

# STEREO

## 关于深度估计的学习报告



上海大学 王聪豪

## 深度估计的学习报告

深度估计是指估计场景中对象到观察者的距离。有多种方法可以完成任务。一些流行的方法，包括飞行时间（time-of-flight）设备，结构光相机（structured light cameras）和多视图几何（multiple-view geometry），其中多视图几何结构的方法成本最低。在多视图几何方法中，我们关心如何使用普通相机来估计场景深度。我将逐步解决 Project\_Stereo 中的各个问题，设计程序来使用相机估算深度。

有一些问题我可能会合并到一节的内容中去，一些重要的点或是没理解的部分会这样**强调**，以便日后查看。

本项目相关代码开源在 <https://github.com/BindyAtobe/stereo>，本人水平有限，如有任何错误或您对我有任何建议，非常希望您能联系我。

上海大学 王聪豪

邮箱：919818192@qq.com

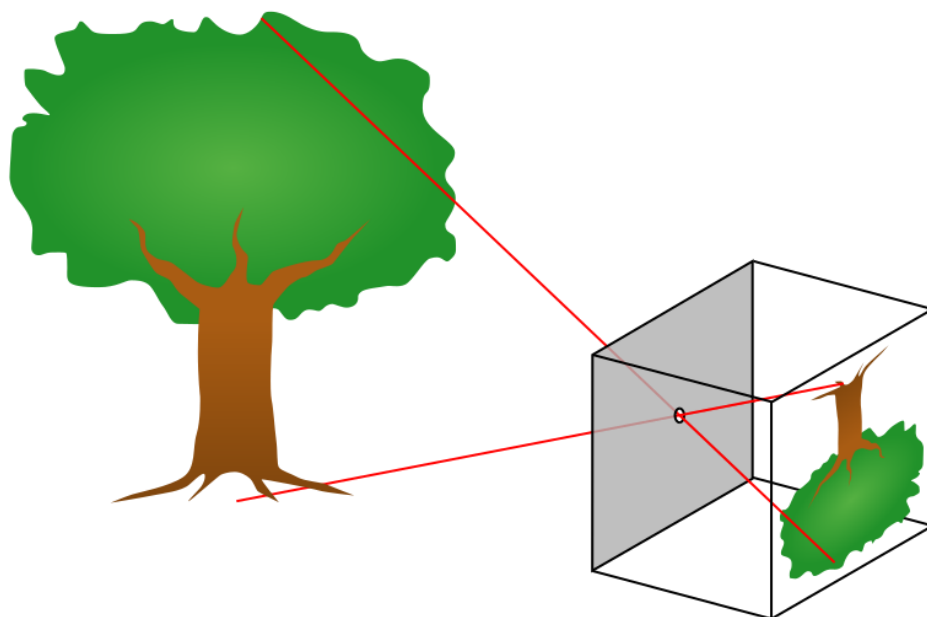
# 目录

<b>1. 相机基础 .....</b>	<b>1</b>
1.1. 相机矩阵，内参、外参 .....	1
1.2. 图像上一点对应的 3D 形状 .....	5
1.3. 图像畸变 .....	6
1.4. 相机标定，畸变校正 .....	8
1.5. 张正友标定 .....	9
<b>2. 双目基础 .....</b>	<b>13</b>
2.1. 更新中... ..	13
<b>3. 立体匹配 .....</b>	<b>14</b>

