

非线性优化

一、状态估计问题

1.1 最大后验与最大似然

经典SLAM模型由一个状态方程和一个运动方程构成，如下式所示：

$$\begin{cases} x_k = f(x_{k-1}, u_k) + w_k \\ z_{k,j} = h(y_j, x_k) + v_{k,j} \end{cases}$$

由于运动方程在视觉SLAM中没有特殊性，目前不讨论，将观测方程具体化，由针孔相机成像原理得到：

$$sz_{k,j} = KT y_i$$

其中， s 为像素点的距离（真实点到相机坐标系的距离 Z ）， K 为相机内参， T 为相机位姿矩阵， y_i 为相机观测到的路标点在世界坐标系下的坐标， z 为在 x_k 处观测 y 时对应的图像像素位置。

考虑噪声影响，假设两个噪声项 w_k 、 $v_{k,j}$ 满足零均值的高斯(正态)分布。在噪声的影响下，我们希望通过带噪声的数据 z 和 u ，推断位姿 x 和地图 y （以及它们的概率分布），这构成了一个状态估计的问题。

状态变量：我们把所有待估计的变量放在一个“状态变量”中：

$$x = \{x_1, \dots, x_N, y_1, \dots, y_M\}$$

我们要解决的问题是：已知输入数据（传感器读数） u 和观测数据 z （可以认为是像素位置）的条件下，计算状态 x 的**条件概率分布**：

$$P(x|z, u)$$

我们考虑没有测量运动的传感器时，只有图像，那么就计算 $P(x|z)$ ，利用贝叶斯公式，有

$$P(x|z) = \frac{P(z|x)P(x)}{P(z)}$$

直接求后验概率 $P(x|z)$ 较为困难，但是求一个状态最优估计，使得在该状态下，后验概率最大化（Maximize a Posterior, MAP），则是可行的。

$$x_{MAP}^* = \operatorname{argmax} P(x|z) = \operatorname{argmax} P(z|x)P(x)$$

对公式中 arg 的解释： $\operatorname{arg max}$ 是使后面式子达到最大值时变量的取值。例如：函数 $f(x, y)$ ， $\operatorname{arg max}$ 是指当 $f(x, y)$ 取得最大值时，变量 x, y 的取值

在不知道先验概率 $P(x)$ 的情况下，可以求解 x 的**最大似然估计**。即在什么样的状态下，最可能产生现在观测到的数据。

1.2 最小二乘法的引出

根据公式含义：

$$z_{k,j} = h(y_j, x_k) + v_{k,j}$$

上式表示，机器人在 x_k 位置观测到路标 y_j ，产生观测数据（像素） z_{kj} ， v_{kj} 为噪声，服从高斯分布。因此，我们可以得到观测数据的条件概率：

$$P(z_{j,k}|x_k, y_j) = N(h(y_j, x_k), Q_{k,j})$$

N代表服从正态分布，可以看出，这还是一个高斯分布。为了计算使它最大化的 x_k, y_j （极大似然估计），通常使用最小化负对数的方法。

任意高维高斯分布： $x \sim N(\mu, \Sigma)$

$$P(x) = \frac{1}{\sqrt{(2\pi)^N \det(\Sigma)}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

其中， Σ 为相关系数， μ 为均值。

取负对数后：前面系数与变量 x 无关，常数不管，只需要后面的二次型最小即可。可以使用无约束最优化求解（Newton、拟牛顿法、共轭梯度法、最速下降法）。

1.3 总结

首先， $h(y, x)$ 为定值，记为 z_{0kj} ，但是由于噪声，观测值 z_{kj} 与 z_{0kj} 不一样。我们要求的是：在得到观测值 z 的情况下， x 为和值的概率最大（例： $z=(100,230)$ 时， $P(x=(1,2,3)^T) \approx 0.95$ 最大，那么认为 $x=(1,2,3)^T$ ，即求 $P(x|z)$ 。这个很不好求，因此使用贝叶斯公式转换，再极大似然估计，把 x 看作参数，观测位置为已知。现在变成求 x 为何值时，得到这个观测位置的概率最大，那么我们把这个 x 认为位置，即求 $\arg\max_x P(z|x)$ ， x 为自变量。之后，这个极大似然估计的过程被转换为最小二乘问题。

二、非线性最小二乘法

最小二乘问题泛指具有如下形式的问题：

$$\min f(x) = \frac{1}{2} \sum_{j=1}^m r_j^2(x)$$

其中， m 一般指实例个数， r_j 指残差，即目标值和预估值的差。

我们想得到

2.1 一阶梯度与二阶梯度法

求解增量最直观的方式是将目标函数在 x 附近进行泰勒展开：

$$\|f(x + \Delta x)\|_2^2 \approx \|f(x)\|_2^2 + J(x)\Delta x + \frac{1}{2}\Delta x^T H \Delta x$$

