

1 SLAM 中的几何与学习方法
2 Geometric and Learning Approaches Used in
3 SLAM systems



4

5 盐粒

<https://github.com/yanyan-li/SLAM-BOOK>

6

updated 2021 年 6 月 11 日

7	目录	
8	1 书稿概述	2
9	1.1 执笔初衷	2
10	1.2 主要内容和面向读者	2
11	2 Hello SLAM	4
12	2.1 SLAM 研究方向的变迁	4
13	2.1.1 基于特征的 SLAM 框架	4
14	2.1.2 多几何特征	5
15	2.1.3 多传感器	6
16	2.1.4 几何与学习的融合	7
17	2.2 动手来写搭建开发环境	7
18	2.2.1 操作系统与开发软件	7
19	2.2.2 常用的开发库介绍	7
20	3 几何视觉基础	10
21	3.1 图像处理基础	10
22	3.1.1 RANSAC	10
23	3.2 几何特征提取	11
24	3.2.1 点特征提取与匹配	11
25	3.2.2 线特征提取与匹配	14
26	3.2.3 面特征提取与匹配	17
27	3.3 几何特征之间的关系	18
28	3.3.1 直线与直线平行-灭点	18
29	3.3.2 平面与平面平行-灭线	19
30	3.3.3 平面与平面垂直-曼哈顿约束	19
31	3.3.4 点在直线上-共线	20
32	3.3.5 点线在平面上-共面	20
33	4 深度学习视觉基础	22
34	4.1 经典网络回顾	22
35	4.2 优化器与学习率	25
36	4.3 常用的度量方法	27

37	5 基于卷积网络的学习任务	29
38	5.1 点线面的预测	29
39	5.1.1 预测特征点	29
40	5.1.2 预测特征线	29
41	5.1.3 平面实例分割	29
42	5.2 物体检测与分割	29
43	5.2.1 分类学习基础	29
44	5.2.2 分类网络中常见策略	31
45	5.3 深度图与法向量预测	32
46	5.3.1 监督学习	32
47	5.3.2 动手训练一个法向量预测网络	35
48	5.3.3 无监督学习	43
49	6 基于图网络的学习任务	44
50	6.1 图网络简介	44
51	6.2 如何定义一个图网络层	49
52	7 传感器模型	53
53	7.1 相机模型	53
54	7.2 惯性导航模型	54
55	7.3 激光雷达模型	55
56	8 位姿估计的模型	56
57	8.1 2D-2D 位姿估计	56
58	8.2 3D-2D 位姿估计	57
59	8.3 3D-3D 位姿估计	60
60	8.4 位姿解耦估计方法	60
61	8.4.1 曼哈顿和亚特兰大假设	61
62	8.4.2 基于特殊环境结构的无漂旋转估计	62
63	8.5 ORB-SLAM2 中的位姿策略	63
64	8.6 VIO 中的位姿估计	65
65	8.6.1 初始化	66

66	9 地图创建	68
67	9.1 稀疏地图重建	68
68	9.1.1 稀疏点重建	68
69	9.1.2 线特征的重建。	71
70	9.2 稠密地图重建	72
71	9.2.1 Mesh	73
72	9.2.2 TSDF	74
73	9.2.3 Surfel	75
74	10 非线性优化	76
75	10.1 光束法平差模型分析	76
76	10.1.1 点特征的误差方程	76
77	10.1.2 线特征的误差方程	77
78	10.1.3 面特征的误差方程	77
79	10.1.4 构建光束法平差模型	78
80	10.2 经典的光束法平差解算方法	78
81	10.2.1 高斯牛顿方法	78
82	10.2.2 Levenberg-Marquardt	80
83	10.2.3 稀疏结构化 BFGS 方法	81
84	11 优化中的位姿描述	84
85	11.1 四元数	84
86	11.2 李群与李代数	84
87	12 优化中的特征描述	84
88	12.1 点特征-逆深度	84
89	12.2 线特征	84
90	12.3 面特征	88
91	13 经典的纯几何 SLAM	92
92	13.1 ORBSLAM (V1 and V2)	92
93	13.1.1 单目系统的初始化	92
94	13.1.2 姿态估计与优化	93
95	13.1.3 基于词袋的闭环检测	94
96	13.2 VINS 系列	95

97	13.2.1 系统概述	95
98	13.2.2 预积分和视觉约束	95
99	13.3 LSD-SLAM	96
100	13.3.1 整体流程	96
101	13.3.2 光度误差与深度误差	97
102	13.4 SVO	97
103	13.4.1 位姿估计	98
104	13.4.2 深度估计	99
105	14 深度 SLAM	99
106	14.1 Structure-SLAM	99
107	参考文献	103

108

109

鲁迅

110

111

无穷的远方,无数的人们

112

都和我有关

113

114

舟夜书所见

115

渣慎行

116

117

月黑见渔灯,孤光一点莹

118

微微风簇浪,散作满天星

119

120

辋川集·南_田

121

王维

122

123

轻舟南_田去,北_田淼难即

124

隔浦望人家,遥遥不相识

125

126

叶芝

127

128

穿过我说谎的日子,

129

在阳光下抖掉枝条和花朵。

130

我现在可以枯萎而进入真理

131

132 **说明**: 文章中的实验结果图片, 除特殊说明以外, 均来自我们自己的实
133 验。读者可以引用, 引用时, 请说明出处。

```
134 @misc{gl-SLAM,  
135   author = {Yanyan Li},  
136   title = {SLAM中的几何与学习方法},  
137   year = {2020},  
138   publisher = {GitHub},  
139   howpublished = {\url{https://github.com/yanyan-  
140     li/SLAM-BOOK}  
141   }  
142 }  
143
```

144 1 书稿概述

145 1.1 执笔初衷

146 SLAM(Simultaneous Localization and Mapping) 关注传感器在环境中
147 的定位与建图问题，其中定位部分研究的是传感器的六自由度位姿估计问
148 题，而建图则是恢复未知环境中的三维信息。SLAM 系统提供的位姿和环
149 境信息可以应用在增强/虚拟现实、机器人系统和无人驾驶系统中，为机器
150 与环境的交互提供服务。

151 在 SLAM 社区，《Multiple View Geometry in Computer Vision》是
152 一种字典、辞海一样的存在，它提供了丰富的理论知识，当我们对某一个
153 概念不清楚的时候，可以去里面寻找答案。对于我们中文读者而言，《视觉
154 SLAM 十四讲》同样是一本非常精彩、扎实的书籍。高翔博士理论知识扎
155 实、代码功力深厚，在书中深入浅出的把 SLAM 画卷展开在中国读者面前。
156 与这些书籍相同的是，本书稿同样试图将作者自己对 SLAM 的理解重新梳
157 理和归类，用更简单的话术将基础理论介绍给读者。在学习 SLAM 的过程
158 中，作者和周围的一些朋友都曾有一种错觉：

159 SLAM这个领域已经被挖掘完了，很难再取得突破。
160
161

162 产生这种错觉的主要原因可能是在于：我们更多的是在学习 SLAM 的
163 历史和基础理论，而不是去最新的进展中发觉尚未解决的问题。例如：ORB-
164 SLAM 和 ORB-SLAM2 是非常重要的、完整的 SLAM 系统，他们分别发表于
165 2015 年和 2017 年，到 2020 年 10 月，他们的引用量分别是 3300 和 2000。
166 这意味着 SLAM 领域的发展是十分迅速的。

167 基于此，本书稿的一个野心是为读者带去更新、更前沿的 SLAM 知识。
168 通过梳理诸多新颖工作的 motivation 和他们关注的问题，我们在介绍充足
169 的理论之后，努力的去探索这个领域的未来的可能性，以便在 SLAM 前沿
170 研究中找到我们自己方向。

171 1.2 主要内容和面向读者

172 **主要内容** 对于不同的传感器，SLAM 系统在位姿估计和重建中应用的方
173 法也大相径庭。对于依赖相机的（包括 RGB-D）系统，位姿估计可以分为
174 $2D - 2D$ ， $2D - 3D$ 和 $3D - 3D$ 几种方式；对于激光雷达，有 $scan - scan$

175 和 *scan-map* 的位姿估计策略；而对于 IMU 和 GPS，则由预积分和差分
176 法等方法进行位姿估计。而重建部分则包括稀疏，半稠密和稠密几种方式。

177 近年来，深度学习在更高层次的视觉任务中取得瞩目的成绩，如：物体
178 识别，语义分割等，这些课题曾是传统视觉无法或很难解决的任务。正是由
179 于深度学习拓展了我们对高层次视觉任务的想象，越来越多的 SLAM 开始
180 在他们的框架中融合学习的方法，来改进位姿估计的准确程度和环境重建
181 的效果。但是深度学习是一个非常宽广的领域，和 SLAM 相关的课题只是
182 它的一个分支。本书将会挑选和聚焦和 SLAM 相关的深度学习话题，我们
183 希望能通过这本书来介绍 SLAM 系统中使用的几何和深度学习的方法，帮
184 助读者掌握最新的进展。从整体上理解基于几何的 SLAM 系统的技术细节，
185 同时还能了解深度学习相关课题对传统几何 SLAM 的影响。

186 本书首先分析单帧图像入手：单帧图像的几何特征提取，单帧图像的深
187 度学习任务分析。通过相机模型和刚体变换的规则，我们把图像帧联系起
188 来，构建帧与帧之间的关系。然后按照 SLAM 中的每一个重要环节来详细
189 介绍相关的知识。几何与学习特征的处理、优化问题的构建与求解，以及环
190 境地图的创建与维护。在每个环节中，我们都深入的介绍相关的理论与工程
191 细节，力求本书为读者带来快速的提高。几何与学习就像是 SLAM 的一双
192 翅膀，未来的 SLAM 将提供更加的鲁棒、更加丰富的信息。最后，我们一
193 起探讨 SLAM 的发展方向，畅想这个领域未来蓝图。

194 在诸多种类的传感器中，本书稿主要和诸位读者一起探讨以相机为中
195 心的 SLAM 系统，包括单目、双目 SLAM 系统、RGB-D SLAM 以及 VIO。
196 围绕 RGB 相机，本书在介绍基础理论的同时，将深入探讨一下几个问题。

- 197 • 在传统的视觉几何中，如何使用点、线、面特征来提高视觉系统定位
198 的精度和重建的效率。
- 199 • 深度学习在诸多视觉任务中取得了令人瞩目的表现，哪些学习信息可
200 以被融合到 SLAM 系统中。
- 201 • 为关键的理论提供相关的代码，以及完整的 SLAM 系统。

202 **面向读者** 我们希望这本书稿能够给为研究 SLAM 的学生和相关从业者系
203 统介绍几何特征与学习信息在 SLAM 系统中的提取和使用过程。更简单的
204 说，如果诸位对《SLAM 视觉十四讲》感兴趣，那这本书稿就是为你而写的。

205 2 Hello SLAM

206 SLAM(Simultaneous Localization and Mapping) 关注传感器在环境中
207 的定位与建图问题, 其中定位部分研究的是传感器的 6D 的位姿估计问题,
208 而建图则是试图恢复环境的 3D 信息。SLAM 系统提供的位姿和环境信息
209 可以应用在增强/虚拟现实, 无人驾驶系统中, 为机器与环境的交互提供服务
210 务。

211 根据传感器的不同, SLAM 系统在位姿估计和重建中应用的方法也大
212 相径庭。对于相机 (包括 RGB-D), 位姿估计可以分为 $2D - 2D$, $2D - 3D$
213 和 $3D - 3D$ 几种方式; 对于激光雷达, 有 *scan - scan* 和 *scan - map* 的
214 位姿估计策略; 而对于 IMU 和 GPS, 则由预积分和差分法等方法进行位姿
215 估计。而重建部分则包括稀疏, 半稠密和稠密几种方式。

216 2.1 SLAM 研究方向的变迁

217 SLAM 领域经过二三十年的高速发展, 逐渐成熟起来。每一篇划时代的
218 论文都获得了巨大的引用量, 通过阅读这些经典, 我们可以感受整个 SLAM
219 研究方向的变迁。

220 2.1.1 基于特征的特征的 SLAM 框架

221 对于基于特征法的 SLAM 系统, 发表于 2007 年的 PTAM [16] 算是一
222 个重要的里程碑, 论文提出了基于多线程的 tracking 和 mapping 策略, 以
223 及关键帧的发展道路。其中 Mapping 部分是只处理关键帧的特征信息, 并
224 使用了局部光束法平差 (Bundle Adjustment) 来优化路标点和相机位姿。
225 比起 2007 年之前的 SLAM 或视觉里程计系统相比, 它可以实时维护几千
226 个点的地图。

227 同样是基于点特征的系统, 发表于 2017 年的 ORB-SLAM2 [24] 是
228 一个集大成的作品。这个系统是由 Raúl Mur-Artal 博士等人提出的一个完
229 备的基于特征点的纯视觉实时 SLAM 系统, 适用于单目、双目和 RGB-D
230 相机。作者在 2015 年第一次将单目版本 (ORB-SLAM [23]) 发表在 IEEE
231 Transactions on Robotics (T-RO) 期刊上, 这是一个基于纯特征点的单目
232 实时 SLAM 框架。两篇文章涉及的系统是非常完整的系统, 通过多线程实
233 现了稀疏地图、回环和优化环节, 为 SLAM 在工业界和学术界的发展作出
234 巨大贡献。

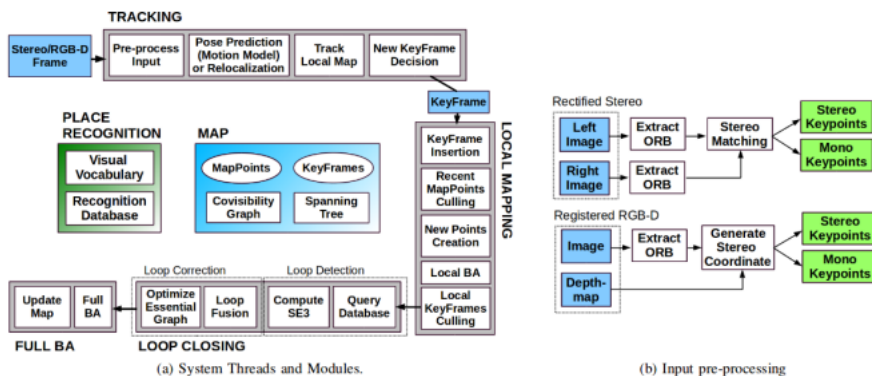


图 1: 特征法 SLAM (ORB-SLAM2) 的经典框架。

235 在 ORB-SLAM (V1/V2) 工程的基础上, 图 5 中的各个模块都吸引研究
 236 人员的注意力, 最常见的就是为系统增加其他的特征信息。众所周知, 点
 237 特征对于光照纹理都比较敏感, 尤其是在纹理比较弱的地方, 提取特征点
 238 的数目是很少的。为了应对低纹理场景中特征点能力的限制, 很多点线, 或是
 239 点线面的系统被提出来满足单目 [26]、双目 [9] 和 RGB-D 相机 [35]。随着
 240 新特征的加入, 关键帧策略、地图和优化环节都会做相应的调整。

241 2.1.1.2 多几何特征

242 ORB-SLAM 系列利用的特征信息都是特征点, 基于这些点特征系统,
 243 后续很都工作都尝试去增加新的特征类型。对于单目/双目传感器, 早期的
 244 工作 PL-SLAM 在系统中增加直线信息; 对于 RGB-D 传感器, 则有工作尝
 245 试增加平面信息。在这些早期的工作中, 增加的直线和平面一般主要是为了
 246 应对特征点较难提取的若纹理区域。由于没有充分利用这些信息的特点, 即
 247 便是系统增加了特征的种类, SLAM 系统仍旧无法较大的提高系统的表现。

248 不同于传统的点-线-面系统, RGB-D SLAM [21] 进一步探索直线与直
 249 线, 平面与平面之间的平行/垂直关系。由于, 这种几何关系满足 Manhattan
 250 World 约束, 这种约束带来了零漂移旋转估计和 3D 平移估计的策略。同时,
 251 在提取平面特征之后, 我们对环境重建的能力也增强了, 可以在 CPU 基础
 252 上对环境进行实时 mesh 重建。

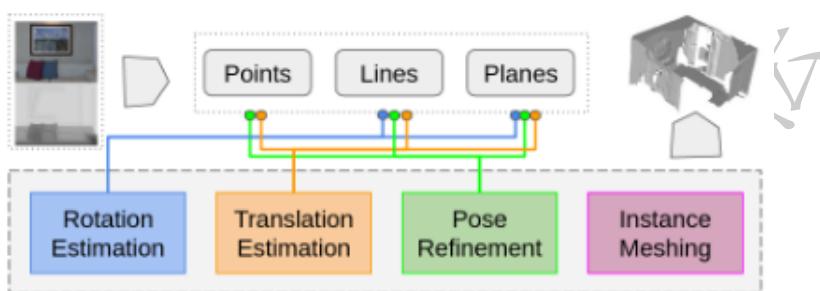


图 2: 多特征 SLAM 系统框架。

2.1.3 多传感器

253 由于运动模糊和位姿剧烈变化，快速运动情况为基于特征法的 SLAM
 254 系统带来较大的挑战，因此研究人员通过增加传感器为系统增加更多的信
 255 息源。惯性测量单元 (IMU) 是一款非常便宜的设备，可以提供高帧率位
 256 姿估计结果。VINS-Mono [28] 是相机和 IMU 紧密耦合的经典系统，由于
 257 IMU 可以高频的记录角速率和加速度信号，通过对高频信号进行预积分操
 258 作，系统可以稳定的获得两帧图像之间的相对位姿，然后通过相机带来的视
 259 觉特征的约束来实现位姿信息的优化。然后也有结合 IMU 与 ORB-SLAM
 260 或 DSO 的工作 [34]，他们都取得了不错的效果。同样是加入线特征的策略，
 261 PL-VIO [10] 在 VINS-Mono 的基础上加入了线特征。
 262

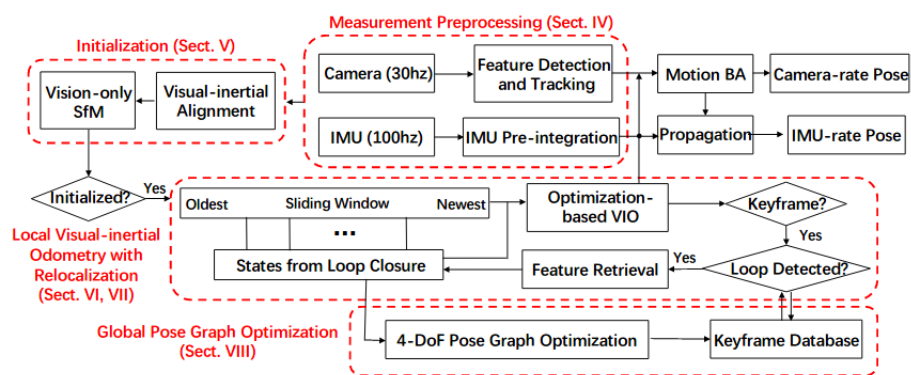


图 3: VIO 系统 (VINS-Mono) 经典框架。

2.1.4 几何与学习的融合

深度学习的出现帮助传统的 SLAM 系统拓展了传感器的边界，针对 SLAM 的目标，结合的目标主要是提升重建和位姿估计。比较早的结合工作有：DVS0、CNN-SLAM 这两篇工作都是发表在 2018 年，并且都是使用深度学习的估计的深度信息去恢复。

除了利用估计的深度信息帮助单目 SLAM 实现密集重建，还有一些新的工作直接利用深度学习提高位姿估计的效果。值得注意的是，深度学习提

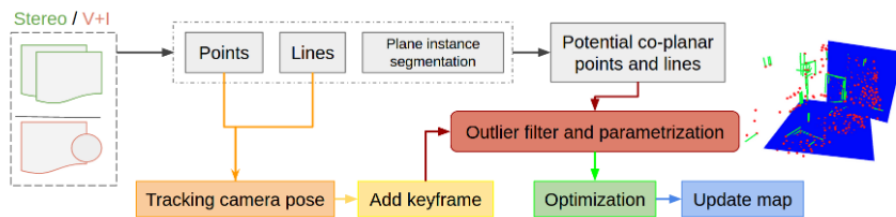


图 4: 基于学习平面的共面参数化方法 (CP-Param [18])。

供的特征并不像传统几何特征那样稳定，因此在为 SLAM 引入学习信息的时候，也要同时加入一个去伪存真的策略。

2.2 动手来写搭建开发环境

2.2.1 操作系统与开发软件

- Ubuntu 系统。很多系统都可以去开发 SLAM，但是方便的系统应该是 Linux 系统。Ubuntu 是 Linux 系统中的一种，系统简洁、操作方便，十分适合开发 SLAM、研究深度学习的朋友们。
- CLion + terminal。CLion 是一款非常优秀的 C++ 集成编译环境。可以在 Ubuntu 环境下使用，对于不熟悉 vim 的小伙伴，它非常容易使用，并且为学生用户提供免费使用的权利。