# 1. 相机基础

## 1.1. 相机矩阵，内参、外参

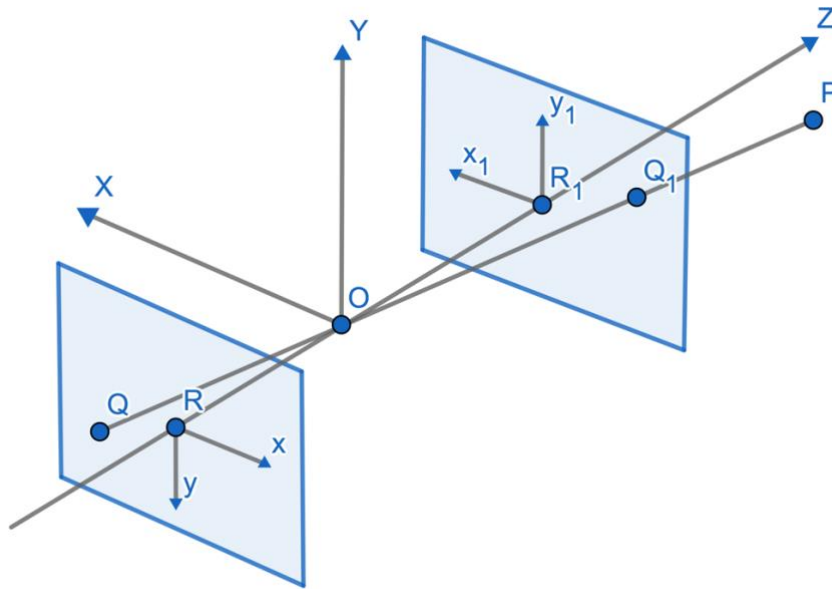
假定我们的相机符合针孔相机模型<sup>1</sup>，那么相机成像大致如下图所示，



我们可以发现物体和投影有倒立相似的关系。

<sup>1</sup> 维基百科-Pinhole camera model , [https://en.wikipedia.org/wiki/Pinhole\\_camera\\_model](https://en.wikipedia.org/wiki/Pinhole_camera_model)

加上 3D 和 2D 坐标系来描述相机成像，其中  $OR=OR_1=f$ （焦距），维基百科中的坐标系是左手的，且和第 2 题描述不同，所以我直接按第 2 题描述的方向来画坐标系，



由于倒像，我们往往会  $180^\circ$  旋转位于  $-f$  位置的图像 2D 坐标系，或者将图像移至  $f$  位置处。

成像的过程其实就是 3D 坐标到 2D 坐标的投影，设图中  $P(X_c, Y_c, Z_c)$ ， $Q(x, y)$ ，由相似易得，

$$\begin{pmatrix} x \\ y \end{pmatrix} = \frac{f}{Z_c} \begin{pmatrix} X_c \\ Y_c \end{pmatrix}$$

写成齐次形式（能用统一形式表示旋转和平移等好处，参考<sup>2</sup>），

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \frac{f}{Z_c} \begin{pmatrix} X_c \\ Y_c \\ \frac{Z_c}{f} \end{pmatrix}$$

<sup>2</sup> 简书-齐次坐标系入门级思考，<https://www.jianshu.com/p/80d0018ed24c>

用矩阵形式来写，

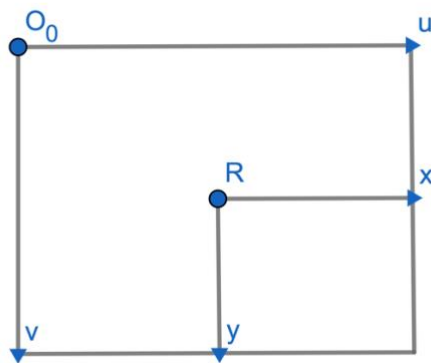
$$Z_c \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{pmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{pmatrix} \quad (1)$$

上述过程只是完成了点 P 从相机坐标到图像坐标的映射。实际上我们在观察点 P 的时候基于一个世界坐标系，它不一定会将原点设在针孔 O，方向也不一定与相机坐标系相同。观察图像时，也往往基于像素坐标系，而不是图中的图像坐标系。建立从世界坐标到像素坐标映射关系的矩阵叫相机矩阵。几个坐标系可以参考<sup>3</sup>，完整的坐标转换如下所示，

世界坐标→相机坐标→图像坐标→像素坐标

所以我们还需要知道如何从世界坐标映射到相机坐标，以及如何从图像坐标映射到像素坐标。

先讲图像坐标映射到像素坐标，下图中  $R(u_0, v_0)$ ，



$$\text{有} \quad \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{1}{dx} & \gamma & u_0 \\ 0 & \frac{1}{dy} & v_0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \quad (2)$$

<sup>3</sup> 百度百科-像素坐标，<https://baike.baidu.com/item/像素坐标/5372225>

其中，dx 和 dy 分别表示 x 和 y 方向上一个像素的宽度，其倒数表示单位长度有多少像素（mm/像素），γ 是 x 和 y 不垂直时的倾斜系数，通常为 0。

