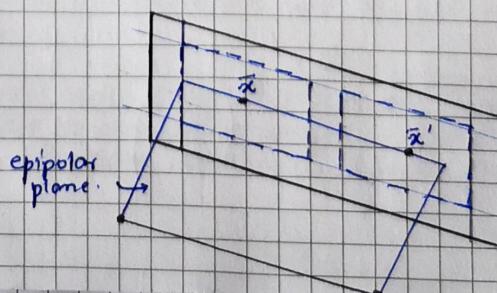
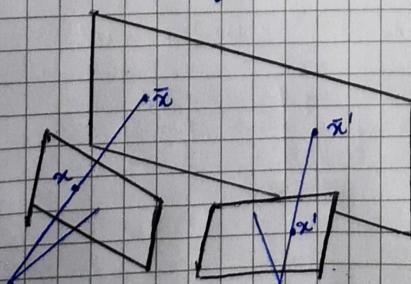
⇒



} Rectifying 1 image.

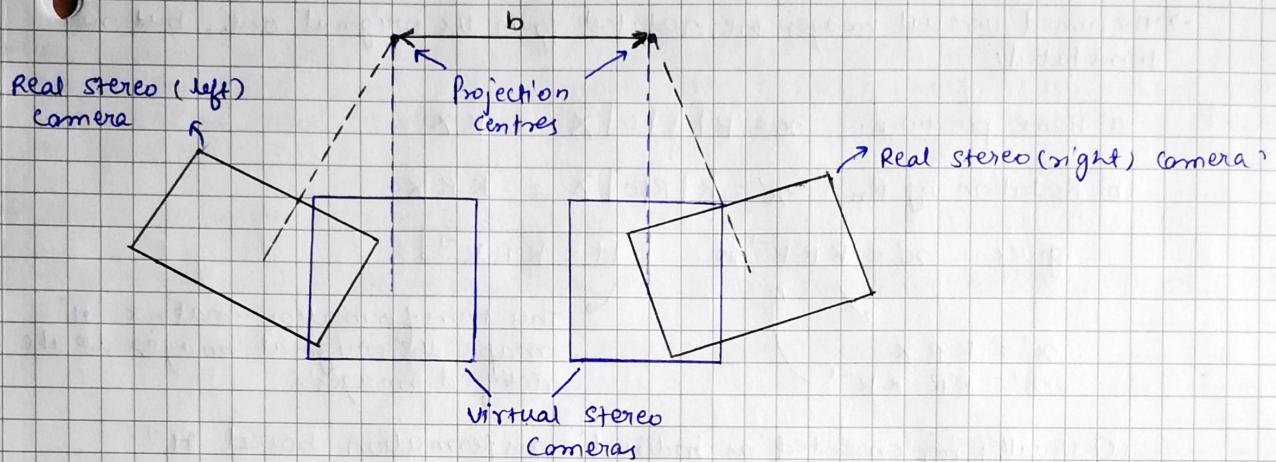
image mapping is a 2D homography.

$$H = KRK^{-1}$$



→ How do we do this?

- ① We have two cameras which are not perfectly aligned. We have to rotate the cameras in such a way that they are aligned with the baseline and also to the y-axis.



- ② Define rotations of the two virtual views that the x-axes of the camera coordinate system are parallel to the transition of the digitized original camera positions.

→ These rotations are distributed equally to the right and the left camera to reduce the resulting distortions.

→ Additionally, the camera matrix representing the intrinsic parameters should be equal for both the virtual views.

- ③ Assuming calibrated stereo cameras

- ④ The intrinsic and extrinsic parameters are given.

⑤ Determination of Image Transformation

- ① The new camera matrix K_n should be equal for both virtual views which can be computed by calculating the mean of the old views.

$$K_n = \frac{K_0^{\text{left}} + K_0^{\text{right}}}{2}$$

(x, y, z)

- ② New rotation matrix R_n which represents the transformation of virtual camera system wrt world coordinates system.