2.2.2 常用的开发库介绍

由于视觉前辈的卓越贡献，我们在编写 SLAM 系统的时候，不需要完全从底层编写。通过使用的常见开发库中的函数接口，我们就可以直接获得我们需要的阶段性结果。这些开发库提高了我们的开发效率，但是这些库是

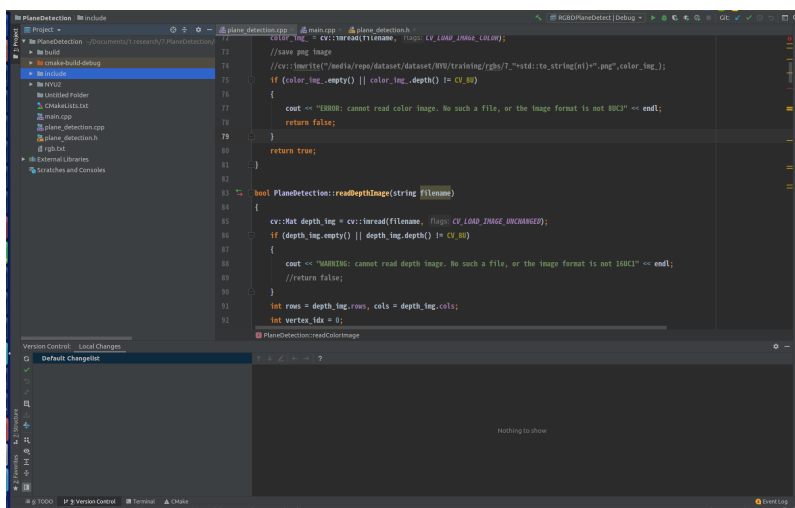


图 5: CLion 软件界面。

284 否可以免费用在商业代码中，还需要读者自己去研读相应库的开发者使用
285 规范。

286 由于每个库本身都是内容极为丰富，例如 OpenCV 本身就有专门的书
287 籍去介绍，因此我们这里只简单介绍与本书稿有关的部分，旨在让诸位对这
288 些库的架构有一个大概的了解。

289 **OpenCV**

290 **PCL**

291 **g2O**

292 **Pangolin**

293 **OpenGL**

294 **ROS** ROS 系统是机器人领域很重要的平台，在大多数没有使用 ROS 的
295 SLAM 系统中，我们大概率都是使用 cmake 和 make 命令来完成代码编译，
296 这种情况在 ROS 相关的系统发生了一些变化，通过对 cmake 进行扩展，ROS
297 系统构建了自己的编译工具 catkin。要用 catkin 编译一个工程或软件包，只

298 需要用 `catkin_make` 指令。一般当我们写完代码，执行一次 `catkin_make` 进
299 行编译，调用系统自动完成编译和链接过程，构建生成目标文件。

300 Package 是 catkin 编译的基本单元，我们调用 `catkin_make` 编译的对象
301 就是一个个 ROS 的 package，也就是说任何 ROS 程序只有组织成 package
302 才能编译。所以 package 也是 ROS 源代码存放的地方，任何 ROS 的代码
303 无论是 C++ 还是 Python 都要放到 package 中，这样才能正常的编译和运
304 行。

盐粒-SLAM 中的几何与数学

305 3 几何视觉基础

306 计算机视觉基础是一个很庞大的课题，为了更贴近 SLAM 的研究主题，
307 我们聚焦在视觉几何的部分，而这一章主要介绍 SLAM 中常用的图像/数据
308 处理方法。

309 3.1 图像处理基础

310 3.1.1 RANSAC

311 随机抽样一致算法 (Random Sample Consensus, RANSAC)¹ 采用迭
代的方式从一组包含外点的数据中估算出需求模型的参数。

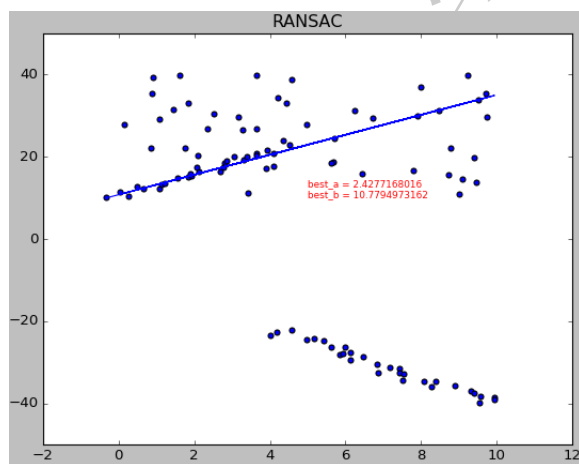


图 6: RANSAC 在 2D 数据中拟合直线。

312

- 313 • 随机从数据 D 中取最小抽样集 S ，其中抽样点个数为 n ，是模型初始
314 化所需的最小样本数，基于抽样集 S 估计出模型 M 。
- 315 • 计算最小抽样集 S 以外的数据与模型 M 之间的距离，距离小于阈值
316 t_d 的数据构成该模型的内点集 I_M 。
- 317 • 若 I_M 中的个数大于阈值 t_n ，则认为是找到模型 M 是该数据包含的正
318 确模型。

¹[https://github.com/yanyan-li/SLAM-BOOK/tree/master/code/chapter 2](https://github.com/yanyan-li/SLAM-BOOK/tree/master/code/chapter%20)

319 我们重复以上步骤多次, 比较哪次计算中内群数量最多, 最多内群点对应的
320 模型就是该数据所要求模型。值得指出的是, 图6中输入为两条直线, 但是
321 上述介绍的方式找到一个直线后就结束了。

3.2 几何特征提取

323 点线面特征被广泛的应用在在 SLAM 中, 这些传统的特征提取与匹配
324 方式, 经过长时间的发展, 已经逐渐稳定。

3.2.1 点特征提取与匹配

326 经常被使用的特征有 SIFT(Scale Invariant Feature Transformation
327) [19]、SURF(Speeded up robust feature) [20], ORB(Oriented fast and
328 rotated BRIEF) [21] 等。这些特征通过不同的算法获得, 也有不同的应用场
329 景。

330 **SIFT 特征** SIFT 特征匹配算法是早期非常著名的特征描述和匹配算子,
331 为计算机视觉的发展起到了巨大的推动作用。这种特征可以很好的应对图
332 像对之间的平移、旋转、尺度变换, 众多的实验表明了 SIFT 在上述情况下
333 的鲁棒性。使用 SIFT 算子惊醒特征提取工作主要包括两个步骤: 分别是特
334 征描述向量的生成, 如图 4 所示; 另一个是特征描述向量之间的匹配。同时
335 这种两步方法也被后来众多的特征匹配算子继承, 甚至影响着使用深度学习
完成匹配工作的网络框架。SIFT 算法可以很好的应对尺度、旋转情况下的

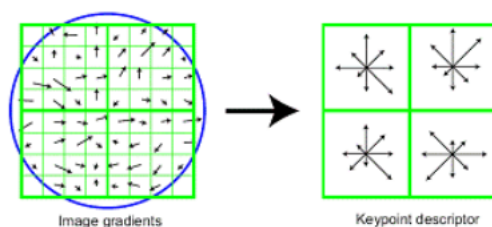


图 7: SIFT 特征提取的 128 维向量

336
337 匹配, 但是该算法建立在两幅图像满足相似变换这一假设。因此该算法适合
338 处理相两幅图像邻角度小于 30° 的图像, 对于较大倾角的倾斜图像, SIFT
339 只能提取很少量的特征, 同时伴随着较大比例的错误匹配, 因此 ASIFT 算
340 子 [22] 被设计出来解决大倾角图像的匹配问题。

341 **SURF 特征** SURF 算法同样是经历局部特征点的提取、描述和特征点的
 342 匹配,但它一种稳健的具有尺度和旋转不变性的局部特征点检测和描述算
 343 法,它降低了特征描述向量的维度,同时这个算子更加容易实现并行处理。
 344 这种算子是对 SIFT 进行改进,同样使用到图像空间金字塔等结构。不同于
 345 SIFT 算法,它使用黑塞矩阵 (Hessian Matrix) 和盒式滤波器去描述兴趣点
 346 (Faster-Hessian Detector)。在特征的描述环节,该方法首先是在兴趣点周围
 347 的圆形区域计算方向,然后再根据选择的方向构建方形区域,并从中提取特
 348 征点描述符。SURF 特征是一种具有较高准确度和速度的特征算子。

349 **ORB 特征** 由于 ORB-SLAM 系列和相关论文的影响,ORB 这种特征变
 350 得非常流行,它运行速度非常的快,我在比较好算力的台式机上从 $480 * 664$
 351 图片提取 ORB 特征并建立描述子需要 8ms 左右。那么 ORB 中有那些操
 352 作呢?

```
353 ComputePyramid(image);
354 ComputeKeyPointsOctTree(allKeypoints);
355 _descriptors.create(nkeypoints, 32, CV_8U);
356
357
```

358 首先,我们将图片建立金字塔,在金字塔每层图像上进行 FAST 特征提取。
 359 对于感兴趣的点,如上图所示的 p 点,以该点为中心,3 像素为半径,获得
 360 圆周上的 16 个像素。如果圆周上有连续 n 个像素点的灰度值大于或小于该
 点的灰度值,则认为该点为 FAST 特征。

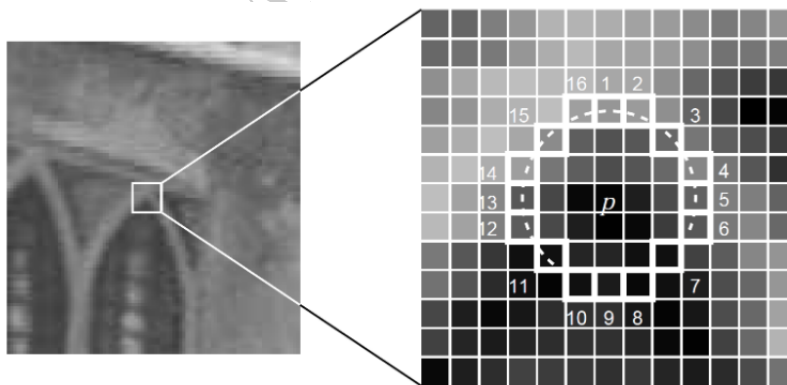


图 8: 感兴趣点 p 和圆周上的 16 个像素。 [33]

361

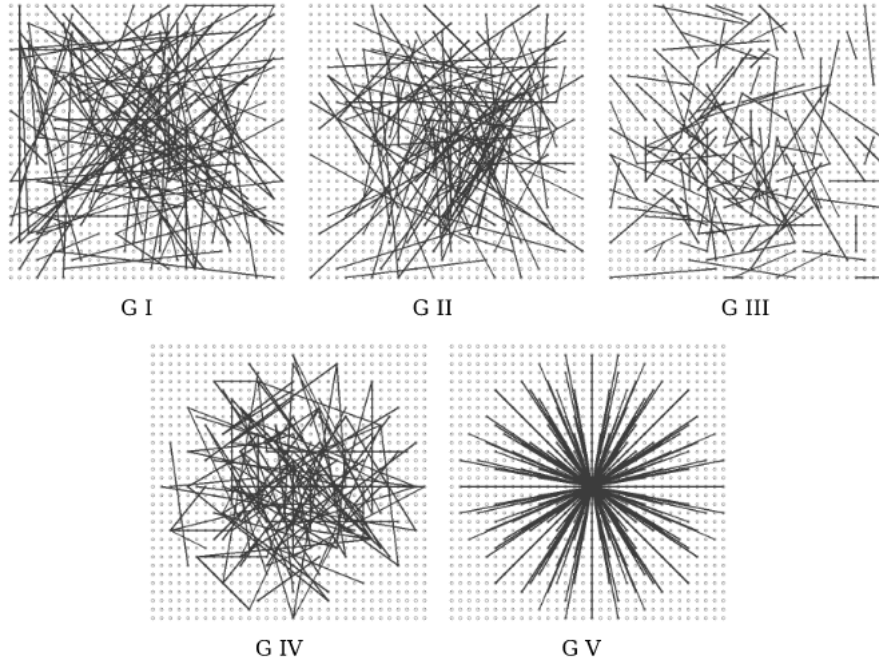


图 9: BRIEF [1] 不同选择点的方式。1. x, y 方向平均分布采样; 2. x, y 均服从 $\text{Gauss}(0, S^2/25)$ 各向同性采样; 3. x 服从 $\text{Gauss}(0, S^2/25)$, y 服从 $\text{Gauss}(0, S^2/100)$ 采样; 4. x, y 从网格中随机获取; 5. x 一直在 $(0, 0)$, y 从网格中随机选取。

362 其次，我们要去除不好的 FAST 点。FAST 特征点和该点周围的 16 个
363 像素输入决策树，筛选出来高质量的 FAST。对于临近的特征点，选择具有
364 较大区分度的特征点，删除临近的其余点。

365 为了实现特征点的旋转不变性。对于特征点和 r 半径内的像素点，计算
366 r 半径范围内的质心，并建立从特征点位置到质心位置的向量，这个向量就
367 是该特征点的方向。具体的计算过程如下所示：

$$m_{pq} = \sum_{x,y} x^p y^q I(x, y) \quad (3.1)$$

368 这里的 x 和 y 表示特征点为中心， r 为半径的圆内点的坐标，他们范围是
369 $(-r, r)$ 。利用上面的动量公式：我们可以计算重心：

$$C = \left(\frac{m_{10}}{m_{00}}, \frac{m_{01}}{m_{00}} \right) \quad (3.2)$$

370 因此围绕这个特征点的图像块的方向可以表示出来了：

$$\theta = \text{atan2}(m_{01}, m_{10}) \quad (3.3)$$

371 上述公式中的 $\text{atan2}()$ 是象限敏感 (quadrant-aware) 的反正切函数。

372 在获得特征点的位置和方向之后，我们一起来看看如何为这个特征点建
373 立描述子。ORB 在 BRIEF 的基础上改进出来自己的描述方法，我们先来
374 回顾原始的 BRIEF 的过程。对于 $n = 256$ 个点 (p_i, p_j) ，每一对之间的
375 灰度值 I 高低就产生了一个二值化结果。

$$r(p_i, p_j) = \begin{cases} 1 & I(p_i) < I(p_j) \\ 0 & I(p_i) > I(p_j) \end{cases} \quad (3.4)$$

376 因此 256 个像素对就构成了一个 256 维的向量。那么 BRIEF 是怎么选择
377 这些点对呢？如图9所示，原论文给出了 5 种不同的思路，但是作者使用第
378 二种：两个点服从高斯 $(0, S^2/25)$ 各向同性采样。

379 与传统 BRIEF 相比，ORB 为描述子增加了方向信息，其被称为 Steered
380 BRIEF。原先的 n 个像素对构成了矩阵 S ，

$$S = \begin{pmatrix} p_i^0 & \cdots & p_i^1 \\ p_j^0 & \cdots & p_j^1 \end{pmatrix} \quad (3.5)$$

381 把方向信息考虑在整体的描述子中，我们可以获得 ORB 特征构建的具
382 有旋转不变性的描述子：

$$S_\theta = R_\theta S \quad (3.6)$$

383 上述公式中的 $R_\theta = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix}$ 。

384 3.2.2 线特征提取与匹配

385 LSD 是比较流行的一个直线提取方法，KNNMatch 做特征描述。同时，
386 还有特征质量的简单筛选。

```
387
388 using namespace cv;
389 using namespace std;
390 using namespace cv::line_descriptor;
391 struct sort_descriptor_by_queryIdx
```

```
392 {
393     inline bool operator()(const vector<DMatch>&
394         a, const vector<DMatch>& b){
395     return ( a[0].queryIdx < b[0].queryIdx );
396     };
397     struct sort_lines_by_response
398     {
399         inline bool operator()(const KeyLine& a,
400             const KeyLine& b){
401             return ( a.response > b.response );
402         }
403     };
404
405     void ExtractLineSegment(const Mat &img, const Mat
406         &image2, vector<KeyLine> &keylines, vector<
407         KeyLine> &keylines2)
408     {
409         Mat mLdesc, mLdesc2;
410
411         vector<vector<DMatch>> lmatches;
412
413         Ptr<BinaryDescriptor> lbd = BinaryDescriptor
414             ::createBinaryDescriptor();
415         Ptr<line_descriptor::LSDDetector> lsd =
416             line_descriptor::LSDDetector::
417             createLSDDetector();
418
419         cout<<"extract lsd line segments"<<endl;
420         lsd->detect(img, keylines, 1.2,1);
421         lsd->detect(image2, keylines2, 1.2,1);
422         int lsdNFeatures = 50;
423         cout<<"filter lines"<<endl;
424         if(keylines.size()>lsdNFeatures)
```

```
425     {
426     sort(keylines.begin(), keylines.end(),
427         sort_lines_by_response());          keylines.resize(lsdNFeatures);
428         for( int i=0; i<lsdNFeatures; i++)
429             keylines[i].class_id = i;
430     }
431     if(keylines2.size()>lsdNFeatures)
432     {
433         sort(keylines2.begin(), keylines2.end(),
434             sort_lines_by_response());
435         keylines2.resize(lsdNFeatures);
436         for(int i=0; i<lsdNFeatures; i++)
437             keylines2[i].class_id = i;
438     }
439     cout<<"lbd describe"<<endl;
440     lbd->compute(img, keylines, mLdesc);
441     lbd->compute(image2, keylines2, mLdesc2); // 计算
442     特征线段的描述子
443     BFMatcher* bfm = new BFMatcher(NORM_HAMMING,
444         false);
445     bfm->knnMatch(mLdesc, mLdesc2, lmatches, 2);
446     vector<DMatch> matches;
447     for(size_t i=0; i<lmatches.size(); i++)
448     {
449         const DMatch& bestMatch = lmatches[i][0];
450         const DMatch& betterMatch = lmatches[i
451             ][1];
452         float distanceRatio = bestMatch.distance
453             / betterMatch.distance;
454         if (distanceRatio < 0.75)
455             matches.push_back(bestMatch);
456     }
457
```

```

458     cv::Mat outImg;
459     std::vector<char> mask( lmatches.size(), 1 );
460     drawLineMatches( img, keylines, image2, keylines2, matches, outImg, S
461     Scalar::all( -1 ), mask,
462                 DrawLinesMatchesFlags::
463                 DEFAULT );
464
465     imshow( "Matches", outImg );
466     waitKey();
467
468 }
469

```

470 最后提取的结果是下图展示的。最后大家觉得这种线提取方式，有什么
 优缺点，可以在下方评论或私信我，一起交流讨论。



图 10: LSD 直线提取与匹配结果

471

472 3.2.3 面特征提取与匹配

473 **平面提取** 平面提取一般是针对点云而言的，文章 [11] 总结了一个从点云
 474 到平面分割的基本方式。

- 475 • 首先是从 3D 点云计算 normal map.
- 476 • 第一步聚类，具有相似的 normal 方向时，分在一类 S。
- 477 • 空间点 \times normal 产成距离，根据距离，再把这类等开成不同的 in-
 478 stance。

479 和上述方法 [11] 相比, AHCP [6] 提取平面的效果更好一些, 当然计算也更
480 复杂一些。每三个连续的 2D 点组成一个节点, 就凑成了 a,b,c,d,e,...。通过
481 AHC 方法 (块状聚类), 找到 g 具有最小的 MSE, 在把 g 与左右进行融合。
482 通过将融合信息与阈值相比, 删除误差大的组合。这种每三个点, $\text{step}=1$,
483 就很容易通过 MSE 找出来边界。

484 AHC 是一个比较好的开源平面提取方案。这个方法分为四个阶段, 首
485 先是建立 Initial graph; AHC 分为两个 A, B 两个部分; 粗糙的提取面, 提
486 取精细的面; 如果不需要精细平面的时候, 对于 $640*480$ 的图片, 可以达到
50Hz。



图 11: 平面提取结果。左图是深度图对应的 RGB 图像, 中图是 AHC 方法提取的平面联通区域, 右图是平面分割结果。

487

488 **平面匹配** 不像是点线特征, 可以用一种特征来描述。平面的匹配大多是用
489 用平面的法向量与原点 to 平面的距离来判断两平面的参数是不是一样。

490 3.3 几何特征之间的关系

491 在这个小节中, 我们将特征 SLAM 系统中的约束罗列出来, 主要的目
492 标是介绍如何来形成这样的约束, 至于如何在 SLAM 系统中利用这些约束,
493 我们会在接下来的章节中逐一介绍。