由(1)(2)可得，

$$Z_c \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{f}{dx} & 0 & u_0 & 0 \\ 0 & \frac{f}{dy} & v_0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{pmatrix}$$

将其写作，

$$Z_c \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} f_x & 0 & c_x & 0 \\ 0 & f_y & c_y & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{pmatrix} \quad (3)$$

其中  $f_x$ ， $f_y$  分别为 x，y 方向上焦距长度（像素）， $\begin{pmatrix} f_x & 0 & c_x & 0 \\ 0 & f_y & c_y & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$  被称为内参矩阵。

再讲从世界坐标映射到相机坐标，世界坐标系可以通过旋转和平移转换到相机坐标系，将这两个变换复合，

$$\begin{pmatrix} I & t \\ \mathbf{0} & 1 \end{pmatrix} \begin{pmatrix} R & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} = \begin{pmatrix} R & t \\ \mathbf{0} & 1 \end{pmatrix}$$

$\begin{pmatrix} R & t \\ \mathbf{0} & 1 \end{pmatrix}$  被称为外参矩阵，那么有，

$$\begin{pmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{pmatrix} = \begin{pmatrix} R & t \\ \mathbf{0} & 1 \end{pmatrix} \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix} \quad (4)$$

由(3)(4)可得，

$$Z_c \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} f_x & 0 & c_x & 0 \\ 0 & f_y & c_y & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} R & t \\ \mathbf{0} & 1 \end{pmatrix} \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}$$

可以整理一下，

$$Z_c \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix} \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}$$

$$Z_c \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix} (\mathbf{R} \quad \mathbf{t}) \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix} \quad (5)$$

第 2 题答案应该就是(5)式。

这一节的推导过程主要参考了博客<sup>4</sup>和维基百科<sup>5</sup>。因为网上很多资料坐标系方向、坐标符号不同，所以自己画了坐标系，重打了公式，表述更加清晰一点。对于齐次坐标，和内参矩阵中的 $\gamma$ 系数还没理解透彻，在这记录下。

## 1.2. 图像上一点对应的 3D 形状

给定一个在像素坐标系中的 2D 点 $(u, v)$ ，求它在相机坐标系中对应的形状。

可以看到在 1.1 中的式子前面总含有  $Z_c$ ，我们可以把它看做是一个比例因子。因为空间到图像的映射是多对 1 的，联系实际，远处一个大的物体和近处一个小的物体在图像上投影确实有可能一样，这个比例因子就反映这种缩放关系。

不妨使(3)式的各分量相等，

$$\begin{cases} Z_c u = f_x X_c + c_x Z_c \\ Z_c v = f_y Y_c + c_y Z_c \end{cases}, Z_c > 0$$

<sup>4</sup> CSDN-深入解读相机矩阵，<https://blog.csdn.net/lingchen2348/article/details/83052214>

<sup>5</sup> 维基百科-Camera\_matrix，[https://en.wikipedia.org/wiki/Camera\\_matrix](https://en.wikipedia.org/wiki/Camera_matrix)



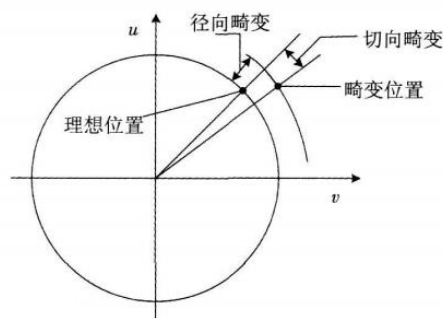
$$\begin{cases} X_c = \frac{Z_c}{f_x}(u - c_x) \\ Y_c = \frac{Z_c}{f_y}(v - c_y), Z_c > 0 \\ Z_c = Z_c \end{cases}$$