- Virtual x-axes is set parallel to the image baseline $\rightarrow \underline{x_1}$
- Virtual y-axes is set orthogonal to the virtual x-axes and orthogonal to old z-axes of left camera. or
- orthogonal to the mean of old z-axes of two cameras $\rightarrow \underline{x_2}$
- Virtual z-axes, orthogonal to the virtual x-axes (baseline) and virtual y-axes $\rightarrow \underline{x_3}$

translation of world
camera on left.

$$\underline{x}_1 = \frac{\underline{x}^{\text{left}} - \underline{x}^{\text{right}}}{\|\underline{x}^{\text{left}} - \underline{x}^{\text{right}}\|}$$

$$\underline{x}_2 = 0.5(\underline{x}_{30L} + \underline{x}_{30R}) \times \underline{x}_1$$

$$\underline{x}_3 = \underline{x}_1 \times \underline{x}_2$$

$$R_n = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \underline{x}_3^T \end{bmatrix}$$

→ Computation of Homography → finding the mapping of old image to the new image.

① Link transformations w/ Homography "H"

→ The aimed virtual images are rotated from the original one, but not translated:

$$(a) \text{Home position: } x = K[I|0]x = Kx$$

$$(b) \text{Rotation by } R_n: x' = K[R|0]x = KRx$$

$$\text{gives } x' = KRK^{-1}x$$

$$x' = KRx$$

$$x' = KRxK^{-1}$$

$$H = KRK^{-1}$$

→ This transformation matrix "H" maps the original images to the rectified images.

(c) Usually implemented as indirect transformation based H^{-1}

(d) Right image is calculated in a similar way

$$H^{-1}x' = x$$

pixels in the new image projected back into the original image.

We don't go from old image to new image, but rather from new image (virtual) to the old image (real).
This is done to take care of interpolation errors.

→ Direct and indirect geometric Transformation

Direct Transformation

① For each pixel of the input image, the corresponding coordinate of the output image are computed.

② These coordinates will usually have non-integer values (decimal values). Also, some pixels might have multiple or non-assignments.

③ This requires additional computational effort during assignments of grey values.

④ Not recommended.

→ step 2 → Triangulation with Parallel cameras

→ Triangulation reduces to a simple geometric problem

$$K = K' = \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R = I$$

$$t = \begin{pmatrix} tx \\ ty \\ tz \end{pmatrix}$$

Indirect Transformation

① The inverse transformation is used to calculate (for each pixel) of the output image, the corresponding (non-integer) position of input image.

② The grey values of the surrounding pixels are then used for simple interpolation (bi-linear interpolation)

→ Any other interpolation method is also possible.

→ Then $y = y'$ and depth z can be computed from disparity:

$$z = \frac{fx}{d} \quad d(\text{disparity}) = x' - x$$

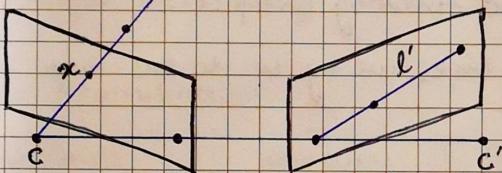
→ refer to Lecture 7
page 3

→ As opposed to the matching for structure for motion where we use feature based mapping with (Harris) lines, edges, SIFT keypoints), here we are more interested in dense reconstruction therefore, we try to match as many pixels as possible.

⇒ Stereo Reconstruction → Dense matching.

→ for every point in the image, there are many possible matches on the other image.

→ there is high-ambiguity since many points look similar.



$$l' = Fx$$

$$F = K'^{-T} [t] \times R K^{-1}$$

$$[x'^T F x = 0]$$

→ We can use the Camera Geometry i.e. finding the correspondences becomes a 1D search problem. But this solution also yields multiple solutions for a point.

→ Therefore we use the pixel-neighbourhood information – correlation based approaches (also called window-based approaches)

→ To do this window based matching, we have 3 approaches:

(a) Normalized cross-correlation (NCC) → Calculating the correlation of 2 signals and also normalizing it.

$$\frac{\sum_{(i,j) \in W} I_1(i,j) \cdot I_2(x+i, y+j)}{\sqrt{\sum_{(i,j) \in W} I_1(i,j)^2 \cdot \sum_{(i,j) \in W} I_2(x+i, y+j)^2}}$$

(b) Sum of squared differences (SSD)

$$\sum_{(i,j) \in W} (I_1(i,j) - I_2(x+i, y+j))^2$$

window.