494 我们从一帧图像上提取了点线面信息, 通过探索单帧图像中特征之间的
495 关系, 我们可以获得更多的几何约束。

496 3.3.1 一直线与直线平行-灭点

497 由于 3D 场景投影到图片的过程中, 会产生投影变换, 因此 3D 场景
498 中平行的直线, 在图片中会聚集于一点, 这个点就成为是灭点 (vanishing
499 point)。

500 灭点在 SLAM 中有一些特别的用处，早期一些工作使用他们来做相机
501 参数的计算，现在有些工作利用不同方向上的灭点的正交关系来进行结构化场景中的相机位姿追踪。

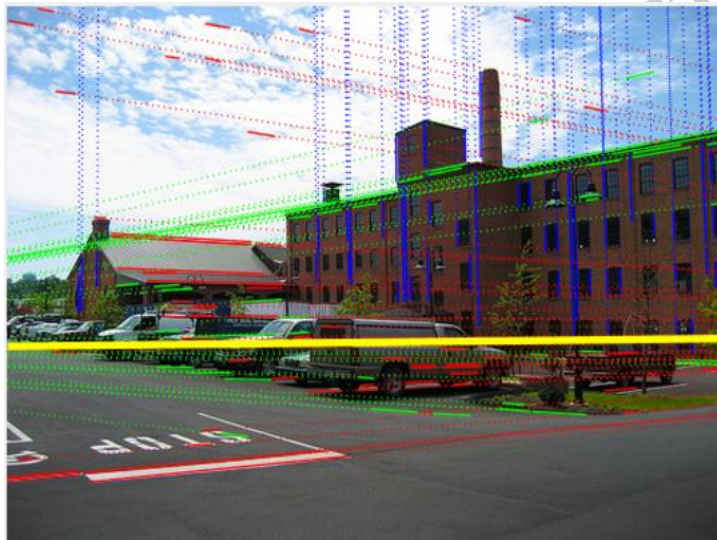


图 12: 灭点与平行直线。红绿蓝颜色代表三个方向上的平行直线，黄色直线是灭线。

502 我们已经知道，2D 图像上的灭点是平行直线的交点。对于 2D 直线的
503 方程 $ax + by + c = 0$ ，我们可以通过联合两条直线的方程求解两条直线的
504 交点 $p = (x, y)$ 。如果另外直线和上述两条直线平行，那么 p 点距离该直线的
505 距离应该会比较小的。由于像素误差的存在，平行直线不会是完全交于一
506 点，因此存在很多方法（RANSAC, J-Linkage 等）[31] 去求出最合理的一个
507 点。
508

509 3.3.2 平面与平面平行-灭线

510 3.3.3 平面与平面垂直-曼哈顿约束

511 如果在一副图像中找到两个垂直的平面，那这两个平面法向量就能够
512 通过叉乘，获得（虚构）与这两个平面垂直的第三个平面，因此构成了一个
513 正交平面系统，我们可以称其为曼哈顿世界（Manhattan World）。

514 与前面特征构成的约束不同，曼哈顿约束可以整体作为一种约束，这种
515 约束在图像与图像之间可以发挥一作用。

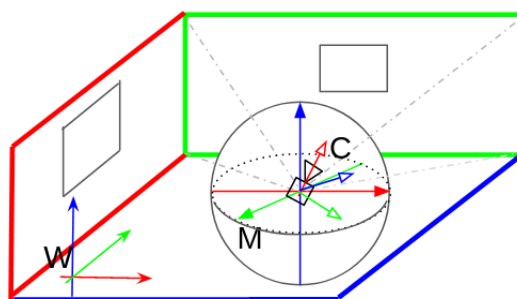


图 13: 曼哈顿世界。不同面之间具有垂直关系。

516 除了图像与图像之间的曼哈顿系统的约束之外，我们还可以将构成曼
517 哈顿系统中的每个平面元素单独做出来平面级（非图像级）的约束。

518 \sqrt{d}

519 3.3.4 点在直线上-共线

520 3.3.5 点线在平面上-共面

521 点线在平面上带来的约束可以被成为共面约束，如图 14所示，空间点
 P_i, P_j 和直线 (n_l 与 n_π 的交线) 都在平面 (法向量 n_π 上)。

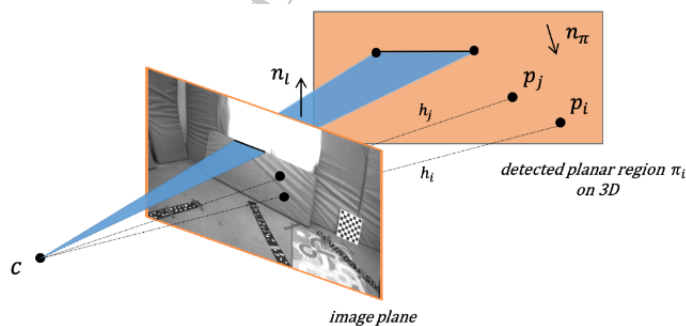


图 14: 点线共面几何 [18]。

522 共面约束产生的方程也是比较直接的，无非就是点和线的描述要满足
523 平面的方程。
524

$$\begin{cases} n_x P_x + n_y P_y + n_z P_z + d = 0 \\ n_l \times n_\pi = \text{dir}(l) \end{cases} \quad (3.7)$$

525 上面的公式分别用简单直接的方式介绍点线构成的共面关系。这里点的表
526 示是简单的 (P_x, P_y, P_z) 方式，直线的表示是用普朗克参数形式 (n_i, d) 的形
527 式表示。值得说明的是，不同参数都是可以相互转化的，因此你可以将你喜
528 欢的参数化形式用在共面约束中。

盐粒-SLAM 中的几何与学习

529 4 深度学习视觉基础

530 随着传感器技术的发展和计算能力的提高，价格便宜的 RGB-D 相机
531 已经在很多场景中已经取得了很好的精度。其中 depth map 可以直接提供
532 稠密的 3D 点云，他们对于场景理解，尤其是近距离场景，起到了重要的作
533 用。但是这种传感器仍旧有很多限制，包括场景的大小、深度图空洞、边缘
534 精度较差等。与此同时，深度学习迅速崛起，不仅在图片分类、物体检测、
535 语义分割等任务上取得了突出的效果，而且在密集重建、场景补全、姿态估
536 等传统的 SLAM 问题上爆出巨大能力。当然很多人还会对深度学习嗤之以
537 鼻，觉得网络缺少解释性。在笔者看来，大可不必这样。这种数据驱动的方法，
538 虽然在很大的部分仍然未得到充分解释，但是他为计算机视觉带来的巨大
539 想象空间是不容忽视的。

540 不管从传感器的价格，还是应用场景（高温、高压、极寒、较远距离）
541 的通用性来说，单目相机永远保持着他独特的优势。当我们无法使用主动光
542 获取深度图时，能否利用这种学习的方法来获得深度图呢？实际上，国内外的
543 研究学者已经在这个问题上取得了很不错的效果，本文就是要对众多的
544 理论、方法进行梳理，分析其中主要方法的逻辑和优缺点，以便后来的读者
545 可以更简单的理解问题、贡献智慧。

546 4.1 经典网络回顾

547 当前深度学习已经在很多任务取得了很多人瞩目的成果，其实早期
548 学习任务是从分类开始的，然后基于学习任务很多网络结构被提出来，这些
549 经典的网络深度的探索了基于卷积神经网络的特征提取工作，如 VGG 系
550 列，ResNet 系列，DenseNet 系列等。这些经典的网络如今被用在各种新
551 的学习任务中，支持网络的优秀表现。因此我们先从分类网络开始，重新理
552 解网络的设计与一些 module 的设计初衷，汲取炼丹经验。

553 **输入：** $32 \times 32 = 1024$ 的手写字体图片。这些手写字体包含 0-9 数字，
554 这个问题也就是相当于是一个 10 个类别的图片。输入一个数字，找到合适
555 的分类。

556 激活函数 ReLU

557 *tanh()* 函数：饱和的非线性函数在梯度下降的时候要比非饱和的非线
558 性函数慢得多，在 AlexNet 中使用 ReLU 函数作为激活函数，达到相同准确
559 率的耗时更少。

LeNet-5 Architecture

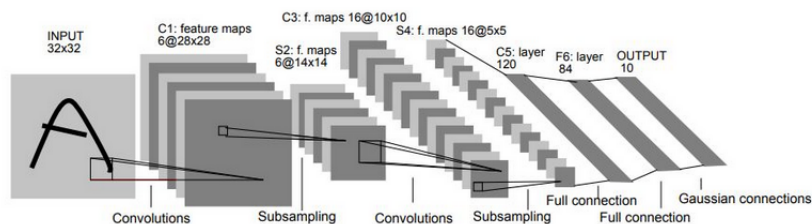


图 15: LeNet-5

卷积	pooling	全连接层
c1: $5 \times 5 * 6$	s2: MaxPool, $s=2$	F6
c3: $5 \times 5 * 6$	s4: MaxPool, $s=2$	output
c5: $5 \times 5 * 120$		

表 1: LeNet-5 网络参数

560 但是：由于 ReLU 函数没有一个有限的值域区间，所以在 ReLU 之后
561 需要进行归一化处理，LRN (**Local Response Normalization**)

562 可以看到 Vgg-16 有 16 个隐藏层（13 个卷积层和 3 个全连接层），还
563 有 Vgg-19 和 Vgg-64 版本。整个网络都使用了同样大小的卷积核（ 3×3 ）和
564 max pooling（ 2×2 ）。

565 论文探讨了一个 7×7 的卷积核与 3 个 3×3 的区别：增加感受野 + 减
566 少参数量。后面很多论文都改用了 3×3 卷积核。

567 GoogleNet 是一个很深的网络，它通过 inception 结构堆叠起来，如下
568 图所示：

569 里面有不同尺度的卷积核，也是第一次看到 1×1 的卷积核，这种操作
570 有调整维度的效果，在 mobilenet 中，也可以看到。

571 Vgg 系列和 GoogleNet 给大家的感觉就是更深的网络，学习能力更强，
572 实际由于梯度消失/爆炸问题，单纯的加深是不能很好提高学习能力的。

573 为了解决因网络深度增加而产生的性能下降问题，residual module 被
574 提出来。**思路是：**不是通过堆叠几层网络直接去拟合一个函数 $H(x)$ ，而是
575 拟合一个残差项，相当于把一个复杂的函数拆成了两项，这样学习起来就变
576 得简单了。

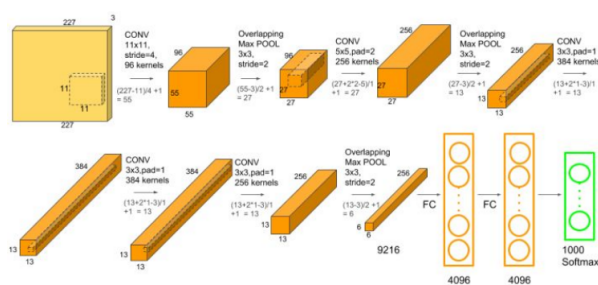


图 16: AlexNet

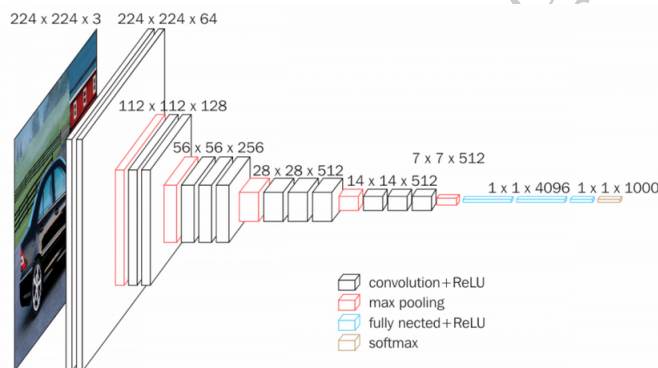


图 17: Caption

577 F(x)+x 能够很容易通过” shortcut connections” 来实现。可以看到 In-
 578 ception 结构更关注网络的宽度，而 ResNet 关注的更多的是网络的深度。
 579 DenseNet 认为上述直接相加的方式会 ** 阻碍信息在网络的流动 **，所以
 580 很简单的实现就是通过一个模块将一段时间内所有的 feature map 直接拼
 581 接起来。

582 Bottleneck Layers

583 不断的拼接 feature map 就会导致后面通道数过多，参数过多。这里
 584 就可以用到 1×1 先降维度，在 3×3 卷积操作，就能维持参数的平稳增加。
 585 BN-ReLU-Conv 1×1 -BN-ReLU-Conv 3×3

586 **Batch Normalization** 在神经网络训练开始前，我们都会对输入数据
 587 做一个归一化处理，其原因是训练过程本质就是为了学习数据分布的过程，
 588 一旦训练数据与测试数据的分布不同，那么网络的泛化能力也大大降低。

589 同样的道理，对于每一个 batch，如果他们的分布各不相同，那么网络

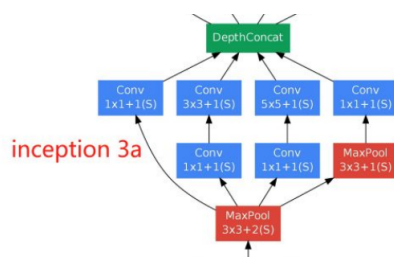


图 18: Vgg-16

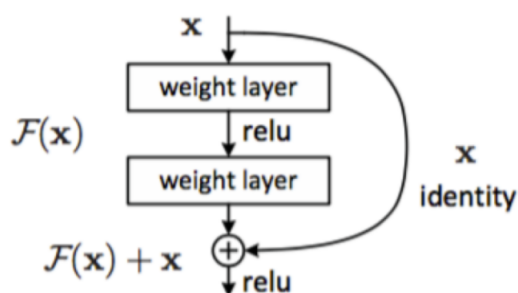


图 19: Caption

590 就要在每次迭代都去学习适应不同的分布，这样将会大大降低网络的训练
 591 速度。所以在网络每一层输入的时候，我们可以考虑插入 BN 层。基于相机的
 592 的 RGB 影像可以提供环境中的结构信息和纹理信息，他们被广泛应用在计
 593 算机视觉和摄影测量的众多问题中。机器人和无人驾驶汽车利用相机传感
 594 器采集图像序列，在未知环境中的姿态估计和环境重建工作。

595 4.2 优化器与学习率

596 前面章节，我们描述了经典的网络结构。网络为什么能完成这些视觉任
 597 务呢？其实一个网络是有很多“未知数/待优化数”够成的，卷积层有卷积
 598 核，全连接层有权重与偏置，激活函数可以调整非线性程...。把输入信息经
 599 过网络前向传播，就可以得到网络预测的结果。如果这个结果离我们的预
 600 期有一些距离 (loss)，我们就把这个距离反向逐层的调整那些网络“未知
 601 数/待优化数”，以减少网络预测结果和实际结果之间的差距。这个反向调整
 602 网络参数的过程就是网络训练的过程。

603 面对这个复杂的优化过程，常见的深度学习框架 (PyTorch/TensorFlow)

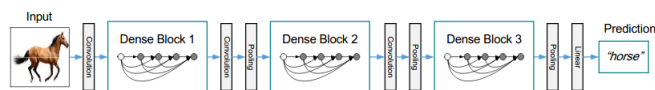


图 20: DenseNet

604 已经帮我们写好了求导的过程，我们需要在网络构建之处选择合适的优化
605 器，并设置相关的参数就可以了。接下来，我们来看几种常见的优化器：

606 **梯度下降 (Gradient Descent)**. 梯度是有所优化问题都会见到的，它是
607 所有参数的偏导构成的一个向量，这个向量指引着极值方向。在 SLAM 中，
608 我们同样有优化问题，不过由于 SLAM 存在稀疏性 [20]，所以我们有特殊
609 的求解器。

610 随机梯度下降方法 (SGD, Stochastic Gradient Descent) 是梯度下降
611 中常见的一种优化器。在参数更新时，该方法不需要计算整个样本的梯度，
612 而是只随机选取其中一个样本计算梯度，因此这个方法训练有很快的训练
613 速度。但是这种简单操作显然也有很多问题，包括：学习率的选择需要经过
614 多次尝试，容易陷入极小值等。

615 为了减少震荡和避免局部最优，我们可以为 SGD 加上 Momentum。在
616 梯度方向不变的维度上，加速参数更新；当梯度方向有所改变时，通过增加
617 动量使参数更新速度变慢。

$$v_{t+1} = \mu * v_t + g_{t+1} \quad (4.1)$$

$$p_{t+1} = p_t - lr * v_{t+1} \quad (4.2)$$

618 这里 p , g , v 以及 μ 分别表示待更新的参数，梯度，下降速度和动量。

```
619
620     if cfg.method == 'sgd':
621         optimizer = torch.optim.SGD(filter(lambda p:
622             p.requires_grad, parameters), lr=cfg.lr,
623             momentum=0.9, weight_decay=cfg.
624             weight_decay)
625
```

626 我们可以在优化的过程中让学习率逐渐下降，例如我们可以让每经历 5 个
627 epoch, 学习率就可以得到原来学习率的 80%。


```
628     if epoch%5 == 0:
629         for p in optimizer.param_groups:
630             p['lr'] *= 0.8
631     print(optimizer.state_dict()['param_groups'][0]['
632           lr'])
633
634
```

635 **学习率自适应优化算法** 为优化器选择合适的学习率是非常困难的问题，因
636 此学界提出很多自适应的优化器，如：Adam 算法、RMSProp 算法，以及
637 AdaDelta。

638 Adam (Adaptive Moment Estimation) [15] 是应用很广泛的方法，它有
639 动量和自适应学习率两个特性，一般全局学习率设置为 $lr = 1e - 4$ 。

```
640     if cfg.method == 'adam':
641         optimizer = torch.optim.Adam(filter(lambda p:
642             p.requires_grad, parameters), lr=cfg.lr,
643             weight_decay=cfg.weight_decay)
644
645     elif cfg.method == 'rmsprop':
646         optimizer = torch.optim.RMSprop(filter(lambda
647             p: p.requires_grad, parameters), lr=cfg.lr
648             , weight_decay=cfg.weight_decay)
649     elif cfg.method == 'adadelta':
650         optimizer = torch.optim.Adadelta(filter(lambda
651             p: p.requires_grad, parameters),
652             lr=cfg.lr, weight_decay=cfg.weight_decay)
653
654
```

655 4.3 常用的度量方法

656 不同的学习问题，有不同的规则去衡量模型的效果。这里我们介绍几
657 种在分类问题中常见的度量方式。我们通过一个例子，来介绍精度 True-
658 Negative, False-Negative, True-Positive, False-Positive, 精确率 (Precision),
659 召回率 (Recall), 准确率 (Accuracy)。

660 对于二分类问题，这些统计概念都很好理解，但是对于多分类问题，有
661 点模糊。在这个三分类问题中，我们分别有 class1, class2, class3, 并获得
网络的预测值，以及数据集的真实值。然后我们就可以统计出以下三张表：

		Predicted		
		class1	class2	class3
Actual	class1	852	83	268
	class2	74	850	273
	class3	216	337	505

图 21: 分类结果统计。

662

	TP	FN	FP
class1	852	83+268	74+216
class2	850	74+273	83+337
class3	505	337+216	268+273

图 22: 统计 TP, FN 和 FP 对应的信息。

	Precision (TP/(TP+FP))	Recall TP/(TP+FN)
Parallel	0.74	0.71
Perpendicular	0.66	0.71
non-relationship	0.48	0.48

图 23: 计算精确率和召回率。

663 5 基于卷积网络的学习任务

664 5.1 点线面的预测

665 5.1.1 预测特征点

666 GCNv2: Efficient Correspondence Prediction for Real-Time SLAM
667 2019 ral. 网络提取特征点, 并提供描述子。实验来看, 比 orb 能更好的
668 应对快速运动。

669 5.1.2 预测特征线

670 Robust Line Segments Matching via Graph Convolution Networks. 刚
671 出来的文章, 建议搭建关注 graph, 可能一大波论文正在袭来。这篇文章已
672 经开源。

673 5.1.3 平面实例分割

674 5.2 物体检测与分割

675 物体检测 (object detection) 可能涉及到物体的识别、分割等。对于识
676 别问题, 也就是判断物体的种类, 是一种分类问题。而物体的分割, 有基于
677 bounding box 的分割, 也有当下流行的实例分割 (instance segmentation)。

678 5.2.1 分类学习基础

679 分类和回归被认为是两种基础的学习任务, 在研究一个学习任务之前,
680 我们都要思考, 这个问题可以被理解成一种什么任务。对于分类问题, 一般
681 是对不同的物体进行概念上的归类, 可以针对整个图片, 也可以针对图片中
682 的局部图像块。