显然这是一条端点为点 O 的射线。引入齐次坐标，并写成矩阵形式，

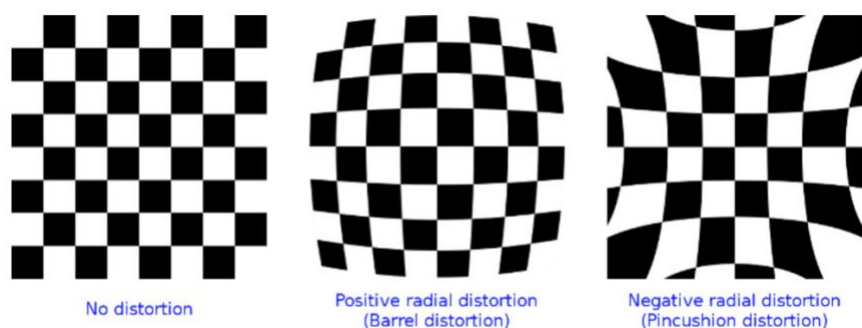
$$\begin{pmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{Z_c}{f_x} & 0 & \frac{-Z_c c_x}{f_x} \\ 0 & \frac{Z_c}{f_y} & \frac{-Z_c c_y}{f_y} \\ 0 & 0 & Z_c \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix}, Z_c > 0$$

### 1.3. 图像畸变

受镜头制造精度的影响，图像会出现不同程度的畸变，这种畸变可以分为径向畸变和切向畸变两种。下图来自百度百科<sup>6</sup>，



径向畸变的影响远大于切向畸变，OpenCV 文档显示了其两种常见类型，



<sup>6</sup> 百度百科-切向畸变，<https://baike.baidu.com/item/切向畸变/4947159?fr=aladdin>

看了这篇博客<sup>7</sup>才更好地理解畸变，其实就是正确的像素值产生了偏移，本属于点(u,v)的像素值，由于畸变，跑到了点(u<sub>d</sub>,v<sub>d</sub>)上，所以为了校正畸变，需要找到点的映射关系，将点(u<sub>d</sub>,v<sub>d</sub>)的像素值赋给点(u,v)。

以下参考 OpenCV 文档<sup>8</sup>，1.1 中的转换等效于以下公式，

$$\begin{pmatrix} X_c \\ Y_c \\ Z_c \end{pmatrix} = \mathbf{R} \begin{pmatrix} X_w \\ Y_w \\ Z_w \end{pmatrix} + \mathbf{t}$$

$$x' = X_c / Z_c$$

$$y' = Y_c / Z_c$$

$$u = f_x x' + c_x$$

$$v = f_y y' + c_y$$

存在畸变时（只考虑 k<sub>1</sub>,k<sub>2</sub>,p<sub>1</sub>,p<sub>2</sub>，k<sub>1</sub>,k<sub>2</sub>是径向畸变系数，p<sub>1</sub>,p<sub>2</sub>是切向畸变系数），

$$x'' = (1 + k_1 r^2 + k_2 r^4) x' + 2p_1 x' y' + p_2 (r^2 + 2x'^2)$$

$$y'' = (1 + k_1 r^2 + k_2 r^4) y' + p_1 (r^2 + 2y'^2) + 2p_2 x' y'$$

$$\text{其中, } r^2 = x'^2 + y'^2$$

$$u_d = f_x x'' + c_x$$

$$v_d = f_y y'' + c_y$$

给定一个点(u,v)，可以算出 x'和 y'，代入上式可以得到点(u<sub>d</sub>,v<sub>d</sub>)。该点像素值是畸变前本应属于点(u,v)的，重新赋给它，就达到了校正的目的。这篇

<sup>7</sup> CSDN-OpenCV 学习(5): 图像畸变校正，

<https://blog.csdn.net/xholes/article/details/80599802>

<sup>8</sup> OpenCV-Camera Calibration and 3D Reconstruction，

[https://docs.opencv.org/2.4/modules/calib3d/doc/camera\\_calibration\\_and\\_3d\\_reconstruction.html](https://docs.opencv.org/2.4/modules/calib3d/doc/camera_calibration_and_3d_reconstruction.html)

博客<sup>9</sup>提到畸变后的坐标  $u_d$  和  $v_d$  往往不是整数，所以需要利用相邻的四个像素点，通过双线性插值来得到点  $(u_d, v_d)$  的像素值，这种方法比较简单，网上资料讲的也很清楚，这里不再具体展开。