(c) Sum of absolute differences (SAD)

$$\sum_{(i,j) \in W} |I_1(i,j) - I_2(x+i, y+j)|$$

(d) Census transform.

→ Based on comparison of radiance values within the window
→ The effect of outliers is low.

* Now we have the steps for matching.

① Search along the epipolar line for the best match.

② Best match = (a) maximum NCC
(b) minimum SSD or SAD.

→ Enhanced matching

→ Estimating the disparity solely based on matching leads to errors.

→ There are 2 approaches of optimizing correspondences jointly.

→ Scanline at a time with dynamic programming techniques
→ full 2D grid with graph cuts.

Scanline - stereo → Coherently match the pixels on the entire scan-line
* Scanline stereo generates streaking artifact.

→ Instead of doing this line by line, we try to do this with 2D.

→ Energy minimization →

(a) In order to find more reliable correspondences, an energy function can be minimized.

(b) It can be adapted to the situation and becomes arbitrarily complicated.

$$\min \underbrace{E_{\text{data}}}_{\text{data term}} + \lambda E_{\text{reg}} \rightarrow \begin{array}{l} \text{regularization} \\ \text{term} \end{array} \quad (\text{smoothness using interpolation})$$

(c) usually consists of a data term (consistency of matching information)
→ NCC, SSD, SAD

(d) A regularization term (smoothness (total variation regularization))

(e) Also different patch sizes can be included.

(f) can be minimized using graph cuts

→ Dense matching → challenges

Scene elements do not always look the same in the images.

(A) Camera-related problems.

(a) Image Noise (b) Lens distortion (c) color/chromatic aberration

(B) Viewpoint related problems

(a) Perspective distortion (b) Occlusions (c) Specular reflection

(C) Scene-related problems

(a) Illumination ~~problems~~ changes (b) moving objects

(d) matching ambiguity

(a) low texture regions (b) Repetitive texture patterns.

→ Assumptions in Dense matching.

- ① Appearance → the projections of a scene (3D) patch should have similar appearances in all images.
- ② Uniqueness → A point in an image will only have one match in another image
- ③ Ordering → the order of appearance is not altered by the camera displacement.
- ④ Smoothness → Depth varies smoothly (objects have smooth surfaces)

→ choice of camera setup during Dense matching.

short baseline

- ① matching is relatively easy
- ② Higher depth uncertainty

wide baseline

- ① matching is hard
- ② lower depth uncertainty ✓

⇒ Multi View Reconstruction.

- ① more constraints
- ② less occlusions
- ③ more robust
- ④ Appearance constraint can be relaxed

However.

- ① more complex
- ② strong perspective distortions
- ③ propagation of calibration uncertainties



Dense matching along many images is still hard.

→ Triangulation in multiple images.

→ Assuming we had found the correct matches along "n" images, we triangulate them exactly as before by stacking the equations for all camera views. i.e.

$$\begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_n \end{bmatrix} x = 0$$

→ Multi View Reconstruction - Algorithms.

* We have many different algorithms but all of them have some assumptions

- ① Lambertian surfaces (no specular reflections)
- ② static scene
- ③ camera parameters are available (projection matrices are known)
- ④ Underlying photometric consistency (NCC, etc)

→ Scene Representation → Voxel Grid (binary representation)

- Divides the scene into cubes (voxels)
- Need to sample the 3D space
- Voxels are marked as occupied / empty according to a photometric consistency measure (NCC, SSD, etc)
- Accuracy depends on size of voxels
- No assumptions on the scene required.
- memory allocation issues.

(b) Depth map.

- usually 1 depth map per view (by matching and triangulating features in small sets of images)
- NO sampling of 3D space
- simple
- Suitable for multi-core
- Depth maps have to be merged.

(c) Point cloud.

- Local planar approximation of the surface at each point
- Usually a normal vector is associated to each point.
- Sparse to dense reconstruction.
- Suitable for multi-core
- may be noisy.

(d) mesh.

- models the surface as a connected set of planar facets
- Usually triangular meshes
- Triangles are locally good approximations of the surface
- Good for optimization
- may become difficult to handle for complex surfaces.