683 手写数字识别是一个基本的分类问题, 对于输入的数字图片, 我们要通
684 过网络来识别数字。我们从这个问题入手, 来理解分类任务。

685 MNIST 数据集是一种识别手写数字常用的数据集, 每张图片的大小是
686 28*28, 并且每个图片都有对应的 ground truth label, 即图片对应的真实数
687 字。这种进点的数据集, 已经被集成在大多数深度学习框架中了, 如 PyTorch,
688 TensorFolow, 大家可以在工程中通过简单的方式直接下载使用。对于这种
689 简单的任务, 我们可以建立一种多层感知机 (Multi-layer perception) 来完成



图 24: Minist 数字手写体数据集

690 工作。下面的代码中可以看出，网络有 3 层全链接层。将图像展开成 1×784
 691 维的向量，经过隐藏层，然后输出 10 个值。

```

692 class MLP(torch.nn.Module):
693     def __init__(self):
694         super(MLP, self).__init__()
695         self.fc1 = torch.nn.Linear(784,512)
696         self.fc2 = torch.nn.Linear(512,128)
697         self.fc3 = torch.nn.Linear(128,10)
698
699
700     def forward(self, input):
701         input = input.view(-1,28*28)
702         dout = torch.nn.functional.relu(self.fc1(
703             input))
704         dout = torch.nn.functional.relu(self.fc2(
705             dout))
706         return self.fc3(dout)
707
708 model = MLP().cuda()
709 print(model)

```

710 基于上述的网络结构，我们就建立了数据之间的传递方式。除了数据传递的
 711 方式，我们还要明确优化器、学习率、损失函数等。

```

712 # loss func and optim
713 optimizer = torch.optim.SGD(model.parameters(), lr
714     =0.01, momentum=0.9)
715 criterion = torch.nn.CrossEntropyLoss().cuda
716
717
718 for x in range(10):

```

```
719     for i, data in enumerate(train_dataset):
720         optimizer.zero_grad()
721         (inputs, labels) = data
722         inputs = torch.autograd.Variable(inputs).
723             cuda()
724         labels = torch.autograd.Variable(labels).
725             cuda()
726
727         outputs = model(inputs)
728
729         loss = criterion(outputs, labels)
730         loss.backward()
731
732         optimizer.step()
733
```

734 MLP 是一种简单的网络结构，对于复杂的任务，可以通过在它的前面
735 增加卷积层来提高网络的学习能力。

736 5.2.2 分类网络中常见策略

737 1. LeNet-5: 每一个卷积核都会形成一个特征图，3 个通道则是每个通道
738 是不同的卷积核，但是最后是将三通道的卷积值相加，最后变成一个通道，
739 所以 $5*5*64$ 的卷积核，感知范围是 $5*5$ ，最后出来是 64 层（个特征图）。
740 每个卷积核都包括 w （权重）和 b （bias 偏置）。LeNet-5 最初用于手写数字
741 识别。

742 2. AlexNet: 5 个卷积层和 3 个全连接层最后输出层是 1000 类的 Soft-
743 max。

744 使用如下创新：

- 745 1) ReLU 非线性激活函数
- 746 2) 多 GPU 训练
- 747 3) 局部响应归一化
- 748 4) 重叠池化

749 3. 降低过拟合的方法：

- 750 1) 数据增强：对图像数据进行变换

751 2) Dropout 随机的关闭一定比例的节点 (神经元)

752 4.VGGnet: 最后输出为 4096 维, 全连接层是 1000 维 (1000 类别)。最

753 后一层是 Softmax 层

754 隐层使用 ReLU

755 5.GoogleNet:2014ILSVRC 分类任务冠军。22 层

756 创新点: 用全局平均池化层取代全连接层, 借鉴了 NIN (network in

757 network) 的做法

758 MLP 网络能够更好的拟合局部特征, 也增强了输入局部的表达能力,

759 NIN 不在分类层前使用全连接, 而是采用全局平均池化。

760 提高深度网络效果的方式是: 增大网络尺寸。但是带来了更多的参数和

761 计算资源需求。

762 改善方案:

763 1) 引入稀疏性

764 2) Hebbin 原则: 两个神经元同步激发, 则他们之间的权重增加, 如果

765 单独激发, 则权重减少。

766 GoogleNet 就是利用 Inception 自动构建非一致结构的神经网络。

767 ResNet: 残差网络 shortcut 一定程度上解决了过深模型梯度发散导致

768 无法训练的问题。

769 DenseNet 从第一层开始每层都作为后面各层的输入。

770 Dual Path Network: 双通道神经网络

771 5.3 深度图与法向量预测

772 5.3.1 监督学习

773 FCN [22] 是一篇无法避开文章, 他的 google scholar 引用已经达到

774 9600 左右。这是一个划时代的工作, 他的出现为深度学习进入像素领域的

775 操作实现了思想解放, 大量的基于像素的语义分割的工作随之而来。回头

776 来看 FCN 提出来的反卷积 (deconvolution) 操作是对特征图的上采样, 实

777 验效果也远后来的方法 [CRF]] 刷新, 但是他对领域的贡献, 已经反应在引

778 用量上了。回到深度预测的话题, David et al. [4] 早在 2014 就提出使用

779 深度学习做这项工作, 他使用了非常基础的网络结构, 如图一所示。框架分

780 为两个部分: coarse net 和 fine net, 前一个网络通过使用较大卷积核子和

781 全链接层试图获得更加全局的结构 (structure) 信息, 第二个网络完全使用

782 卷积操作, 并在把 coarseNet 的结果也作为一种 feature map 输入。这两个

783 网络也是分开训练的。在这个框架基础上, David et al. 提出了更加复杂的
784 框架 [3], 来同时处理 depth prediction, Surface normals(SN) and semantic
785 segmentation. CoarseNet 和 FineNet 经过少许变化之后, 输入 Scale3 以获
786 得更细致的结果。

787 另外, 论文 [3] 在 [4] 深度估计的 Loss 基础上进行修改, 直接使用深度
788 之差 (predicted and GT) 代替对数化的差, 同时考虑差在水平、竖直图像
789 上的梯度。

790 作者表示, 这种设计不仅可以使预测图像接近 GT, 同时还有助于保持
791 局部结构的相似。

792 除此之外, [3] 同样给出了 SN 和 semantic labels 的 LOSS 表示: 对于
793 SN, 作者计算每个像素处的点积 (predicted and GT), 并对整体图像求平
794 均。对于 Semantic 部分, 作者使用 softmax 分类器进行逐项素预测, 不同
795 channels 之间使用 cross-entropy 给出最后预测。数据在监督学习中起着重
796 要的作用, David 的工作是在 NYU-Depth 上训练的。Wang et al.[8] 提出了一
797 种结构来进行 normal estimation. 不同的是, 他们考虑了环境特点 (man-
798 made, Manhattan World) 和中间表述 (房间布局等), 并提出一种全局和
799 局部结合的网络, 并使用融合策略统一他们。[37] 提出了一种在 Physically-
800 Based Rendering 的渲染方法, 并利用 SUNCG [30] 产生很多 synthetic 数
801 据, 作者使用这些数据去训练网络, 并通过 NYU 数据做 finetuning 之后,
802 效果更好。这种渲染方法被说成是以后实时渲染的趋势, 我不太了解, 还是
803 回到数据集和模型。这篇文章的框架是比较简单, 使用 VGG16 提取特征,
804 然后是全卷积网络进行 normal 估计, Loss 则是沿用了 David Eigen 的设
805 计, 训练方法是 RMSprop。同时这个工作也涉及到目标边界检测, 并取得不
806 错的效果。接着这个工作, 作者发表了为 RGB-D 数据进行深度补全 (deep
807 completion) 的论文 [36]。论文首先是利用之前工作 [30], 通过单帧图像获
808 得 SN 和物体边界 (occlusion boundaries), 然后通过一个优化函数, 可以
809 训练 SN, OB 与深度之间的关系。作者指出, 对整体图像进行训练比仅训
810 练洞处的像素效果更好。

811 同样使用 FCN 结构的还有, Iro [17] 的论文, 今天还在 chair 遇见她。
812 这是她和 Christian 完成的工作。Christian 已经去 VGG 做 PostDoc 了, 接
813 触比较少, 但是很喜欢他的安静与真诚。他们的工作是使用 ResNet-50 来
814 提取特征, 并且获得了很好的效果。后来, Keisuke 和 Iro 连接了 Iro-Net
815 and LSD-SLAM 一起做了 CNN-SLAM [32], 虽然网络可以很好的估计深度

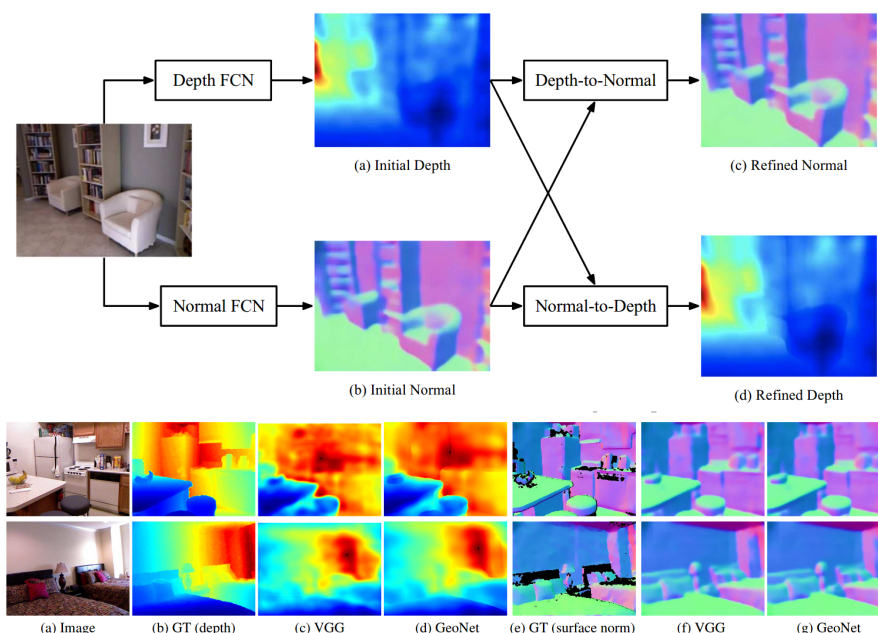


图 25: Geo-Net [27] 的网络框架与训练结果。

816 图，但是计算速度却达不到 real-time，因此 CNN-SLAM 只对 keyframe 做
817 深度估计。

818 上面这些工作确实是比较久远了，网络结构没有那么好，显得比较笨重。
819 我们来看一下比较轻量级的网络：mobileNet V2 [29].mobileNetV2 提供了一个
820 语义分割的框架。论文指出在与其他网络表现相近的情况下，mobileNet
821 V2 的速度和超参的数量有很大的优势。而原因主要是，该网络将标准卷
822 积分解成一个深度卷积和一个点卷积 (1×1)。对于实时性要求的 SLAM 领
823 域，这个网络结构具有较大的意义。论文 [29] 指出，联合 DeepLabv3 和
824 MNetV2 feature map 可以较 MNetV1 减少参数，并保持相似表现。这篇文
825 章没有估计深度图，但是网络结构可以有很好的借鉴意义。Geo-Net [27] 和
826 Deep ordinal [7] 是朋友推荐的两篇，第一篇是港中文 Prof. Jiaya Jia 的团
827 队，第二篇是悉尼大学的 Prof. Dacheng Tao。GeoNet 考虑了 depth map
828 与 normals 之间的紧密关系，建立两个 steam 的网络结构。

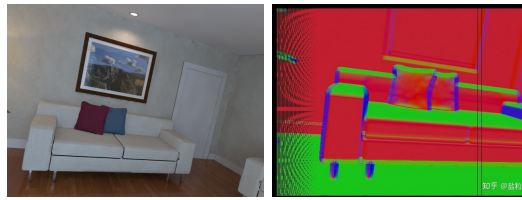


图 26: 制作数据集

829 5.3.2 动手训练一个法向量预测网络

830 在 David Eigen 的 depth prediction, surface normals prediction 的论文
831 网页中, 给出了 NYU depth 的 normal 训练集, 场景也比较丰富, 不过数
832 据量太少了。

833 **制作数据集** 这里是我们自己做的训练集合, 使用 ICL 的深度图来制作自
834 己的法向量数据集。

```

835 void Frame::EstimateSurfaceNormalGradient(const
836     cv::Mat &imDepth, const double &timeStamp, cv::
837     Mat &K)
838 {
839     //ofstream surfaceNormal;
840     //surfaceNormal.open("surfaceN.txt", ios::out
841         | ios::trunc );
842     //PointCloud::Ptr tmp( new PointCloud() );
843     vector<SurfaceNormal> surfaceNormalVector;
844
845     int cellsize=10;
846     cv::Mat mPointCLOud=cv::Mat::zeros(imDepth.
847         rows, imDepth.cols, CV_32FC3); //x,y,z
848
849     for ( int m=0; m<imDepth.rows; m++ )
850     {
851         for ( int n=0; n<imDepth.cols; n++ )
852         {
853

```

```

854         float d = imDepth.ptr<float>(m)[n];
855     if(d<0.001){ continue;}           mPointCloud.at
856         <cv::Vec3f>(m,n)[2]=d;
857
858         mPointCloud.at<cv::Vec3f>(m,n)[0]= (n
859             - K.at<float>(0,2))*d / K.at<
860             float>(0,0);
861
862         mPointCloud.at<cv::Vec3f>(m,n)[1]= (m
863             - K.at<float>(1,2))*d / K.at<
864             float>(1,1);
865     }
866 }
867
868 //读取中心点 和周围的四个点
869 //创建Mat
870 cv::Mat mTangeMask=cv::Mat::zeros(imDepth.
871     rows,imDepth.cols,CV_8U);
872 cv::Mat tangeMaskIntegralTemp;cv::Mat
873     mUTangeMapIntegralTemp;Mat
874     mVTangeMapIntegralTemp;
875
876 cv::Mat tangeMaskIntegral;cv::Mat
877     mUTangeMapIntegral;Mat mVTangeMapIntegral;
878
879 cv::Mat mUTangeMap=cv::Mat::zeros(imDepth.
880     rows,imDepth.cols,CV_32FC3);
881 cv::Mat mVTangeMap=cv::Mat::zeros(imDepth.
882     rows,imDepth.cols,CV_32FC3);
883 for ( int m=1; m<imDepth.rows-1; m++ )
884 {
885     for ( int n=1; n<imDepth.cols-1; n++ )
886     {

```

```

887         //根据深度，获取有效点
888     }    }
889     //求integral
890     integral(mTangeMask, tangeMaskIntegralTemp);
891     integral(mUTangeMap, mUTangeMapIntegralTemp,
892         CV_32FC3);
893
894     integral(mVTangeMap, mVTangeMapIntegralTemp,
895         CV_32FC3);
896     //去掉积分图的padding
897
898     tangeMaskIntegralTemp(cv::Range(1, imDepth.
899         rows-3), cv::Range(1, imDepth.cols-3)).
900         copyTo(tangeMaskIntegral);
901     mUTangeMapIntegralTemp(cv::Range(1, imDepth.
902         rows-3), cv::Range(1, imDepth.cols-3)).
903         copyTo(mUTangeMapIntegral);
904     mVTangeMapIntegralTemp(cv::Range(1, imDepth.
905         rows-3), cv::Range(1, imDepth.cols-3)).
906         copyTo(mVTangeMapIntegral);
907
908     int cellCount=0;
909     for ( int v=cellsize+1; v<imDepth.rows-3;v++
910         )
911     {
912         //处理每个cell
913     }
914     vSurfaceNormal=surfaceNormalVector;
915 }
916

```

917 **完成一个神经网络** 其实我的工作不是一个语义分割的任务，由于目前正在准备 IROS2019，如果后面能接受，我在详细的介绍。训练集和测试集是
918 相同 size 的图片，一个是 RGB，一个是 mask。
919

920 训练集是 RGB 图像和 mask, 把图像尺寸压缩成 160×160 大小。

```
921 transform = transforms.Compose([
922     transforms.ToTensor(),
923     transforms.Normalize(mean=[0.485, 0.456,
924         0.406], std=[0.229, 0.224, 0.225])])
925
926
927 class BagDataset(Dataset):
928
929     def __init__(self, transform=None):
930         self.transform = transform
931
932     def __len__(self):
933         return len(os.listdir('bag_data'))
934
935     def __getitem__(self, idx):
936         img_name = os.listdir('bag_data')[idx]
937         imgA = cv2.imread('bag_data/'+img_name)
938         imgA = cv2.resize(imgA, (160, 160))
939         imgB = cv2.imread('bag_data_msk/'+
940             img_name, 0)
941         imgB = cv2.resize(imgB, (160, 160))
942         imgB = imgB/255
943         imgB = imgB.astype('uint8')
944         imgB = onehot(imgB, 2)
945         imgB = imgB.transpose(2,0,1)
946         imgB = torch.FloatTensor(imgB)
947         #print(imgB.shape)
948         if self.transform:
949             imgA = self.transform(imgA)
950
951         return imgA, imgB
952
953 bag = BagDataset(transform)
```

```

954
955     train_size = int(0.9 * len(bag)) test_size = len(
956     bag) - train_size
957     train_dataset, test_dataset = random_split(bag, [
958         train_size, test_size])
959
960     train_dataloader = DataLoader(train_dataset,
961         batch_size=4, shuffle=True, num_workers=4)
962     test_dataloader = DataLoader(test_dataset,
963         batch_size=4, shuffle=True, num_workers=4)
964

```

965 在开始网络结构之前，还是要回顾一下 feature map 的尺寸是如何变化的，
966 我们需要保持两段 size 的一致。

$$Out_{height} = \frac{Input_{height} - kernel + 2 \times padding}{stride} + 1 \quad (5.1)$$

967 我们还是使用预训练好的 VGG16 做为卷积层，这里要输出五个 size 的
968 feature map, 用来反卷积回来。

```

969
970     class VGG16(nn.Module):
971         def __init__(self):
972             super(VGG16, self).__init__()
973             self.conv1_1 = nn.Conv2d(3, 64, 3,
974                 padding=100)
975             self.relu1_1 = nn.ReLU(inplace=True)
976             self.conv1_2 = nn.Conv2d(64, 64, 3,
977                 padding=1)
978             self.relu1_2 = nn.ReLU(inplace=True)
979             self.pool1 = nn.MaxPool2d(2, stride=2,
980                 ceil_mode=True) # 1/2
981
982     class FCNs(nn.Module):
983
984         def __init__(self, pretrained_net, n_class):
985             super().__init__()