## 1.4. 相机标定，畸变校正

需要安装 OpenCV，这篇博客<sup>10</sup>讲了 macOS 下的安装过程以及会碰到的问题，cmake 不一定要像他一样用 GUI，用命令也很方便。再记录一个知乎上看到的 tip：从 github 上 git clone 速度很慢，可以先 fork 到自己的 github 仓库，然后导入到国内的码云仓库，从码云上 git clone 的速度超级快，亲测有效。

在网上找资料的过程中，发现官方提供了标定和畸变校正的 sample，所以直接拿来学习，一开始看不懂，放到 CLion 中 debug 单步跟一遍，一边查一下函数，就差不多了。所有函数都能在官方文档<sup>11</sup>检索到，这篇报告就不再对函数的使用展开说明。

相机标定就是根据几组 3D 世界坐标-2D 像素坐标的点，来求出相机矩阵。我们经常利用黑白相间的棋盘来标定，因为标准棋盘的制作成本低，将世界坐标系的原点设在棋盘的左上角，容易得到所有角点的世界坐标；至于各角点的像素坐标可以调用 OpenCV 的函数 `findChessboardCorners` 来得到。

再调用 `calibrateCamera` 就可以算出内参，外参和畸变参数，可能与旧版本的 OpenCV 提供的接口不同，但功能一样。

<sup>9</sup> CSDN-[图像]畸变校正详解，<https://blog.csdn.net/humanking7/article/details/45037239>

<sup>10</sup> 简书-OpenCV macOS：编译安装 OpenCV4+Opencv Contrib，<https://www.jianshu.com/p/162f2cdf4f88>

<sup>11</sup> OpenCV 官方文档，<https://docs.opencv.org/master/>

原本的 sample 的使用<sup>12</sup>还是有些麻烦，需要事先生成图像列表，查到一个 glob 函数可以读取一个文件夹下的所有图片路径，利用它可以使程序的使用更简单。原本的 sample 还考虑了影像，我也对此做了简化。

关于畸变校正，相机标定后我们已经得到了内参矩阵（OpenCV 的示例和源码中都将其命名为 cameraMatrix，但我认为相机矩阵应该是内参矩阵和外参矩阵的乘积）和畸变参数，回顾 1.3 畸变校正。可以利用 `undistort` 函数，一步到位得到校正好的图像；也可以先利用 `getOptimalNewCameraMatrix` 函数得到考虑了畸变的参数矩阵，再用 `initUndistortRectifyMap` 函数得到畸变前后点的映射关系，最后用 `remap` 得到校正好的图像。

## 1.5. 张正友标定

张正友标定的原文<sup>13</sup>写的很详细，有了之前的基础并不难看懂，这里不把公式一一列举了，大致理一下实现思路：

1. 一张图（至少 4 个角点）可计算出一个单应性矩阵  $H$ （8 自由度）
2. 每个  $H$  可得到相应的  $v$ ，并组成  $V$
3. 用  $Vb=0$  解出  $b$ （没有其他条件的话，至少 3 张图才有非 0 解）
4. 由  $b$  算出各内参，即算出了论文中的矩阵  $A$
5. 由  $A$  和之前每张图的  $H$  算出每张图对应的外参
6. 用 Levenberg-Marquardt 算法（以下简称 LM 算法）优化参数

<sup>12</sup> 百度经验-OpenCV：相机标定示例程序的使用，

<https://jingyan.baidu.com/article/7f41ecec5877eb593d095ce9.html>

<sup>13</sup> Z. Zhang, "A flexible new technique for camera calibration," IEEE Transactions on pattern analysis and machine intelligence, vol. 22, no. 11, pp. 1330–1334, 2000.

