

Lecture 8, Sep 26, 2025

Camera Pose Estimation (Perspective-n-Point)

- The *pose estimation problem* is the estimation of the pose of an object (or camera) from 3D-2D correspondences
 - Typically we know where the points lie on the object in 3D, and we also know where those points appear in our image
 - This is known as the *perspective-n-point* or PnP problem
 - * For 3 points, this is known as P3P
 - * 3 points is the minimum number of points we need to get a solution, but often we want to use more points to remove ambiguity and reject noise
 - This process can be used to estimate the camera pose for extrinsic calibration; intrinsic parameters can be estimated at the same time
 - * This is how camera calibration with a checkerboard works
- There are linear and nonlinear algorithms for this
 - The linear case locally linearizes the problem and may not be very accurate if the initial guess is bad
 - The nonlinear algorithm uses an initial guess and iterates to find the solution
 - We often start from the linear algorithm and then use nonlinear algorithm to refine
- Solving for $\mathbf{P} = \mathbf{K} [\mathbf{C} \quad \mathbf{t}]$ requires at least 6 correspondences (since there are 6 degrees of freedom), but if we have intrinsics already we only need 3
- The linear algorithm for solving PnP is the *direct linear transform* (DLT), which stacks a system of equations
 - We have $\mathbf{P} \in \mathbb{R}^{3 \times 4}$ (the combination of the extrinsic and intrinsic matrices) with 12 unknowns
 - For each correspondence we can construct 2 equations from it; with 6 correspondences we can solve the whole system
 - * $x_i = \frac{p_{00}X_i + p_{01}Y_i + p_{02}Z_i + p_{03}}{p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}}$
 - * $y_i = \frac{p_{10}X_i + p_{11}Y_i + p_{12}Z_i + p_{13}}{p_{20}X_i + p_{21}Y_i + p_{22}Z_i + p_{23}}$
 - * Note the division due to normalization
 - Because the intrinsics matrix \mathbf{K} is upper-triangular, we can do a QR factorization on \mathbf{P} to recover the separate intrinsic and extrinsic matrices
- However, DLT does not impose constraints on the structure of the resulting matrices, namely the structure of the rotation matrix, so after QR factorization we often end up with a result that does not fit into our camera models, forcing us to make an approximation; this is why DLT is not an exact solution

Nonlinear Least Squares

- Nonlinear least squares (Gauss-Newton optimization) is an optimization approach we can use to solve this system
 - For nonlinear least squares, we wish to optimize $E(\mathbf{x}) = \frac{1}{2} \mathbf{e}(\mathbf{x})^T \mathbf{e}(\mathbf{x})$ where $\mathbf{e}(\mathbf{x})$ is some nonlinear function
 - * Note $\mathbf{e}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) - \mathbf{y}$, i.e. the prediction minus the observation; the order is important, otherwise we end up maximizing the error instead!
 - * In the linear case we can substitute $\mathbf{e}(\mathbf{x}) = \mathbf{A}\mathbf{x} - \mathbf{b}$ and expand the error function, then take a derivative to obtain the normal equation
 - For the nonlinear case we linearize around an initial guess (the operating point) \mathbf{x}_{op}
 - $\mathbf{e}(\mathbf{x}) \approx \mathbf{e}(\mathbf{x}_{op}) + \mathbf{J}_e \delta \mathbf{x}$ where $\mathbf{J}_e = \left. \frac{\partial \mathbf{e}}{\partial \mathbf{x}} \right|_{\mathbf{x}_{op}}$ is the Jacobian and $\delta \mathbf{x}$ is a small deviation
 - * Now substitute this back into the error function and notice we get an expression very similar to the linear form, so we can use the same techniques to solve this