```

```
986         self.n_class = n_class
987     self.pretrained_net = pretrained_net
988     self.relu = nn.ReLU(inplace=True)
989     self.deconv1 = nn.ConvTranspose2d(512,
990         512, kernel_size=3, stride=2, padding
991         =1, dilation=1, output_padding=1)
992     self.bn1 = nn.BatchNorm2d(512)
993     self.deconv2 = nn.ConvTranspose2d(512,
994         256, kernel_size=3, stride=2, padding
995         =1, dilation=1, output_padding=1)
996     self.bn2 = nn.BatchNorm2d(256)
997     self.deconv3 = nn.ConvTranspose2d(256,
998         128, kernel_size=3, stride=2, padding
999         =1, dilation=1, output_padding=1)
1000     self.bn3 = nn.BatchNorm2d(128)
1001     self.deconv4 = nn.ConvTranspose2d(128,
1002         64, kernel_size=3, stride=2, padding
1003         =1, dilation=1, output_padding=1)
1004     self.bn4 = nn.BatchNorm2d(64)
1005     self.deconv5 = nn.ConvTranspose2d(64, 32,
1006         kernel_size=3, stride=2, padding=1,
1007         dilation=1, output_padding=1)
1008     self.bn5 = nn.BatchNorm2d(32)
1009     self.classifier = nn.Conv2d(32, n_class,
1010         kernel_size=1)
1011     # classifier is 1x1 conv, to reduce
1012     # channels from 32 to n_class
1013
1014     def forward(self, x):
1015         output = self.pretrained_net(x)
1016         x5 = output[ 'x5' ]
1017         x4 = output[ 'x4' ]
1018         x3 = output[ 'x3' ]
```

```
1019         x2 = output['x2']
1020     x1 = output['x1']
1021     score = self.bn1(self.relu(self.deconv1(
1022         x5)))
1023     score = score + x4
1024     score = self.bn2(self.relu(self.deconv2(
1025         score)))
1026     score = score + x3
1027     score = self.bn3(self.relu(self.deconv3(
1028         score)))
1029     score = score + x2
1030     score = self.bn4(self.relu(self.deconv4(
1031         score)))
1032     score = score + x1
1033     score = self.bn5(self.relu(self.deconv5(
1034         score)))
1035     score = self.classifier(score)
1036     return score
1037
```

```
1038
1039 def train( epo_num=50, show_vgg_params=False ):
1040
1041     device = torch.device('cuda' if torch.cuda.
1042         is_available() else 'cpu')
1043
1044     vgg_model = VGGNet(requires_grad=True,
1045         show_params=show_vgg_params)
1046     fcn_model = FCNs(pretrained_net=vgg_model,
1047         n_class=2)
1048     fcn_model = fcn_model.to(device)
1049     criterion = nn.BCELoss().to(device)
1050     optimizer = optim.SGD(fcn_model.parameters(),
1051         lr=1e-2, momentum=0.7)
1052
```

```
1053     all_train_iter_loss = []
1054 all_test_iter_loss = []
1055     # start timing
1056     prev_time = datetime.now()
1057     for epo in range(epo_num):
1058
1059         train_loss = 0
1060         fcn_model.train()
1061         for index, (bag, bag_msk) in enumerate(
1062             train_dataloader):
1063
1064             bag = bag.to(device)
1065             bag_msk = bag_msk.to(device)
1066
1067             optimizer.zero_grad()
1068             output = fcn_model(bag)
1069             output = torch.sigmoid(output) #
1070                 output.shape is torch.Size([4, 2,
1071                 160, 160])
1072             loss = criterion(output, bag_msk)
1073             loss.backward()
1074             iter_loss = loss.item()
1075             all_train_iter_loss.append(iter_loss)
1076             train_loss += iter_loss
1077             optimizer.step()
1078
1079             output_np = output.cpu().detach().
1080                 numpy().copy() # output_np.shape =
1081                 (4, 2, 160, 160)
1082             output_np = np.argmax(output_np, axis
1083                 =1)
1084             bag_msk_np = bag_msk.cpu().detach().
1085                 numpy().copy() # bag_msk_np.shape
```



```
1086         = (4, 2, 160, 160)
1087     bag_msk_np = np.argmin(bag_msk_np, axis=1)
1088     if np.mod(index, 15) == 0:
1089         print('epoch {}, {}/{}'.format(
1090             epo, index, len(
1091                 train_dataloader.iter_loss)
1092             )
1093             # plt.subplot(1, 2, 1)
1094     # plt.imshow(np.squeeze(bag_msk_np[0,
1095         ...]), 'gray')
1096     # plt.subplot(1, 2, 2)
1097     # plt.imshow(np.squeeze(output_np[0,
1098         ...]), 'gray')
1099     # plt.pause(0.5)
1100
```

1101 5.3.3 无监督学习

1102 monodepth 和 monodepthv2。

1103 6 基于图网络的学习任务

1104 图网络对于做 SLAM 研究的同学有着特殊的意义，因为我们很早就利
 1105 用 graph 中的节点和边来表示相机（路标点）和观测，相比与 SLAM 优化中
 1106 的图，这里介绍的图网络相同点，也有不同之处。一方面他们处理的数据结
 1107 构是相同的，所谓数据结构就是任意可能的连接关系导致的非结构化数据。
 1108 另一方面，不想 SLAM 中单一的节点形式，图网络中的节点由 deep feature
 1109 形成，并且可以表示任意的内容。由于图网络的资料不像卷积网络那么多，
 1110 因此我们先花些功夫来介绍基础知识，然后再运用到具体的视觉任务中。

1111 6.1 图网络简介

1112 不同于图像数据（像素结构化排列），图（graph）是由节点与边构成的
 1113 非结构化的数据形式，他们形态各异，是一种非结构化的数据。图网络的提
 1114 出就是处理这样的数据结构。PyTorch geometric 是一个基于 pytorch 的图
 网络处理库，里面封装了处理图网络需要用到的基础结构。基于这个图网络



图 27: PyTorch 是基于 PyTorch 的图网络库

1115
 1116 库，我们介绍相关的细节信息。

1117 **1. 输入数据** 单个的图被描述为 `torch_geometric.data.Data`，他有以下属
 1118 性：

1119 我们对上面 data 中的变量进行赋值，就可以创建一个简单的图网络：

```
1120 import torch
1121 from torch_geometric.data import Data
1122
1123 # 边的连接情况： 0->1  1->0  1->2  2->1
1124
```

data.x (节点的特征)	shape: [num_nodes,num_node_features]
data.edge_index (图的连接情况)	shape: [2, num_edges]
data.edge_attr (边的特征)	shape: [num_edges, num_edge_features]
data.y	
data.pos	shape: [num_nodes,num_dimensions]

表 2: 图网络中的参数

```

1125 edge_index = torch.tensor([[0, 1, 1, 2], [1, 0,
1126     2, 1]], dtype=torch.long)
1127 # 这是每个节点的特征x = torch.tensor([[ -1], [0],
1128     [1]], dtype=torch.float)
1129
1130 # 构建Graph的数据
1131 data = Data(x=x, edge_index=edge_index)
1132 >>> Data(edge_index=[2, 4], x=[3, 1])
1133

```

1134 **Data 常用属性** Data 类中已经为大家封装好了一些常用的属性，可以调用相应的函数获得你想要的信息：

```

1136 print(data.keys)
1137 >>> ['x', 'edge_index']
1138
1139 # 节点的特征
1140 print(data['x'])
1141 >>> tensor([[ -1.0],
1142     [0.0],
1143     [1.0]])
1144
1145
1146 for key, item in data:
1147     print("{} found in data".format(key))
1148 >>> x found in data
1149 >>> edge_index found in data

```

```
1150
1151     'edge_attr' in data>>> False
1152
1153     data.num_nodes
1154     >>> 3
1155
1156     data.num_edges
1157     >>> 4
1158
1159     data.num_node_features
1160     >>> 1
1161
1162     data.contains_isolated_nodes()
1163     >>> False
1164
1165     data.contains_self_loops()
1166     >>> False
1167
1168     data.is_directed()
1169     >>> False
1170
1171     # Transfer data object to GPU.
1172     device = torch.device('cuda')
1173     data = data.to(device)
1174
```

如何把点云数据转成图数据 ShapeNet dataset (包含了 17,000 3D 形状的点云和每个点被分成 16 类), 可以通过 `torch_geometric.datasets`

```
1175
1176     from torch_geometric.datasets import ShapeNet
1177
1178     dataset = ShapeNet(root='/tmp/ShapeNet',
1179                       categories=['Airplane'])
1180
```

```
1181 dataset [0]
1182 >>> Data(pos=[2518, 3], y=[2518])
1183
```

1184 如何从点云转成 graph 呢? geometric 也提供了数据转换接口。

```
1185
1186 import torch_geometric.transforms as T
1187 from torch_geometric.datasets import ShapeNet
1188
1189 # 产生最近邻方式产生 graph
1190 dataset = ShapeNet(root='/tmp/ShapeNet',
1191                   categories=['Airplane'],
1192                   pre_transform=T.KNNGraph(k=6)
1193                   )
1194
1195 dataset [0]
1196 >>> Data(edge_index=[2, 15108], pos=[2518, 3], y
1197         =[2518])
1198
1199 # 使用随机增强出来 graph
1200 dataset = ShapeNet(root='/tmp/ShapeNet',
1201                   categories=['Airplane'],
1202                   pre_transform=T.KNNGraph(k=6)
1203                   ,
1204                   transform=T.RandomTranslate
1205                   (0.01))
1206
1207 dataset [0]
1208 >>> Data(edge_index=[2, 15108], pos=[2518, 3], y
1209         =[2518])
1210
```

1211 可以看到 Data 的形状就转过去了。

1212 **构建第一个图网络** 前面在知乎“SLAM 烹饪术”上,我们介绍了基础的图
1213 网络概念,这里我们来实现一个两层 GCN,可以看到,我们直接调用库提
1214 供的 GCN 层就可以了,非常的便捷。在前向传输的过程中,我们为第一层

1215 增加一个 ReLU 和 dropout，来增加网络的线性能力和泛化能力。

```
1216 import torch
1217 import torch.nn.functional as F
1218 from torch_geometric.nn import GCNConv
1219
1220
1221 class Net(torch.nn.Module):
1222     def __init__(self):
1223         super(Net, self).__init__()
1224         self.conv1 = GCNConv(dataset.
1225             num_node_features, 16)
1226         self.conv2 = GCNConv(16, dataset.
1227             num_classes)
1228
1229     def forward(self, data):
1230         # 数据data，获得x，edge_index
1231         x, edge_index = data.x, data.edge_index
1232
1233         x = self.conv1(x, edge_index)
1234         x = F.relu(x)
1235         x = F.dropout(x, training=self.training)
1236         x = self.conv2(x, edge_index)
1237
1238         return F.log_softmax(x, dim=1)
1239
```

1240 写好了网络结构，下面是训练代码的主体，和卷积网络一样，里面包含
1241 了优化器选择与设置等内容。

```
1242 device = torch.device('cuda' if torch.cuda.
1243     is_available() else 'cpu')
1244 model = Net().to(device)
1245 data = dataset[0].to(device)
1246 optimizer = torch.optim.Adam(model.parameters(),
1247     lr=0.01, weight_decay=5e-4)
1248
1249
```

```

1250 model.train()
1251 for epoch in range(200):    optimizer.zero_grad()
1252     out = model(data)
1253     loss = F.nll_loss(out[data.train_mask], data.
1254         y[data.train_mask])
1255     loss.backward()
1256     optimizer.step()
1257
1258 # 模型测试部分
1259 model.eval()
1260 _, pred = model(data).max(dim=1)
1261 correct = int(pred[data.test_mask].eq(data.y[data
1262     .test_mask]).sum().item())
1263 acc = correct / int(data.test_mask.sum())
1264 print('Accuracy: {:.4f}'.format(acc))
1265 >>> Accuracy: 0.8150
1266

```

6.2 如何定义一个图网络层

在本节开讲之前，我们再回顾一下图网络的独特之处

Generalizing the convolution operator to irregular domains is typically expressed as a neighborhood aggregation or message passing scheme.

那么如何实现邻近聚合和信息传递策略，就是本文讲述的内容。

$$x_i^k = r^k(x_i^{k-1}, \rho^k(x_i^{k-1}, x_j^{k-1}, e_{ji})) \quad (6.1)$$

在 PyTorch Geometric 提供的 MessagePassing 基类中，已经实现了消息自动传递机制，作者只需要去定义少数的几个函数，如 message(), update(), ，以及信息聚合的策略，如 aggr="add", aggr="mean" 或 aggr="max"。

1280 **GCN 代码实现** 下面我们来实现一个下面公式表示的 GCN 层，

$$x_i^{(k)} = \sum_{j \in N(i)} \frac{1}{\sqrt{\deg(i)} \cdot \sqrt{\deg(j)}} (\Theta \cdot x_j^{(k-1)}) \quad (6.2)$$

1281 这里 Θ 表示权重矩阵 (weight matrix)。上述公式可以清楚的看出相邻节点
1282 是如何传递到该节点的，

1283 step1. 首先使用权重矩阵来处理 $x_j^{(k-1)}$ ，

1284 step2. 然后消除 degree 的影响，

1285 step3. 最后再相加。

```

1286 import torch
1287
1288 from torch_geometric.nn import MessagePassing
1289 from torch_geometric.utils import add_self_loops,
1290     degree
1291
1292 class GCNConv(MessagePassing): ##继承自
1293     MessagePassing
1294     def __init__(self, in_channels, out_channels)
1295     :
1296         super(GCNConv, self).__init__(aggr='add')
1297         # 使用"Add"策略 step3
1298         self.lin = torch.nn.Linear(in_channels,
1299             out_channels)
1300
1301     def forward(self, x, edge_index):
1302         # x 是输入graph输入, N个节点, 每节点大小
1303         in_channels: 形状 [N, in_channels]
1304         # edge_index 是连接关系, shape [2, E]
1305
1306         # 自我连接, Add self-loops to the
1307         adjacency matrix.
1308         edge_index, _ = add_self_loops(edge_index
1309             , num_nodes=x.size(0))
1310
1311         # Step 1: 线性转换这个 节点特征矩阵.

```



```

1312         x = self.lin(x)
1313     # Step 2: 计算连接度, 并消除影响           row, col
1314         = edge_index
1315         deg = degree(col, x.size(0), dtype=x.
1316             dtype)
1317         deg_inv_sqrt = deg.pow(-0.5)
1318         norm = deg_inv_sqrt[row] * deg_inv_sqrt[
1319             col]
1320
1321         # 搞定 Start propagating messages.
1322         return self.propagate(edge_index, x=x,
1323             norm=norm)
1324
1325     def message(self, x_j, norm):
1326         # x_j has shape [E, out_channels]
1327
1328         # Step 4: Normalize node features.
1329         return norm.view(-1, 1) * x_j
1330

```

1331 当我们调用 `propagate()` 函数时, 系统会自动去调用 `message()`, `aggre-`
1332 `gate()` 和 `update()` 函数。上面的 `message()` 函数会把输出节点特征限制在
1333 $(-1,1)$ 之间。

1334 上述过程就是一个完整的 GCN 信息传递和更新过程。这个简单的 GCN
1335 可以直接在你的网络中调用。

```

1336     conv = GCNConv(16, 32)
1337     x = conv(x, edge_index)
1338
1339

```

1340 **边代码实现** 我们还是先来看边卷积层的数学定义,

$$x_i^{(k)} = \max_{j \in N(i)} h_{\Theta}(x_i^{(k-1)}, x_j^{(k-1)} - x_i^{(k-1)}) \quad (6.3)$$

1341 这里 h_{Θ} 可以看作是一个 MLP。我们同样使用 `messagepassing` 类来写
1342 这个边卷积层。

```
1343 import torch
1344
1345 from torch.nn import Sequential as Seq, Linear,
1346 ReLU
1347 from torch_geometric.nn import MessagePassing
1348
1349 class EdgeConv(MessagePassing):
1350     def __init__(self, in_channels, out_channels)
1351     :
1352         super(EdgeConv, self).__init__(aggr='max
1353         ') # "Max" aggregation策略.
1354         # 定义MLP
1355         self.mlp = Seq(Linear(2 * in_channels,
1356                             out_channels),
1357                         ReLU(),
1358                         Linear(out_channels,
1359                             out_channels))
1360
1361     def forward(self, x, edge_index):
1362         # x has shape [N, in_channels]
1363         # edge_index has shape [2, E]
1364
1365         return self.propagate(edge_index, x=x)
1366
1367     def message(self, x_i, x_j):
1368         # x_i has shape [E, in_channels]
1369         # x_j has shape [E, in_channels]
1370
1371         tmp = torch.cat([x_i, x_j - x_i], dim=1)
1372         # tmp has shape [E, 2 * in_channels]
1373         return self.mlp(tmp)
```

1374 在 message 函数中, 我们使用 mlp () 函数处理获得的 tmp, 这个 tmp
1375 包含该节点本身的 feature, 以及相邻节点的一个关系特征。

1376 7 传感器模型

1377 7.1 相机模型

1378 简单来说，相机模型就是把环境中的东西转变成观察信息的过程。不同的相机种类也有不同的相机模型，我们在这里介绍最简单的小孔成像模型。

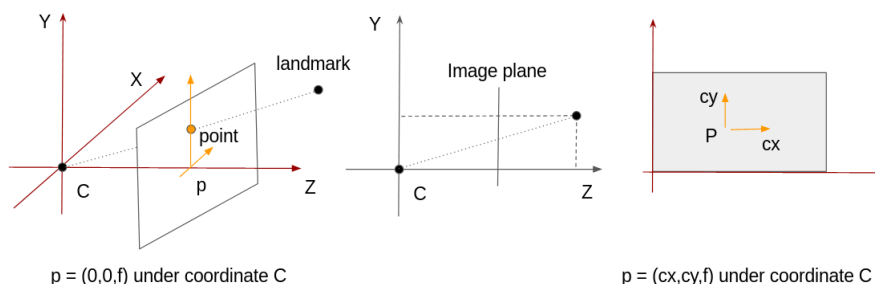


图 28: 小孔成像相机模型。

1379 如图28所示，平面 $z = f$ 被成为是图像平面 (image plane or focal
1380 plane)。把 landmark 点 $(X, Y, Z)^T$ 投影在图像平面上 (图28中):

$$(X, Y, Z)^T \rightarrow (fX/Z, fY/Z)^T \quad (7.1)$$

1382 由于图片中心的偏移, 图片上投影点更加普遍的表达方法是 $(fX/Z + c_x, fY/Z +$
1383 $c_y)$ 。使用齐次坐标表示:

$$(X, Y, Z, 1)^T \rightarrow (fX + Zc_x, fY + Zc_y, Z)^T = \begin{bmatrix} f & s & c_x & 0 \\ 0 & f & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{bmatrix} \quad (7.2)$$

1384 在上述公式中就可以看到相机的内参矩阵 K ，其中 s 是偏斜参数 (skew
1385 parameter) 在大多数开源数据集中都是 0。当 landmark 的表达是在世界坐
1386 标系时，就可以通过刚体变化来转换坐标系。在图29中，我们使用向量关系，
1387 $\vec{w}_c + \vec{c}_x = \vec{w}_x$ ，来描述 landmark 在不同坐标系下的转化关系。

$$X_c = R_{cw}(X_w - C_w) \quad (7.3)$$

1388 X_w 和 C_w 表示世界坐标系下的 landmark 和相机像主点坐标。 $X_w - C_w$ 可
 1389 以表示向量 \vec{CX} 在世界坐标系观点下的形式，这个向量在相机坐标系下的
 观测通过 R_{cw} 来转换。

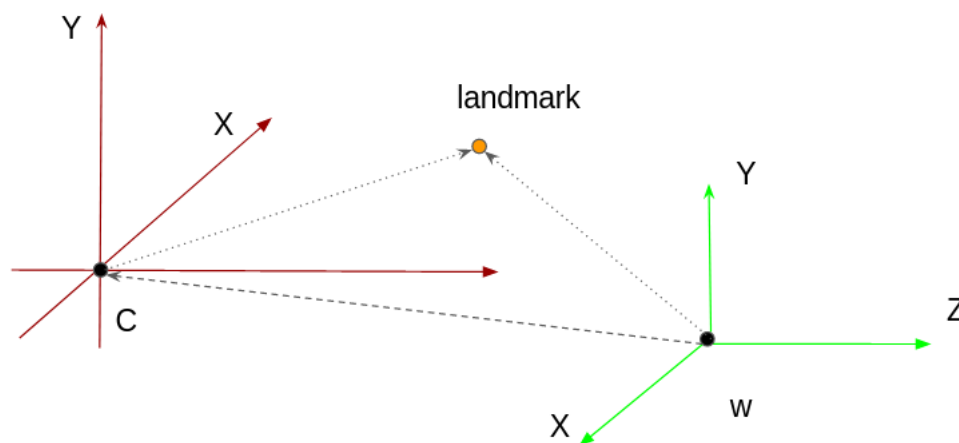


图 29: 路标点在不同坐标系下的表示。

1390

1391

因此，我们利用齐次坐标表示这个过程。

$$X_c = \begin{bmatrix} R & -R_{3 \times 3} C_{w3 \times 1} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad (7.4)$$

1392 通过把上述公式与公式7.2联系起来，我们就可以获得世界坐标系到图像平
 1393 面的过程。

$$x = K X_c = K R [I \quad | \quad -C] X_w \quad (7.5)$$

1394 因此，我们就构建了针孔相机模型 $P = K X_c = K R [I \quad | \quad -C]$ 。

1395 7.2 惯性导航模型

1396 IMU (Inertial Measurement Unit) 的中文名称是惯性测量单元，它由
 1397 加速度计和陀螺仪这两部分组成，我们先抛开坐标系的干扰，可以简单的感
 1398 受到加速度计提供的是加速度信号，而陀螺仪提供的是角速度信号。由此我
 1399 们大概可以理解加速度随着时间积分成速度，然后获得位移运动，角速度则

1400 在时间的作用下产生了旋转运动，而上述两者的结合貌似可以算出来携带
1401 IMU 载体的六自由度位姿。

1402 上述朦胧的感受暗示这个思路貌似可以实现，但是我们都知道，离开坐
1403 标系来谈位姿估计就是要流氓。所以，我们进一步去看这种所谓的加速度和
1404 角速度是在哪里定义的。随着我们揭开这个系统输入信息的面纱，就能十分
1405 自然的过度到传统的捷联导航和预计分的核心算法。