我们可以看出，只保留一阶导数时(J 雅可比矩阵)，增量方向显然与梯度方向相反：

$$\Delta x^* = -J^T(x)$$

而当保留二阶导数时，我们想得到增量 Δx 使得上式值最小，因此对 Δx 求导数并且使其为0，可以得到：

$$H \Delta x = -J^T$$

第一种方法称为最速下降法，第二种方法称为牛顿法

2.2 Gauss-Newton法(只能处理最小二乘问题)

整体思想：将 $f(x)$ 进行一阶泰勒展开：

$$f(x + \Delta x) \approx f(x) + J(x)\Delta x$$

这里 $J(x)$ 是 $f(x)$ 关于 x 的导数，我们用一阶泰勒展开式代替原来的函数，再用得到的函数进行最小二乘求解。我们对一阶展开的上式求最小值：求 Δx 使得 $\|f(x+\Delta x)\|_2^2$ 达到最小：

$$\Delta x^* = \operatorname{argmin}_{\Delta x} \frac{1}{2} \|f(x) + J(x)\Delta x\|^2$$

我们对 Δx 求导(不是对 x 求导), 先化简得到:

$$\begin{aligned} & \frac{1}{2} \|f(x) + J(x)\Delta x\|^2 \\ &= \frac{1}{2} (f(x) + J(x)\Delta x)^T (f(x) + J(x)\Delta x) \\ &= \frac{1}{2} (\|f(x)\|_2^2 + 2f(x)^T J(x)\Delta x + \Delta x^T J(x)^T J(x)\Delta x) \end{aligned}$$

对上式求导, 并令导数为0:

$$2J(x)^T f(x) + 2J(x)^T J(x)\Delta x = 0$$

我们可以看到, 上式提供了 Δx 的方程:

$$\begin{aligned} H(x)\Delta x &= g(x) \\ \text{其中: } H &= J(x)^T J(x) \\ g &= -J(x)^T f(x) \end{aligned}$$

我们简化了牛顿法中计算Hesse矩阵的步骤, 使用J代替Hesse矩阵。因此, 高斯牛顿法步骤如图:

1. 给定初始值 x_0 。
2. 对于第 k 次迭代, 求出当前的雅可比矩阵 $J(x_k)$ 和误差 $f(x_k)$ 。
3. 求解增量方程: $H\Delta x_k = g$ 。
4. 若 Δx_k 足够小, 则停止。否则, 令 $x_{k+1} = x_k + \Delta x_k$, 返回 2。

2.3 Levenberg-Marquadt法 (阻尼牛顿法)

根据泰勒公式, 只有当 Δx 足够小的时候, 近似等式才有较好的效果。因此, 我们应该给 Δx 添加一个**信赖区域(Trust Region)**, 不能让它太大而使得近似不准确。这种思想叫做: **信赖区域方法(Trust Region Method)**。在信赖区域中, 我们认为近似有效; 否则无效。

利用实际函数与近似函数的变化量的比值来确定信赖区域:

$$\rho = \frac{f(x + \Delta x) - f(x)}{J(x)\Delta x}$$

分子为实际函数的增量 Δf , 分子为一阶泰勒近似增量(可见2.2)。比值越接近于1, 说明近似效果越好; 比值太小, 则应减小信赖区域范围; 比值太大, 则可以放大信赖区域范围。

1. 给定初始值 \mathbf{x}_0 ，以及初始优化半径 μ 。
2. 对于第 k 次迭代，求解：

$$\min_{\Delta \mathbf{x}_k} \frac{1}{2} \|f(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k) \Delta \mathbf{x}_k\|^2, \quad s.t. \|D \Delta \mathbf{x}_k\|^2 \leq \mu, \quad (6.24)$$

这里 μ 是信赖区域的半径， D 将在后文说明。

3. 计算 ρ 。
4. 若 $\rho > \frac{3}{4}$ ，则 $\mu = 2\mu$ ；
5. 若 $\rho < \frac{1}{4}$ ，则 $\mu = 0.5\mu$ ；
6. 如果 ρ 大于某阈值，认为近似可行。令 $\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta \mathbf{x}_k$ 。
7. 判断算法是否收敛。如不收敛则返回 2，否则结束。

上面是一个约束最优化问题，我们使用**Lagrange乘子**将其转换为一个无约束优化问题：

$$\Delta x^* = \operatorname{argmin}_{\Delta x} \frac{1}{2} \|f(x) + J(x)\Delta x\|^2 + \frac{\lambda}{2} \|D\Delta x\|^2$$

类似高斯牛顿法展开，我们简单的得到：

$$(H + \lambda D^T D) \Delta x = g$$

可解方程。如果我们把D取为单位矩阵，那么有：

$$(H + \lambda I) \Delta x = g$$

当 λ 较大时，L-M法更接近于一阶梯度下降法；当 λ 较小时，L-M法更接近于高斯牛顿法。因此，LM法可以一定程度上避免线性方程组系数矩阵的奇异和病态问题。