- We can solve the linearized system, and add the $\delta \mathbf{x}$ to our initial operating point \mathbf{x}_{op} to get the next linearization point
- Do this until our Jacobian becomes sufficiently small which means we have converged
- We are essentially approximating the cost function as a quadratic at each step and optimizing the quadratic for a local solution
- Important notes for nonlinear least squares:
 - Choosing good initial guesses is important, otherwise the optimization process can get trapped in a local minimum
 - The states we are solving for must exist in a vector space (i.e. we cannot apply constraints, since then the vector space is no longer closed)
- For multiple errors, we sum over all the errors, so in the normal equation we sum over $\mathbf{J}^T \mathbf{J}$ and $\mathbf{J}^T \mathbf{e}$ and multiply in the end
 - $E(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^N \mathbf{e}_i(\mathbf{x})^T \mathbf{e}_i(\mathbf{x})$
 - $\delta \mathbf{x}^* = - \left(\sum_{i=1}^N \mathbf{J}_{e_i}^T \mathbf{J}_{e_i} \right)^{-1} \left(\sum_{i=1}^N \mathbf{J}_{e_i}^T \mathbf{e}_i(\mathbf{x}_{op}) \right)$
 - Note we need to reevaluate the Jacobian for each measurement i , since the linearization point is all different!
- Often we have associated uncertainties for each error, so we can do a weighted version of least squares, so $E(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^N \mathbf{e}_i(\mathbf{x})^T \mathbf{W}_i \mathbf{e}_i(\mathbf{x})$
 - Now we have $\delta \mathbf{x}^* = - \left(\sum_{i=1}^N \mathbf{J}_{e_i}^T \mathbf{W}_i \mathbf{J}_{e_i} \right)^{-1} \left(\sum_{i=1}^N \mathbf{J}_{e_i}^T \mathbf{W}_i \mathbf{e}_i(\mathbf{x}_{op}) \right)$
 - \mathbf{W}_i are symmetric matrices, one describing the weight for each measurement
 - Often we want scalar weights so \mathbf{W}_i for each measurement i is just a multiple of the identity
 - For vector-valued observations, we can use the inverse of the measurement covariance matrix Σ , known as the *information matrix*
- Note nonlinear least squares is equivalent to a maximum likelihood estimate if we assume our data has additive noise drawn from IID multivariate zero-mean Gaussians; weights are the information matrices of each noise Gaussian in this case
- To use NLS for pose estimation, we optimize the reprojection error $\mathbf{e}_i = \mathbf{x}_i - f(\mathbf{p}_i; \mathbf{C}, \mathbf{t}, \mathbf{K})$
 - f is our projection function which takes \mathbf{p}_i , converts it to the camera frame using the extrinsic calibration, then to pixel space using the intrinsics and normalizes it
 - But, rotation matrices have constraints, so we can't use Gauss-Newton as-is; we need to either express \mathbf{C} in a way that ensures it is a rotation matrix, or modify our update so that \mathbf{C} remains orthogonal
- This leads to the *Wahba problem*, which involves identifying a rotation matrix \mathbf{C} between frames given corresponding unit vector measurements $\mathbf{u}_i, \mathbf{v}_i$ in two frames
 - The cost function is $E(\mathbf{C}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{C} \mathbf{u}_i - \mathbf{v}_i)^T (\mathbf{C} \mathbf{u}_i - \mathbf{v}_i) = \frac{1}{2} \sum_{i=1}^N \mathbf{e}_i(\mathbf{C})^T \mathbf{e}_i(\mathbf{C})$
 - Approach 1: Euler angles
 - * Optimize over 3 Euler angles instead
 - * However, this runs the risk of Gimbal lock – if our solution ends up near points with Gimbal lock, we run into numerical sensitivity issues or our solution may collapse entirely
 - This can work if we have good initial guesses
 - * Let $\mathbf{C}(\boldsymbol{\theta}) = \mathbf{C}_3(\theta_3) \mathbf{C}_2(\theta_2) \mathbf{C}_1(\theta_1)$, then $\mathbf{e}_i(\boldsymbol{\theta}_{op} + \delta \boldsymbol{\theta}) = \mathbf{e}_i(\boldsymbol{\theta}_{op}) + \mathbf{J}_{e_i} \delta \boldsymbol{\theta}$
 - * The Jacobian $\mathbf{J}_{e_i} = \frac{\partial \mathbf{C} \mathbf{u}_i}{\partial \boldsymbol{\theta}}$ can be computed in each column as follows:
 - $\frac{\partial \mathbf{C} \mathbf{u}_i}{\partial \theta_3} = [\mathbf{1}_3]_{\times} \mathbf{C}_3(\theta_3) \mathbf{C}_2(\theta_2) \mathbf{C}_1(\theta_1) \mathbf{u}_i$

- $\frac{\partial \mathbf{C} \mathbf{u}_i}{\partial \theta_2} = \mathbf{C}_3(\theta_3) [\mathbf{1}_2]_{\times} \mathbf{C}_2(\theta_2) \mathbf{C}_1(\theta_1) \mathbf{u}_i$
 - $\frac{\partial \mathbf{C} \mathbf{u}_i}{\partial \theta_1} = \mathbf{C}_3(\theta_3) \mathbf{C}_2(\theta_2) [\mathbf{1}_1]_{\times} \mathbf{C}_1(\theta_1) \mathbf{u}_i$
 - Note $\mathbf{1}_i$ is a zero vector with 1 in the i th spot
- Approach 2: use axis-angle
- * To avoid gimbal lock we keep the operating point \mathbf{C}_{op} in matrix form, and consider a perturbation $\mathbf{C}(\delta\phi)$, so the update becomes $\mathbf{C}_{op} \leftarrow \mathbf{C}(\delta\phi) \mathbf{C}_{op}$
 - * Recall the Rodrigues formula: $\mathbf{C}(\phi) = 1 + \sin \phi [\hat{\mathbf{n}}]_{\times} + (1 - \cos \phi) [\hat{\mathbf{n}}]_{\times}^2$
 - For small angles ϕ we approximate $\sin \phi = 0, \cos \phi = 1$
 - $\mathbf{C}(\delta\phi) \approx 1 + \delta\phi [\hat{\mathbf{n}}]_{\times} = 1 + [\delta\phi]_{\times}$
 - * Substitute the approximation for $\mathbf{C}(\delta\phi)$:
 - $\mathbf{e}_i(\delta\phi) = \mathbf{C} \mathbf{u}_i - \mathbf{v}_i$

$$= \mathbf{C}(\delta\phi) \mathbf{C}_{op} \mathbf{u}_i - \mathbf{v}_i$$

$$= (1 + [\delta\phi]_{\times}) \mathbf{C}_{op} \mathbf{u}_i - \mathbf{v}_i$$

$$= \mathbf{C}_{op} \mathbf{u}_i - \mathbf{v}_i + [\delta\phi]_{\times} \mathbf{C}_{op} \mathbf{u}_i$$

$$= \mathbf{e}_i(\mathbf{C}_{op}) - [\mathbf{C}_{op} \mathbf{u}_i]_{\times} \delta\phi$$

$$= \mathbf{e}_i(\mathbf{C}_{op}) + \mathbf{J}_{e_i} \delta\phi$$
 - * Therefore the Jacobian is $\mathbf{J}_{e_i} = -[\mathbf{C}_{op} \mathbf{u}_i]_{\times}$, and with this we can use the normal equation to compute the update for $\delta\phi^*$ and update the operating point