1406 相较于传统的相机传感器，IMU 是一种高帧率设备，在获得每一帧 IMU
1407 数据的时候，加速度计 $a_{i,i+1}$ 获得值是从前一帧到当前帧的加速度信息。不
1408 仅如此，加速度同样受到了重力的影响。对习惯视觉 VO/SLAM 的同学来
1409 说，位姿估计一般是以第一帧图像的单位阵为世界坐标系，关系是发生在
1410 这个假定的世界坐标系中的，因此完全不用考虑和真实世界的关系。但是
1411 IMU 中加速度计采集信息是受到重力影响的，因此这中加速度也被称之为
1412 **包含反向重力加速度的比力**。除了上述问题之外，传感器都存在噪声，噪声
1413 在长时间累积之后，同样或带来巨大的漂移问题。

1414 陀螺仪和加速度计测量模型：

$$\begin{cases} \hat{w}_{wb}^b(t) = w_{wb}^b(t) + b_g(t) + \theta_g(t) \\ \hat{a}^b(t) = R_b^{wT}(a^w - g^w) + b_a(t) + \theta_a(t) \end{cases} \quad (7.6)$$

1415 其中 $b_g()$ 和 $b_a()$ 是随时间变化的偏差， $\theta_g(t)$ 和 $\theta_a(t)$ 是白噪声项。

1416 7.3 激光雷达模型

1417 8 位姿估计的模型

1418 8.1 2D-2D 位姿估计

1419 完成图片的匹配同时,我们就在相邻的 RGB 图像之间获得了同名点的
 1420 信息,一对同名点对应着 3D 空间中同一个路标点。图 30 可以看出,图像的
 1421 匹配结果中还是存在着大量的错误匹配线。如果使用这些错误匹配计算,那
 1422 么将使得后面图像的相对位姿估计出现错误。因此我们必须去除错误的匹
 1423 配对,而常用的方法是 RANSAC 方式。

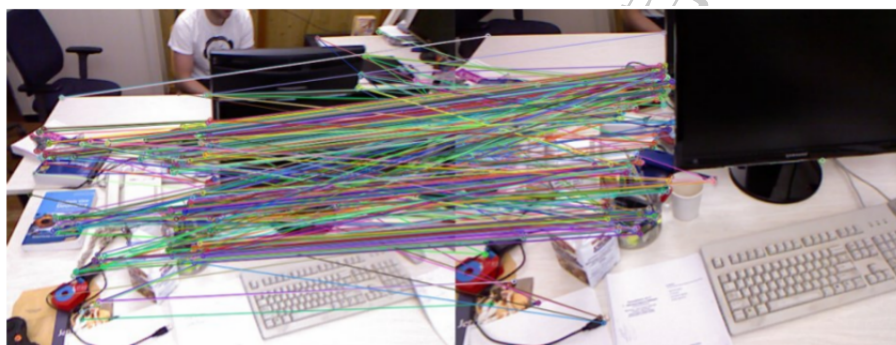


图 30: 彩色通道的特征提取与匹配

1424 对于完成匹配的两张影像之间,我们可以获得两张影像的相对运动情
 1425 况,如下图所示,两张影像的相主点分别是 O_1 和 O_2 ,他们共同观测到空间
 1426 中的点 $X = (x, y, z)$,并且 X 在两张影像上的投影分别是 x_1, x_2 。由于已
 1427 经知道 x_1 与 x_2 的匹配关系以及相机之间存在刚体变换,假设 O_2 通过旋转
 1428 矩阵 R 和平移矩阵 t ,就可以构建起两张影像之间的关系。

1429 **基础矩阵** 相机的内参矩阵是 K ,对于 P_1 和 P_2 ,我们可以计算出 X 在
 1430 两幅影像上的重投影点: $x_1 = KX$ 和 $x_2 = K(RX + t)$ 。因此, $K^{-1}x_2 =$
 1431 $(RK^{-1}x_1 + t)$ 。等式两边同时与 t 左叉乘,则

$$t \times K^{-1}x_2 = t \times RK^{-1}x_1 \quad (8.1)$$

1432 把 $t \times$ 记做 t_{\square} ,则公式变形为 $t_{\square}K^{-1}x_2 = t_{\square}RK^{-1}x_1$,同时两边再点成
 1433 $(K^{-1}x_2)^T$,

$$x_2^T K^{-T} t_{\square} R K^{-1} x_1 = 0 \quad (8.2)$$

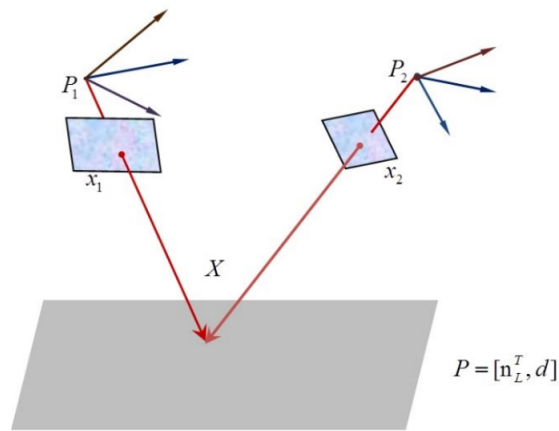


图 31: 空间点属于平面区域

1434 上面的公式构建起了一对同名点的关系。我们把 $K^{-T}t_0RK^{-1}$ 称为基础矩
1435 阵 F , F 矩阵的大小是 3×3 。

$$x_2^T F x_1 = 0 \quad (8.3)$$

1436 通常情况下, F 矩阵中的包含 8 个独立参数, 因此至少八个同名点对就可
1437 完成 F 矩阵的求解, 如果特征点对的数量多于 8, 我们就可以通过奇异值
1438 (SVD) 分解的方法获得最优的旋转矩阵 R 和平移矩阵 t 。

1439 **单应矩阵** 如图31所示, 如果 x_1 和 x_2 对应的路标点落在平面区域, 那么
1440 他们之间的变化关系将会退化成单应矩阵:

$$x_1 = H_\pi x_2 \quad (8.4)$$

1441 通过矩阵变化, 我们构建起了最初的两张影像的位姿变化关系, 同时可
1442 以把两个相机的共同观测作为系统的局部地图维护起来, 当接下来的图像进
1443 入系统时, 可以通过新图像与局部点云之间的转换关系, 进行新的图像帧的
1444 位姿判定的条件。这样就会大大增加系统对于场景的鲁棒性。

1445 8.2 3D-2D 位姿估计

1446 局部地图中的路标点具有世界坐标系, 他们是最近多帧图像所观察的一
1447 个场景。由于机器人运动的连续性, 当前帧和最近多帧具有较多的相同观测

1448 信息, 因此, 我们可以通过局部地图构建一个 3D 点云和 2D 图像之间的关
1449 系。同时确定的关键帧可以为局部地图带来新的路标点, 实现局部地图的更
1450 新。

1451 PnP(Perspective-n-Point) 算法就是一个常用来解算 3D 和 2D 关系的
1452 方法, 这个算法不需要满足上面的对极约束, 可以通过点云和图像点的齐次
1453 坐标之间的关系来获得旋转矩阵 R 和平移矩阵 t 。

1454 对于 PnP 算法的解法, 我们选择是用直接线性变换进行求解。我们首
1455 先来推导 3D 点云和 2D 图像特征点的关联过程, 然后获得求解关系。公
1456 式 8.5 中使用齐次坐标系表示,

$$\begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = K \begin{pmatrix} \hat{X}_c \\ \hat{Y}_c \\ 1 \end{pmatrix} \quad (8.5)$$

1457 这里 $(\hat{X}_c, \hat{Y}_c, 1) = (X_c/Z_c, Y_c/Z_c, 1)$, 其中 (X_c, Y_c, Z_c) 是 3D 点 P 在相机
1458 C 坐标系下的坐标。K 是相机的内参。

1459 由于局部点云中提供的 3D 点的坐标是在世界坐标系下, 他们和当前相
1460 机坐标系下的关系正是由旋转 R_{cw} 和平移矩阵 t_{cw} 构成, cw 则表示从世界
1461 坐标系 w 到相机坐标系 c 。

$$\begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = R_{cw} \begin{bmatrix} X_w \\ Y_w \\ Z_w \end{bmatrix} + t_{cw} \quad (8.6)$$

1462 因此我们就可以构造出当前帧上的 2D 特征点 $(u, v, 1)$ 与世界坐标系下
1463 对应的 3D 路标点 $(X_w, Y_w, Z_w, 1)$ 的关系。这个关系由相机的旋转矩阵、平

1464 移矩阵和相机内参构成, 它可以表示为 $A = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 \\ t_5 & t_6 & t_7 & t_8 \\ t_9 & t_{10} & t_{11} & t_{12} \end{bmatrix} = \begin{bmatrix} A_1^T \\ A_2^T \\ A_3^T \end{bmatrix}$ 。

1465 在相机内参无法获得的情况下, 我们可以通过求解 12 个未知参数获得, 如
1466 果提供了相机的内参矩阵, 未知数的个数将变成 6 个。

$$s \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 \\ t_5 & t_6 & t_7 & t_8 \\ t_9 & t_{10} & t_{11} & t_{12} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} = \begin{pmatrix} A_1^T \\ A_2^T \\ A_3^T \end{pmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad (8.7)$$

1467 其中 s 表示尺度信息, $s = r_7X + r_8Y + r_9Z + t_3 = A_3^T P$, 故,

$$\begin{aligned} A_1^T P - A_3^T P u &= 0 \\ A_2^T P - A_3^T P v &= 0 \end{aligned} \quad (8.8)$$

1468 上面的公式 8.8 可知, 每一对 3D-2D 关系, 可以获得两个方程, 因此我
1469 们至少需要六个关系即可获得 12 个方程, 并求解出来 R, t 的关系。与对极
1470 几何中的 F 矩阵求解类似, 我们在构成超定方程时, 可以使用 SVD 分解获
1471 得最优的 R 和 t 矩阵。

1472 在 OpenCV 中提供了不同的 PnP 求解的函数, 对于 P3P, EPnP, 以
1473 及需要使用 LM 方法优化 3D-2D 之间的最小投影误差等问题, 可以使用
1474 solvePnP 这个函数。这里 objectPoints 是 3D 点云数组, imagePoints 是与
1475 之对应的 2D 图像像素点数组, cameraMatrix 是相机内参矩阵, distCoeffs
1476 表示相机畸变参数 (空输入默认为畸变为零畸变)。

```
1477
1478
1479 C++: bool solvePnP(InputArray objectPoints ,
1480                   InputArray imagePoints , InputArray
1481                   cameraMatrix , InputArray distCoeffs ,
1482                   OutputArray rvec , OutputArray tvec , bool
1483                   useExtrinsicGuess=false , int flags=ITERATIVE )
1484
1485 Python: cv2.solvePnP(objectPoints , imagePoints ,
1486                      cameraMatrix , distCoeffs [, rvec [, tvec [,
1487                      useExtrinsicGuess [, flags ]]]) → retval , rvec ,
1488                      tvec
1489
```

1490 除此之外, OpenCV 还提供基于 RANSAC 的 PnP 函数 solvePnP Ransac, 这
1491 个方法对含有外点 (outliers) 的输入数据更加鲁棒。可以看到除了 solvePnP
1492 方法中的参数以外, 还有 iterationsCount 表示迭代的次数, reprojection-
1493 Error 表示判断外点的阈值, minInliersCount 表示当内点数超过这个数是,
1494 就认为结果足够可靠, 进而结束迭代。

```
1495
1496
1497 C++: void solvePnP Ransac(InputArray objectPoints ,
1498                          InputArray imagePoints , InputArray
```

```

1499     cameraMatrix, InputArray distCoeffs,
1500     OutputArray rvec, OutputArray tvec, bool
1501     useExtrinsicGuess=false, int iterationsCount
1502     =100, float reprojectionError=8.0, int
1503     minInliersCount=100, OutputArray inliers=
1504     noArray(), int flags=ITERATIVE )
1505 Python: cv2.solvePnPRansac(objectPoints,
1506     imagePoints, cameraMatrix, distCoeffs[, rvec[, tvec[, useExtrinsicGuess[,
1507     reprojectionError[, minInliersCount[, inliers[, flags]]]]]])) → rvec, tv
1508

```

1509 此时我们就可以获得当前帧的位姿关系，并将当前帧存储进入轨迹图，并且
 1510 更新局部地图和全局地图。局部地图的维护主要在于控制路标点的数量，在
 1511 相对位姿信息的确定过程中，随着局部地图的增大而增大。因此，我们在局
 1512 部地图中加入新的路标点的同时，必须剔除过时的路标点，以保持局部地图
 1513 的较轻规模。

1514 在有些里程计和 SLAM 系统中，不仅涉及关键帧的选取同时也有关键
 1515 帧的剔除功能，在确定关键帧之后，如果后面遇到更合适的帧，我们可以
 1516 把冗余的关键帧进行踢出，这样可以进一步减少系统的压力。限于篇幅，本
 1517 文对于关键帧的剔除工作不做进一步的介绍。

1518 8.3 3D-3D 位姿估计

1519 8.4 位姿解耦估计方法

1520 前面的方法都是同时计算出来 6D 的位姿中旋转和平移信息。有些研究
 1521 者称，位姿估计中的漂移问题主要是来自旋转漂移，这里具体的原因，感兴
 1522 趣的同学可以去尝试证明一下，应该是很有意义的一件事。

1523 位姿解耦算法一般是将旋转矩阵和平移矩阵分开来解算，并提供一个
 1524 更加精确的旋转矩阵。如何能提供更加准确的旋转矩阵呢？在结构化场景
 1525 中，如室内场景、走廊等人造环境，存在特殊的几何现象：平行与垂直的直
 1526 线和平面。这种环境被归结成曼哈顿世界（Manhattan World）假设和亚特
 1527 兰大世界（Atlanta World）假设。

1528 [12-14] 提出了 RGB-D 相机下的解耦位姿估计的框架，利用深度图计
 1529 算出来平面的法向量，然后使用 mean-shift 聚类方法，找出来环境中正交
 1530 平面，然后去追踪正交平面与正交平面之间的旋转关系。由于每次计算旋

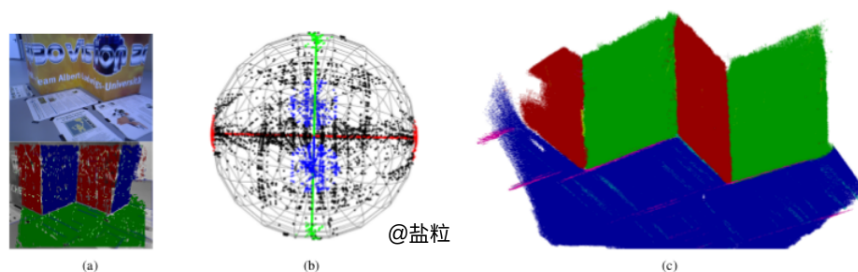


图 32: 曼哈顿假设下的定位与重建。

1531 转, 都是和第一次看到的曼哈顿世界相比较, 因此从理论上说, 这是一种没
 1532 有漂移 (drift-free) 的旋转矩阵估计方式。没有漂移并不是意味着, 没有误
 1533 差, 再每次 rotation 估计中依然会存在误差, 只是这种误差不会积累起来
 1534 了。

1535 8.4.1 曼哈顿和亚特兰大假设

1536 曼哈顿世界是一种存在正交几何信息的一种特殊环境, 如常见的室内
 1537 环境 (如图32 (a) 所示), 墙面之间满足垂直关系, 前面与地板和天花板之间
 1538 同样满足垂直关系。在曼哈顿世界中, 这样的一个正交结构是全局环境中的
 1539 主要结构。与曼哈顿不同, 亚特兰大世界中包含多个曼哈顿世界的场景, 即
 1540 某一个局部满足单一的曼哈顿世界, 然而全局包含很多个不同的曼哈顿世
 界场景, 但是这些场景满足相同的重力方向。

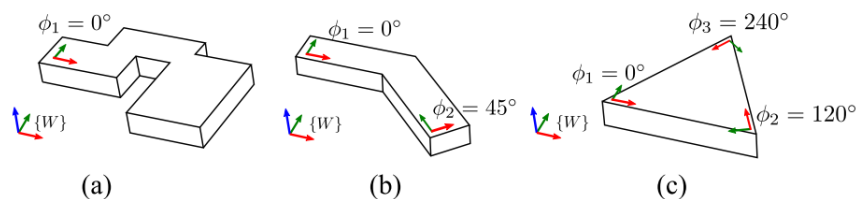


图 33: 环境示意图 [39]。(a) 曼哈顿世界假设; (b) 和 (c) 是亚特兰大世界假设。

1541

1542 **8.4.2 基于特殊环境结构的无漂旋转估计**

1543 前面我们已经介绍了两种常见的环境，在这样的环境运行时，SLAM 系
1544 统可以利用环境中特殊的几何关系来完成相机的追踪和环境重建。

1545 通过什么样的几何元素才能展示这种特殊的环境呢？当然是直线与平
1546 面信息了。对于不同的传感器，我们能使用的信息也受到相应的限制。对
1547 于单目或双面传感器，比较直接的方式就是使用直线和灭点来呈现正交信
1548 息，而 RGB-D 相机，我们就可以同时使用直线和平面信息。下面我们针对
1549 RGB-D 输入，介绍在曼哈顿假设下的无漂旋转矩阵的估计算法。前面章节
1550 已经介绍了直线和平面的提取方法，以及法向量和灭点的计算方法，这里我
1551 们直接使用相关的信息。

1552 首先我们建立一个曼哈顿坐标系 M (如图32所示，并计算当前帧和曼
1553 哈顿坐标系的相对旋转矩阵 R_{C_1M} ，表示从 Manhattan coordinate M 到第
1554 一帧 C_1 。我们可以假设 R_{C_1M} 初始为单位矩阵，然后通过迭代的方式逐渐
1555 优化出来结果。

1556 如图32(b)，首先把 C_1 帧对应的法向量 v_k^s 投影到一个单位高斯球上 [38]
1557 [14]，这个高斯球也同时对应这 R_{C_1M} 的三个方向的正交轴 x, y, z 。然后把
1558 单位法向量投影到每个轴 r_n 对应的切平面上。为了移除噪声，只取切平面
1559 上靠近轴心的圆形区域，并获得 $v_{kn}^{s'}$ 和 $v_{kn}^{d'}$ 。

1560 随后，在切平面上使用 mean-shift 聚类方法，

$$s'_n = \frac{\sum_{in} e^{-c\|m'_{in}\|^2} m'_{in}}{\sum_{in} e^{-c\|m'_{in}\|^2}} \quad (8.9)$$

1561 上述公式中 c 是调整高斯核宽度的参数。完成聚类之后，我们就在切平面上
1562 获得了 r_n 对应的 2D 中心点 s_n 。

$$\hat{r}_n = Qs_n \quad (8.10)$$

1563 这里 $Q = [r_{\text{mod}(n+1,3)}, r_{\text{mod}(n+2,3)}, r_{\text{mod}(n,3)}]$, $\text{mod}()$ 是取模函数。经过迭代，
1564 我们可以获得合理的旋转矩阵 $R_{C_1M} = [\hat{r}_1, \hat{r}_2, \hat{r}_3]$ 。

1565 由于每一帧都是按照上述类似的方式进行旋转估计，这样第 i 帧的旋转
1566 矩阵的估计方式，就不再依赖前一帧或前面的关键帧，而是直接估计 R_{C_iM} ，
1567 因此就实现了漂移的阻断，完成了 SLAM 中常见的漂移问题。估计出来旋
1568 转矩阵之后，我们还是可以使用之前的方式估计相机的平移矩阵。

1569 **8.5 ORB-SLAM2 中的位姿策略**

1570 如上图所示，系统的初始化是在这个函数中完成，初始化根据传感器的不同，分为单目和 stereo 两种方式，其中将 RGB-D 转化出来双目信息。

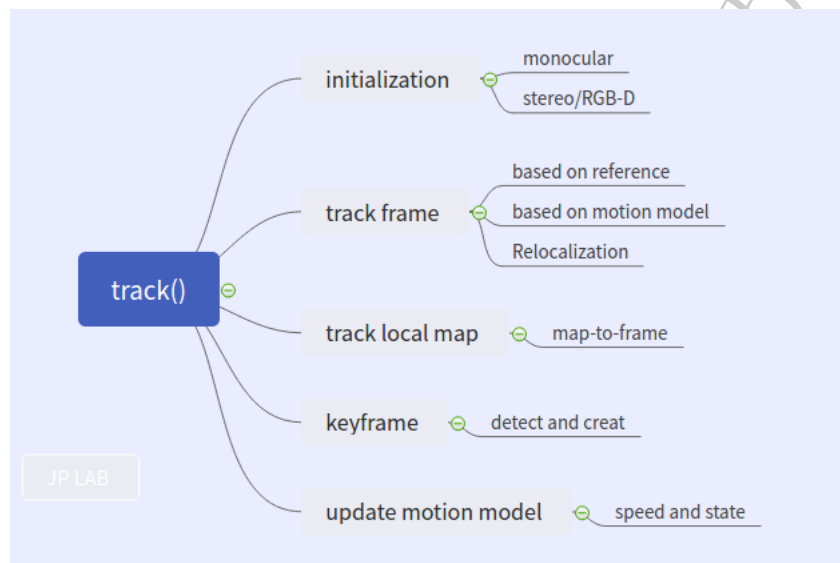


图 34: ORB-SLAM2 中用到的位姿估计模型

1571

1572 MonocularInitialization();

1573 StereoInitialization(); 大家都知道，单目系统为了构建一个稳定的地图，
1574 有比较复杂的初始化过程。双目则是直接启动。后面我们会专门介绍单目初
1575 始化模块。初始化结束之后，会建立关键帧，全局地图。