畸变参数不会很大，不妨就把初值设为 0，和其他参数一起用 LM 算法优化，切向畸变影响比较小，通常不需要考虑。论文中几何解释的部分没怎么看懂，虽不影响实现，但我认为是重要的。

做之前还是先看看 OpenCV 相关函数的源码，用 CLion，发现 debug 不进去，查了资料知道自己之前安装了 release 版本，又重装了一遍 OpenCV，cmake 命令还不是很熟悉，在这记录下：

- `cmake -D CMAKE_BUILD_TYPE=DEBUG -D CMAKE_INSTALL_PREFIX=/usr/local -D OPENCV_EXTRA_MODULES_PATH=/Users/wangconghao/opencv_contrib/modules ..`

源码中经常出现 InputArray 这种参数类型，这样的话用户给出 Mat，Matx 或 vector 类型都能接受，然后再用 getMat 转换成 Mat。

跟到初始化参数时，感觉源码和论文的方法有些出入，这部分我自己按论文实现了，一边查一下 Mat 的基本操作不很困难。有一些细节的实现论文没有给出具体方法，但 OpenCV 提供了许多强大的函数：

- 转换为齐次坐标，*convertPointsToHomogeneous*
- 求单应性矩阵，*findHomography*
- 求解齐次线性方程组的非 0 解，*SVD::solveZ*
- 将旋转矩阵转换为旋转向量，*Rodrigues*

优化阶段，我对于 LM 算法不是很熟悉，看了知乎<sup>14</sup>，再结合比较熟悉的梯度下降，大致了解了它是利用雅克比矩阵来计算当前优化的方向。

<sup>14</sup> 知乎-如何用 LM 算法求解目标函数最小值？，

<https://www.zhihu.com/question/269579938/answer/349205519>

Algorithms	Update Rules	Convergence	Computation Complexity
EBP algorithm	$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{g}_k$	Stable, slow	Gradient
Newton algorithm	$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}_k^{-1} \mathbf{g}_k$	Unstable, fast	Gradient and Hessian
Gauss-Newton algorithm	$\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \mathbf{J}_k^T \mathbf{J}_k \right)^{-1} \mathbf{J}_k \mathbf{e}_k$	Unstable, fast	Jacobian
Levenberg-Marquardt algorithm	$\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \mathbf{J}_k^T \mathbf{J}_k + \mu \mathbf{I} \right)^{-1} \mathbf{J}_k \mathbf{e}_k$	Stable, fast	Jacobian
NBN algorithm [08WC] <sup>a</sup>	$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{Q}_k^{-1} \mathbf{g}_k$	Stable, fast	Quasi Hessian <sup>a</sup>

本想使用 cminpack，但参考资料比较少。找资料的过程中发现 OpenCV 还提供了一个 LMSolve 抽象类，除了官方文档也没什么其他资料，只能在 OpenCV 官方 github 上搜搜看源码对其的使用，但对于用在论文复现上也没什么头绪。

OpenCV 源码中实现相机标定的函数用了 CvLevMarq，所以修改了一下源码拿过来用了，过程中还用到了一些在新版本中已经被封装的内部接口，如 cvProjectPoints2Internal，为了在新版本下能运行，我找到它的声明和定义，抓出来放在自己的项目下，work！

**最终结果**，跑出来的误差是 1.4 中结果的两倍。内参相差不大，有些图的外参相差很大，我发现源码中在初始化外参的过程中也用到了优化技术，应该和这个有点关系。这里不考虑切向畸变，径向畸变也只考虑了  $k_1, k_2$ ，如果考虑和 1.4 中相同的畸变参数，得到结果会稍许相近一些，但影响不大。按论文的方法初始化得到的内参  $\gamma$  不是 0，但相对于其他内参很小，所以将它设为 0 进行后面的工作。

在畸变校正的时候有一些问题，用 1.4 中第二种方法校正后会出现很奇怪的图像，就如同 OpenCV 问答论坛<sup>15</sup>有人提到过的这样，检查了一下发现由

<sup>15</sup> OpenCV 问答论坛-Undistortion at far edges of image，

<https://answers.opencv.org/question/28438/undistortion-at-far-edges-of-image/>

getOptimalNewCameraMatrix ( 参数 $\alpha=1$  ) 得到的新矩阵很奇怪, 不知何原因。用 undistort 表现正常。

OpenCV 源码实现相机标定还是考虑了比较多的, 主要体现在 flags 参数上, 比如, 如果用户设置了 CV\_CALIB\_FIX\_ASPECT\_RATIO 模型, 就是固定了  $f_x/f_y$  的比值, 那么只用其中一个量进行优化计算, 将更准确。相关内容网上有不少的参考资料, 这里不再将每个模型拿来讨论。

## 2. 双目基础

### 2.1. 更新中...



# 3. 立体匹配