1576 然后会进行 frame-to-frame, map-to-frame 的位姿估计，当然如果丢失，
1577 则会进行 relocalization。

1578 完成对的当前帧的 track，就会建立起来当填充前帧中 2D 特征对应的
1579 3D landmark。然后去判断是否要加入关键帧。并对追踪状态进行更新。下
1580 面，详细介绍追踪的过程。

1581 其实，ORB-SLAM2 中的追踪策略是非常值得我们借鉴的。首先是给
1582 当前帧一个初始的位姿。通过 track reference 或 track motion model 来获
1583 得用重投影误差去优化当前帧的 Pose。每个结构里面，按照当前系统状态
1584 是否 OK，分成下面三种情况：

1585 OK:

1586 速度值不存在，或者初始化后的一帧，或者重定位后的一帧，使用参考
1587 帧进行追踪：bOK = TrackReferenceKeyFrame();

1588 在追踪正常的情况下，存在运动的速度值，使用运动模型进行追踪：
1589 TrackWithMotionModel();

1590 ! OK:

1591 Relocalization()

```

1592     bool Tracking::TrackReferenceKeyFrame()
1593     {
1594         // Compute Bag of Words vector
1595         cout<<"track reference keyframe"<<endl;
1596         mCurrentFrame.ComputeBoW();
1597         ORBmatcher matcher(0.7,true);
1598         vector<MapPoint*> vpMapPointMatches;
1599         // 找到 vpMapPointMatches, nmatches
1600         int nmatches = matcher.SearchByBoW(
1601             mpReferenceKF, mCurrentFrame,
1602             vpMapPointMatches);
1603
1604         if(nmatches<15)
1605             return false;
1606         cout<<"存在足够匹配点"<<endl;
1607         mCurrentFrame.mvpMapPoints = vpMapPointMatches;
1608         ;
1609         mCurrentFrame.SetPose(mLastFrame.mTcw);
1610
1611         // 通过优化3D-2D的重投影误差来获得位姿
1612         Optimizer::PoseOptimization(&mCurrentFrame);
1613
1614         cout<<" Discard outliers"<<endl;
1615         // 剔除优化后的 outlier 匹配点 (MapPoints)
1616         int nmatchesMap = 0;
1617         for(int i =0; i<mCurrentFrame.N; i++)
1618             {
1619 
```

```

1620         if (mCurrentFrame.mvpMapPoints [ i ])
1621     {
1622         if (mCurrentFrame.mvbOutlier [ i ])
1623         {
1624             MapPoint* pMP = mCurrentFrame.
1625                .mvpMapPoints [ i ];
1626
1627             mCurrentFrame.mvpMapPoints [ i ]=
1628                 static_cast <MapPoint*>(NULL);
1629             mCurrentFrame.mvbOutlier [ i ]= false;
1630             pMP->mbTrackInView = false;
1631             pMP->mnLastFrameSeen =
1632                 mCurrentFrame.mnId;
1633             nmatches--;
1634         }
1635         else if (mCurrentFrame.mvpMapPoints [ i
1636             ]->Observations () > 0)
1637             nmatchesMap++;
1638     }
1639
1640     return nmatchesMap >= 10;
1641 }
1642

```

1643 tracking 成功后接下来就开始调用 TrackLocalMap() 函数来优化位姿
1644 了。

1645 TrackLocalMap() 包括下面几个函数: UpdateLocalMap(): 更新局部地
1646 图中的 keyframes 和 Local Points searchLocalPoints(): 获得局部地图与当
1647 前帧的匹配 PoseOptimization(): 最小化重投影误差来优化位姿

1648 8.6 VIO 中的位姿估计

1649 随着硬件设备的发展, VIO 相机定位系统也被大家熟知, 其全称是
1650 Visual-Inertial Odometry, 即视觉惯性里程计。由于 IMU 的加入, 让系统
1651 在图像质量不佳, 甚至是无图像状态下仍然能够进行一段时间稳定的位姿

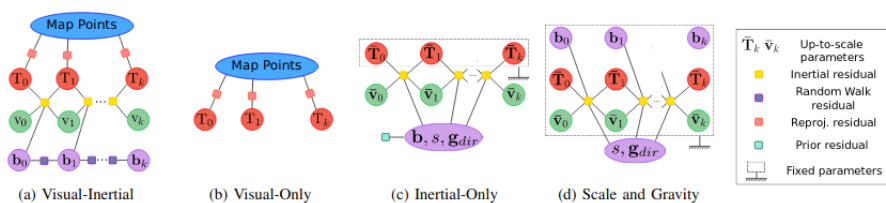


图 35: 不同优化模型中的因子图表示 [from ORBSLAM3 [2]].

1652 估计。在 VIO 中一般有是世界坐标系、相机坐标系和 IMU 坐标系，其中
 1653 IMU 坐标系可以通过一个刚体变化成相机坐标系。VIO 系统初始化的程
 1654 度对后续的定位精度有重要的影响，因此我们首先来看如何将视觉部分和
 1655 IMU 部分通过初始化过程耦合起来。

1656 如图 35所示，不同的追踪系统建立不同的因子图，可以看到 Visual-
 1657 inertial 系统上部分就是常见的视觉图结构，这部分通过重投影误差来连接
 1658 路标点和关键帧，不同的是 (a) 中还有 inertial residual 和 walk residual。
 1659 然后，我们接下来的目标就是构建上述残差项。IMU 速度、偏差和相机位
 1660 姿构成了惯性残差，预计分中的偏差构成了行动残差。这些量就是我们要在
 1661 后续优化过程中计算的量了。

1662 8.6.1 初始化

1663 整个初始化部分分为视觉地图初始化、IMU 中的参数初始化，以及传
 1664 感器耦合三个阶段。基于视觉的地图初始化部分，我们在本章前面的几节已
 1665 经对不同传感器下的位姿估计和地图构建做了详细的介绍。因此，此处我们
 1666 假设地图和刚开始的几帧已经完成了估计，我们直接进行 IMU 的初始化工
 1667 作。

1668 一、角速度偏差的估计。我们在前面的章节介绍的 IMU 的观测方程，
 1669 在角速度和加速度的观测信息包含了偏差 (bias)，因此我们从观测信息
 1670 中把偏差去掉之后，就会更加接近真实值。连续的两帧图像在世界坐标系
 1671 下的旋转矩阵为 $R_{w,i}$ 和 $R_{w,i+1}$ ，因此两帧之间的相对旋转运动就可以获得
 1672 $R_{i,i+1} = R_{w,i}^T * R_{w,i+1}$ 。于此同时，通过预积分我们同样可获得两帧之间的
 1673 旋转运动 ΔR_{imu} ，我们可以通过两个不同模型估计出来的同一个旋转运动

1674 来描述角速度的偏差。

$$R_{i,i+1} = \Delta R_{imu} \begin{bmatrix} 1 \\ \frac{1}{2} J_{b_{wi}}^{\Delta R_{imu}} \Delta b_{wi} \end{bmatrix} \quad (8.11)$$

1675 在初始化过程中，我们有 m 帧图像，因此联立所有相邻图像帧便得到一系列关于角速度偏差的方程，并且可以实现求解。

1677 二、**速度，重力估计**。除了旋转矩阵构成的约束之外，图像和 IMU 之间还有相同的平移矩阵的约束，这个约束用估计每一个图像帧的速度，重力参数 g 和尺度 s 。在 VINS-Mono 中，

$$st_{wb} = st_{wc} + R_{wc} t_{cb} \quad (8.12)$$

1680 利用 IMU 预积分理论，我们可以获得两帧之间的运动速度和位姿表达式，

$$\Delta V_{i,i+1} = R_{wb_i}^T (R_{wb_{i+1}} - R_{wb_i} V_i + g \Delta t) \quad (8.13)$$

$$\Delta t_{i,i+1} = R_{wb_i}^T (t_{wb_{i+1}} - t_{wb_i} - R_{wbk} V_k \Delta t + \frac{1}{2} g \Delta t^2) \quad (8.14)$$

1682 我们把 IMU 和相机之间平移关系（公式 8.12）和 IMU 两帧相对平移联系起来就可以获得：

$$\Delta t_{i,i+1} = R_{wb_i}^T (s(t_{wc_{k+1}} - t_{wc_k}) - (R_{wc_{k+1}} - R_{wbk}) t_{bc} - R_{wbk} V_k \Delta t + \frac{1}{2} g \Delta t^2) \quad (8.15)$$

1684 方程中还剩余的未知量 χ_I 是每个帧位的速度，以及统一的重力和尺度，

$$\chi_I = [v_{b_0}^{b_0}, v_{b_1}^{b_1}, \dots, v_{b_i}^{b_i}, g^{c_0}, s] \quad (8.16)$$

1685 然后，连立公式 8.12 和 8.15，我们构建起来

$$H_{b_{k+1}}^{b_k} \chi_I = z_{b_{k+1}}^{b_k} \quad (8.17)$$

1686 其中 $H_{b_{k+1}}^{b_k}$ 可以表示为：

$$H_{b_{k+1}}^{b_k} = \begin{bmatrix} -I \Delta t_k & 0 & \frac{1}{2} R_{c_0}^{b_k} & R_{c_0}^{b_k} (t_{c_{k+1}}^{c_0} - t_{c_k}^{c_0}) \\ -I & R_{c_0}^{b_k} R_{b_{k+1}}^{c_0} & R_{c_0}^{b_k} \Delta t_k & 0 \end{bmatrix} \quad (8.18)$$

$$z_{b_{k+1}}^{b_k} = \begin{bmatrix} \alpha_{b_{k+1}}^{b_k} - t_c^b + R_{c_0}^{b_k} R_{b_{k+1}}^{c_0} t_c^b \\ \beta_{b_{k+1}}^{b_k} \end{bmatrix}. \quad (8.19)$$

1687 9 地图创建

1688 前面章节介绍了传感器位姿的估计，这里我们介绍 SLAM 另外一个目
1689 标-地图重建。对于未知环境的重建，是机器和环境进行交互的基础。在常
1690 见的 SLAM 系统中，地图主要分为稀疏点云、半稠密点云和稠密点云。

1691 9.1 稀疏地图重建

1692 稀疏点云是指使用一部分特殊几何信息来描述 3D 环境，一般来说这些
1693 稀疏点是梯度变化比较明显的部分。我们先讲解如何利用特征法来进行重
建，然后再讨论直接法。

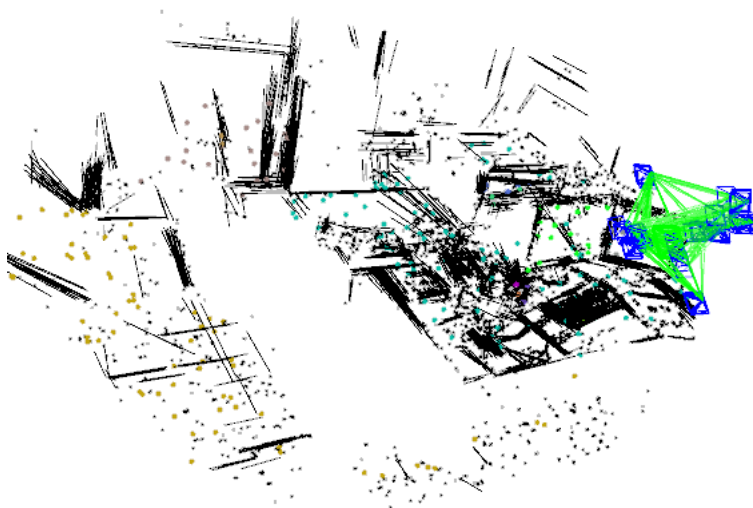


图 36: 点线面构成的稀疏点云。

1694

1695 9.1.1 稀疏点重建

1696 特征法和直接法都可以获得稀疏点特征。特征法是从图片中提取特征
1697 点、线和面特征，并利用特征信息来计算相机位姿和重建环境信息。前面章
1698 节介绍了基本的特征提取与匹配知识，读者可以根据自己的需求选择相应
1699 的特征，因此这里我们直接使用这些特征来完成重建。

1700 **三角化** 如图37所示, 在两幅图像上分别提取同名特征点 x_1 和 x_2 , 而他们
 1701 都是对空间路标点 X 的观测信息。我们的目标就是利用 2D 特征点来恢复
 空间 3D 路标点, 众多的路标点则构成了稀疏点云。现在我们已经获得了同

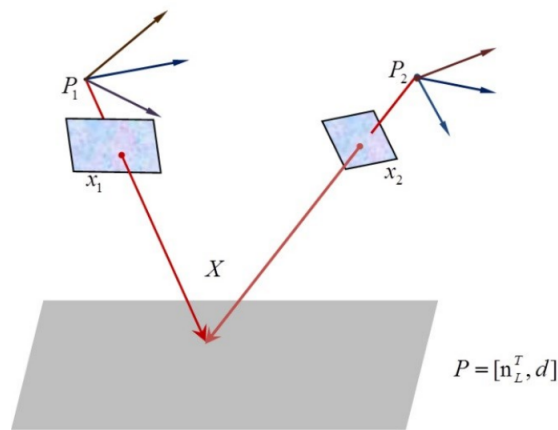


图 37: 特征点的重建

1702
 1703 名点 $x_1 = (u_1, v_1)$ 和 $x_2 = (u_2, v_2)$, 接下来就是通过简单线性三角化方法来
 1704 重建稀疏点 X 。

$$\begin{aligned} x_1 &= P_1 X_w \\ x_2 &= P_2 X_w \end{aligned} \quad (9.1)$$

1705 这里 P_1 和 P_2 是两个相机对应的矩阵。我们构成一个矩阵 A , 并且 $A X_w = 0$ 。

$$A X_w = \begin{bmatrix} x_1 \times P_1 X_w \\ x_2 \times P_2 X_w \end{bmatrix} \quad (9.2)$$

1707 通过对矩阵 A 进行 SVD 分解, 可以获得 3D 点的坐标。

```

1708 cv::Mat A(4,4,CV_32F);
1709 A.row(0) = kp1.pt.x*P1.row(2)-P1.row(0);
1710 A.row(1) = kp1.pt.y*P1.row(2)-P1.row(1);
1711 A.row(2) = kp2.pt.x*P2.row(2)-P2.row(0);
1712 A.row(3) = kp2.pt.y*P2.row(2)-P2.row(1);
1713 cv::Mat u,w,vt;
1714
  
```

```

1715 cv::SVD::compute(A,w,u,vt,cv::SVD::MODIFY_A|cv::
1716 SVD::FULL_UV);
1717 x3D = vt.row(3).t();x3D = x3D.rowRange(0,3)/x3D.
1718 at<float>(3);
1719

```

1720 **视差角** 前文已经从理论上说明任意同名点都可以通过三角化操作实现 3D
 1721 点的重建。事实上，由于特征点和位姿信息都存在误差，当同名点之间视差
 1722 角太小的时候，误差更容易被放大，因此，我们在重建的地图点的时候，尽
 1723 量过滤掉这种情况。

1724 首先，我们获得两个 2D 点在归一化平面上的坐标 $xn1$ 和 $xn2$ ，这两
 1725 个点可以看成是相机到这个点的向量，通过相机到世界坐标系的旋转矩阵
 1726 R_{wc} ，可以获得世界坐标系下的两条光线 $ray1$ 和 $ray2$ ，并计算光线之间的
 余弦值。3D 点在世界坐标系下的光线

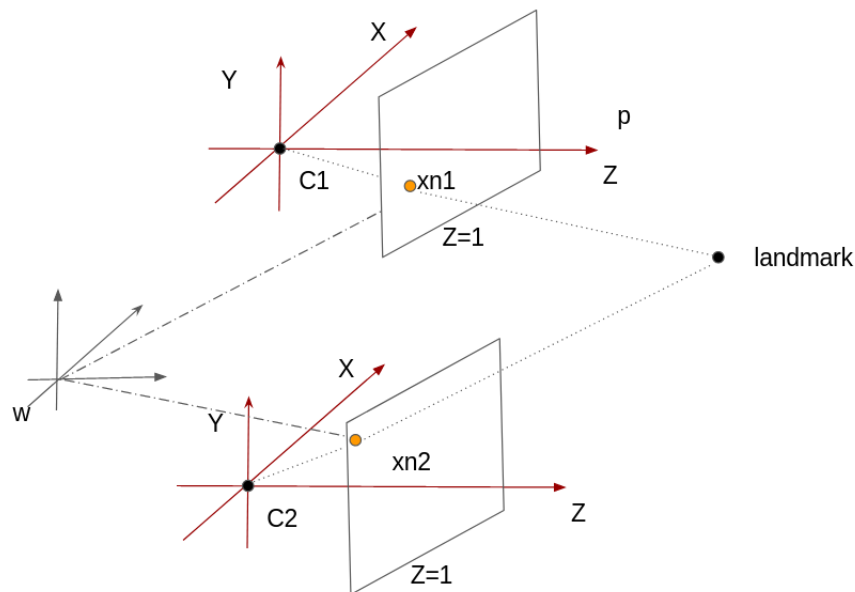


图 38: Caption

```

1727
1728 // 世界坐标系下 观察光线
1729 cv::Mat ray1 = Rwc1*xn1;
1730

```

```

1731 cv::Mat ray2 = Rwc2*xn2;
1732 // 计算在世界坐标系下，两个坐标向量间的余弦值
1733 const float cosParallaxRays = ray1.dot(ray2)
1734 /(cv::norm(ray1)*cv::norm(ray2));
1735
1736 float cosParallaxStereo = cosParallaxRays+1;
1737 //ORB-SLAM使用的判断条件
1738 if(cosParallaxRays<cosParallaxStereo &&
1739     cosParallaxRays>0 && cosParallaxRays<0.9998)
1740

```

1741 9.1.2 线特征的重建。

1742 对于直线的 3D 重建，我们可以使用类似于点的方式。如图39所示，3D
 1743 直线 L 的两个端点 X_s 和 X_e 。利用点的方式，公式9.3，我们可以把两个端
 1744 点投影到图像上，可以获得 $x_s = P_1 X_s$ ，然后点 x_s 在 l_1 和 l_2 上，因此我
 1745 们可以建立

$$\begin{aligned}
 l_1 x_1 &= l_1 P_1 X_s = 0 \\
 l_2 x_2 &= l_2 P_2 X_s = 0
 \end{aligned}
 \tag{9.3}$$

其中 l_1 和 l_2 是 2D 直线的方程，他们可以使用 2D 直线的端点计算出来。

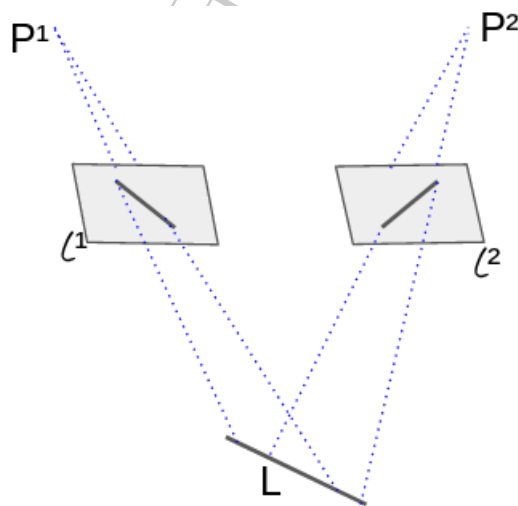


图 39: 特征线的重建

1746
1747 可见公式9.3是不足以求解出来 X_s 的, 我们再把 l_1 (或 l_2) 的端点加入到方
1748 程中, 就像是特征点在公式9.3中的表现一样, $l_{1n} = P_1 X_n$, 这里 n 表示端
1749 点 s 或 e 。

```
1750 // 起始点s
1751 cv::Mat A(4,4,CV_32F);
1752 A.row(0) = k1F1.t()*M1;
1753 A.row(1) = k1F2.t()*M2;
1754 A.row(2) = StartC1.at<float>(0)*Tcw1.row(2)-Tcw1.
1755         row(0);
1756
1757
1758 A.row(3) = StartC1.at<float>(1)*Tcw1.row(2)-Tcw1.
1759         row(1);
1760 cv::Mat w1, u1, vt1;
1761 cv::SVD::compute(A, w1, u1, vt1, cv::SVD::
1762         MODIFY_A| cv::SVD::FULL_UV);
1763
1764 s3D = vt1.row(3).t();
1765
```

1766 **视差角** 与点重建过程中的处理方式相同, 我们计算直线的中点, 并考虑中
1767 点与中点之间的视差角。通过这种方式, 在视差角太小情况下, 我们可以避
1768 免 3D 直线的重建。

1769 9.2 稠密地图重建

1770 上面介绍了稀疏点云中点线重建, 他们都只考虑到一张图片中一小部分
1771 信息。而稠密重建的目标是获得每一个像素的深度信息, 然后利用像素坐标
1772 和深度表示出来。DTAM 是利用 GPU 的一款单目直接法稠密重建系统, 作
1773 者 Newcombe Richard 专注与稠密重建, 并在 2014 年发表了 KinectFusion,
1774 这两篇文章都是 SLAM 历史上的里程碑。这里我们先简单的分析 DTAM
1775 的稠密重建思路, 详细的优化过程推荐阅读论文原文。在直接法中, 我们依
1776 靠光度误差来实现参考帧 I_r 和当前帧 I_m 的位姿估计, 也就是文中的亮度



图 40: 从左到右是深度恢复粗糙到精致的过程。图片来源 DTAM [25]

1777 一致性假设 (brightness constancy assumption)。

$$\rho_r(I_m, u, d) = I_r(u) - I_m(\pi(KT_{mr}\pi^{-1}(u, d))) \quad (9.4)$$

1778 公式中 $\pi()$ 是将 3D 投影成 2D。这样就建立起了 d 和 T_{mr} 的光度误差方
1779 程。

$$C_r(u, d) = \frac{1}{|I(r)|} \sum_m \|\rho_r(I_m, u, d)\| \quad (9.5)$$

1780 将所有关键帧的光度误差联合起来构建优化方程，就可与优化出合适
1781 的地图和位姿信息。

1782 根据地图的表现形式来分，稠密重建可以分为 mesh, SDF, surfel。下面
1783 我们一一介绍。

1784 9.2.1 Mesh

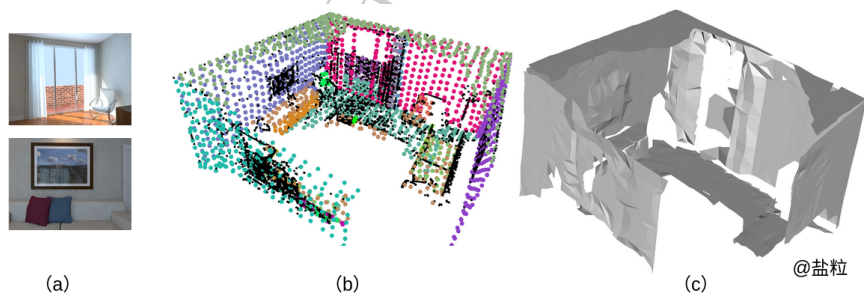


图 41: 基于 mesh 的地图重建。(a) 是场景图片; (b) 是稀疏点云地图; (c) 是 mesh 的密集重建。

1785 mesh 是在点的基础上构网。

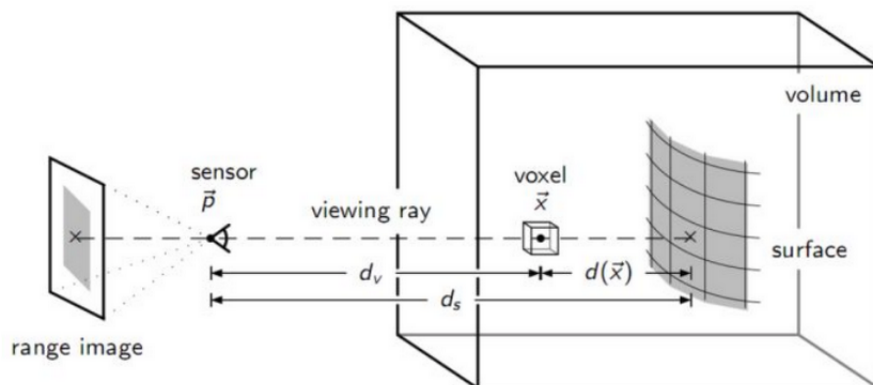


图 42: TSDF 算法示意图

1786 9.2.2 TSDF

1787 TSDF (truncated signed distance function) 是一种密集重建中计算隐
1788 势面常用的表示方式, 其中文全称为“基于截断的带符号距离函数”。TSDF
1789 是 SDF 的一种改进方法, 通过在 SDF 中增加表面到体素的截断距离来实
1790 现计算量的下降。

1791 TSDF 算法把即将重建的空间看成是一个大的立方体区域 (volume), 并
1792 将大的立方体划分成许多小的立方体, 我们称之为体素 (voxel, 可以看出这
1793 这个词就是 volume pixel 的意思)。

1794 在 SLAM 或 SfM 问题中, 我们通过不同的形式来获得空间信息, 如空
1795 间三维点 P_w , 然后将该点初始成一个体素, 很多体素堆放在一起就是我们
1796 的重建目标。

1797 **1. 如何将当前帧融入到模型中** 在讨论这个问题之前, 我们先介绍 TSDF
1798 算法中需要用到哪些信息来定义一个 voxel。

- 1799 • 该 voxel 到最近表面的带符号距离, 记作 $\text{tsdf}(x)$
- 1800 • 体素涉及到更新问题, 以此需要记录每次更新的权重, 记作 w 。

1801 我们抛开计算 3D 点的不同情况, 从比较简单的 RGB-D 数据出发, 深度
1802 信息由第 i 帧深度数据提供, 并获得了相机位姿矩阵 T_{wc} , 然后通过 $T_{wc}^T P_w$
1803 我们可以将物体表面信息 P 在世界坐标系下的坐标转化到相机坐标系下 P_c 。

1804 然后利用相机内参矩阵，找到该点对应在当前帧中的位置信息，一遍求取当
1805 前帧新观测信息到已知表面的距离信息。

1806 具体来说，深度图可以知道当前帧给出的深度值 d ，对因表面上的点距
1807 离当前帧相机原点 O 的距离通过下式计算。

$$d_s = \text{distance}(O - P_c) \quad (9.6)$$

1808 由相机内参矩阵，反投影 P_c 点求深度图像中的对应像素点 v ，我们就可以
1809 获得改像素点对应的深度值 d_v 。因此我们就获得了当前帧观测点的 sdf 距
1810 离。

$$\text{sdf}(p) = d_v - d_s \quad (9.7)$$

1811 与 sdf 不同，TSDF 通过对距离设置一个截断阈值，构造了新的截断函数，

$$\text{tsdf}(p) = \begin{cases} \text{sdf}(d)/|u| & d < u \\ s & \end{cases} \quad (9.8)$$

在计算 sdf 距离之后，我们还需要更新权重 $w(p)$ 的计算公式：

$$w(p) = \cos(\theta)/d_s$$

1812 其中 θ 是投影光线与表面法向量的夹角。

1813 TSDF 的全局更新地图中的体素 x 的信息记为 $TSDF(x)$ 和 $W(x)$ ，这
1814 个点对应在当前帧中的信息是 $\text{tsdf}(x)$ 和 $w(x)$ ，新的帧进入之后，我们对
1815 观测到的信息进行更新：

$$\begin{aligned} TSDF(x) &= \frac{W(x)TSDF(x) + w(x)\text{tsdf}(x)}{W(x) + w(x)} \\ W(x) &= W(x) + w(x) \end{aligned} \quad (9.9)$$

1816 **2. 物体表面** 物体表面在计算 tsdf 中有很重的作用，用来计算新加入信息
1817 的截断距离。一般使用 marching cubes 算法在去寻找 distance 加权和为 0
1818 的等值表面。

1819 9.2.3 Surfel

1820 10 非线性优化

1821 前面的章节介绍了特征的提取匹配、相机姿态求解和全局地图的构建，
1822 但是他们获得的相机姿态和路标点存在误差，并且误差会随着时间的增加
1823 逐渐积累。因此本章介绍光束法平差模型对上述信息进行优化，以获得准确
1824 的相机姿态求解和全局地图，平差模型解算结果会直接影响机器人、无人机
1825 所提供服务的质质量。首先，我们描述基于点线面特征的光束法平差模型，然
1826 后介绍光束法平差模型常用的解算方法。

1827 10.1 光束法平差模型分析

1828 光束法平差 (bundle adjustment, BA) 是三维重建和视觉定位中非常
1829 关键的一步，在获得相机内参数、外参、影像特征点和空间 3D 物点坐标初
1830 始值之后，可以使用 BA 对包括相机的内参和外参在内的上述数据进行整
1831 体的优化，消除误差并获得精确的视觉信息。

1832 光束法平差被认为是摄影测量和计算机视觉中的黄金定律，而数学本
1833 质上来说，光束法平差目标方程模型就是非线性最小二乘模型，并且是关于
1834 相机姿态 (ξ) 与三维路标信息 (*Landmark*) 的最小二乘优化问题。相机姿
1835 态包含平移矩阵 R 和旋转矩阵 t ，而常用的三维路标信息有点 P 、线 L 和
1836 面 Q 。

1837 10.1.1 点特征的误差方程

1838 一组同名像点对应着一个空中的三维点 $P_j = (x, y, z)$ ，通过计算三维
1839 点在影像上的重投影与对应的 2D 特征之间的距离，我们可以获得点的冲投
1840 影误差方程，

$$r_{k,j}^p = p_k - \pi(R_{k,j}P_j + t_{k,j}) \quad (10.1)$$

1841 公式中的 $\pi(\cdot)$ 就是投影过程， f_x 与 f_y 则是相机的内参信息。在优化过程中，
1842 误差方程对于相机位姿的 Jacobian 矩阵可以表示为，

$$\frac{\partial e_{k,j}^p}{\partial \xi} = \begin{bmatrix} \frac{xyf_x}{z^2} & -\frac{z^2+x^2}{z^2}f_x & \frac{yf_x}{z} & -\frac{f_x}{z} & 0 & \frac{xf_x}{z^2} \\ \frac{z^2+y^2}{z^2}f_y & -\frac{xyf_y}{z^2} & -\frac{xf_y}{z} & 0 & -\frac{f_y}{z} & \frac{yf_y}{z^2} \end{bmatrix} \quad (10.2)$$

1843 **10.1.2 线特征的误差方程**

1844 与特征点相比, 直线可以被表示为两个端点 p_{start} 和 p_{end} 。首先我们利
1845 用这两个端点计算直线的方程, 通过单位化处理, 可以显出直线长度带来的
1846 影响。

$$l = \frac{p_{start} \times p_{end}}{\|p_{start}\| \|p_{end}\|} = (a, b, c) \quad (10.3)$$

1847 类似与点的重投影误差, 我们来计算线特征的误差方程。在 2D 图像上
1848 获得直线方程之后, 通过将空间 3D 直线的端点 (P_{start} 和 P_{end}) 重投影回
1849 来。对于每一个端点 $P_x, x \in [start, end]$, 误差方程可以被描述为:

$$e_{k,j}^l = l\pi(R_{k,j}P_x + t_{k,j}) \quad (10.4)$$

1850 同理, 上述误差方程带来的 Jacobian 矩阵为

$$\frac{\partial e_{k,j}^l}{\partial \xi} = \begin{bmatrix} -\frac{f_y^2 z^2 + axyf_x + by^2 f_y}{z^2}, & \frac{az^2 f_x + ax^2 f_x + bxyf_y}{z^2}, & -\frac{ayf_x - bx f_y}{z} \\ \frac{af_x}{z}, & \frac{bf_y}{z}, & -\frac{axf_x + byf_y}{z^2} \end{bmatrix} \quad (10.5)$$

1851 **10.1.3 面特征的误差方程**

1852 平面信息可以在人造环境中随处可见, 常用的平面表示方式是海森表
1853 示 $n_\pi P + d = 0$, 这种方式可以很直观的平面信息, 但是存在过参的问题,
1854 因此我们使用 $q(\pi) = (\phi, \psi, d)$ 来描述平面, 这里 ϕ 和 ψ 是平面法向量的
1855 azimuth 和 elevation, d 表示点到平面的距离。

$$q(\pi) = (\phi = \arctan(\frac{n_y}{n_x}), \psi = \arcsin(n_z), d). \quad (10.6)$$

1856 所以, 我们可以描述两个匹配平面 (π_k 和 π_x) 之间的误差方程,

$$e_{k,\pi_x}^\pi = q(\pi_k) - q(T_{cw}^{-T} \pi_x) \quad (10.7)$$

1857 这里 $T_{cw}^{-T} = (R_{cw}, t_{cw})$ 是从世界坐标系到相机坐标系下的转化矩阵, R_{cw}
1858 用来旋转平面的法向量, t_{cw} 用来处理点到平面的距离。

1859 同理, 我们可以获得面误差方程(10.7)带来的 Jacobian 矩阵,

$$\frac{\partial e_{k,\pi_x}^\pi}{\partial \xi} = \begin{bmatrix} \frac{n_{cx}n_{cz}}{n_{cx}^2 + n_{cy}^2} & \frac{n_{cy}n_{cz}}{n_{cx}^2 + n_{cy}^2} & -1 & 0 & 0 & 0 \\ \frac{n_{cy}}{\sqrt{1-n_{cz}^2}} & \frac{n_{cx}}{\sqrt{1-n_{cz}^2}} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & n_{cx} & n_{cy} & n_{cz} \end{bmatrix} \quad (10.8)$$

1860 **10.1.4 构建光束法平差模型**

$$\begin{aligned}
 t^* = \operatorname{argmin} & \sum_j^M \rho_p \left(e_{k,j}^p T \Lambda_{p_{k,j}} e_{k,j}^p \right) \\
 & + \rho_l \left(e_{k,P_x}^l T \Lambda_{p_{k,P_x}} e_{k,P_x}^l \right) + \rho_\pi \left(e_{k,\pi_x}^\pi T \Lambda_{k,\pi_x} e_{k,\pi_x}^\pi \right)
 \end{aligned} \quad (10.9)$$

1861 在公式10.9中, Λ 是一个协方差矩阵, 同一个传感器观场景进行观测的时候,
 1862 我们认为误差之间是相互独立的, 故通常作为一个对角阵存在, 在实际的
 1863 工作中通常取单位阵。通过对目标方程的解算, 我们最终可以获得优化的结
 1864 果。

1865 **10.2 经典的光束法平差解算方法**

1866 在非线性优化问题的诸多解法中, 牛顿法因为具备二阶收敛速度而得
 1867 到广泛的应用, 但是在迭代求解的过程中, 牛顿法因计算复杂度过高、海森
 1868 矩阵稠密等问题, 无法应用在实际的光束法平差问题中。因此, 我们要有光
 1869 束法平差问题的专门解算方法。

1870 在光束法平差模型解算放法中, 高斯牛顿 (Gauss-Newton, GN) 方法
 1871 通过忽略牛顿方法中产生的高阶导数很好的降低了算法的复杂度, 作为一种
 1872 具有快速收敛性的收敛方法, 被广泛用来解算光束法平差中, 获得精准的相
 1873 机位姿信息和 3D 特征点信息。但是 GN 对初值的精度要求较高, 在初值不
 1874 好的情况下, 将面临发散的危险。列温伯格-马夸尔特 (Leverberg-Marquart,
 1875 LM) 方法通过增加阻尼项, 保持了收敛的鲁棒性, 但是这种方法使用置信
 1876 区间来尝试阻尼的值, 不能准确的获得下降方向, 因此需要多次迭代才能完
 1877 成收敛。这两种方法是最广泛使用的解算方法, 除此之外还有共轭梯度法、
 1878 拟牛顿等解算方法。

1879 **10.2.1 高斯牛顿方法**

1880 高斯牛顿方法从牛顿方法简化而来, 因此在剖析高斯牛顿方法发散原
 1881 因之前, 本文首先介绍牛顿优化方法。通过牛顿法的介绍, 可以更加容易理
 1882 解高斯牛顿发散的数学机理。在数值优化领域, 牛顿法因具有较快的收敛速
 1883 度, 被用来求解无约束优化问题。这个方法的主要思路是在迭代点处进行泰
 1884 勒展开, 使用展开式对目标方程进行近似。因此非线性问题就可以转化为线

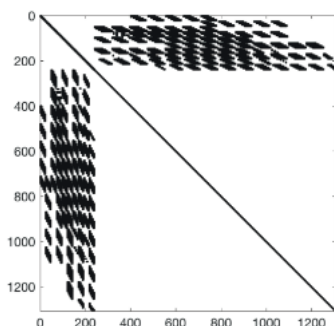


图 43: 海森矩阵的对称稀疏结构。

1885 性问题，并使用求二次模型的极小点的方式来获取新的迭代点。通过多次迭
 1886 代，达到满足的精度为止。对公式10.9进行一阶泰勒展开，这里我们只展开
 1887 公式中特征点的部分

$$z(x) = f(x_{k+1}) + g_k(x - x_{k+1}) + \frac{(x - x_{k+1})^T G(x - x_{k+1})}{2} \quad (10.10)$$

1888 这里 $G = \Delta^2 f(x_{k+1})$, $g_{k+1} = \Delta f(x_{k+1})$ 。对于高斯解法，通过对上述公式
 1889 求导，并令导数为零，就可以获得极值点，

$$\Delta z(x) = \Delta f(f(x_{k+1})) + G(x - x_{k+1}) = 0 \quad (10.11)$$

1890 如果 G^{-1} 是非奇异阵，那么 $x_{k+1} = x_k - G^{-1}g_{k+1}$ ，这样就可以实现迭代更
 1891 新。在之前已经描述过了，牛顿法具有不低于二阶的收敛速度，但是该方法
 1892 在遇到 BA 问题时，求解 Hessian 矩阵工作量巨大，同时不能保证目标函数的
 1893 的 Hessian 矩阵在每个迭代点处都能保持正定。为了利用牛顿法的优势，避
 1894 免不足，出现了 Gauss-Newton 等 BA 解法。

1895 正常情况下 $Hd = J(x)^T F(x)$ ，其中 H 为 Hessian 矩阵， $H = J(x)^T J(x) +$
 1896 $\sum F_i(x)^2$ 。在 Gauss-Newton 方法中，忽略公式的第二项，使用 $J(x)^T J(x)$
 1897 来代替 Hesse 矩阵，求解方程的下降方向。因此，可以看出缺少高阶导数
 1898 信息的高斯牛顿算法在迭代求解过程中必须满足 Hessian 矩阵正定的条件，
 1899 进而要求 Jacobian 矩阵必须是列满秩，不然算法失效。而 Jacobian 矩阵是
 1900 相机与三维特征点之间的表示方式是随着输入参数的变化而变化的。因此
 1901 对 Jacobian 矩阵的要求就变成了对观测初值的要求，这就是高斯牛顿方法
 1902 对初值敏感的根本原因。

Algorithm 1: Gauss-Newton

Input: A vector function $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ $n > m$,
 A measure vector $x \in \mathbb{R}^n$ 初值 P_0
 Output: a vector $P^+ = \operatorname{argmin} \|f(p)\|^2$ $f(p) = x - r(p)$
 $K = 0$ $A = J^T J g = J^T f(p)$ $B_0 = A$

$$\begin{aligned} A d_k^{GN} &= -J_k^T f(x_k) \\ \partial &= 1 \\ x &:= x + d_k^{GN} \end{aligned}$$

图 44: 高斯牛顿法的伪码算法过程

1903 值得说明的是，虽然高斯牛顿法不是非常的鲁棒，但是中对称稀疏的海
 1904 森矩阵形式（如图43），有利的推动了光束法平差模型的使用，因为这种稀
 1905 疏结构可以通过矩阵知识很快速的被计算出来，即使是面对数万维度大小的
 1906 矩阵。

1907 10.2.2 Levenberg-Marquart

1908 有上面内容可知，对于 GN 方法，使用 $J^T J$ 矩阵近似 Hessian 矩阵很
 1909 容易发生奇异，因此，光束法平差问题常常无法解算。为了克服 GN 方法
 1910 的缺点，LM 方法被提出来，并作为应用最为广泛的方法求解 BA 模型。LM
 1911 算法可以表示为公式

$$(J^T J + \lambda I) d_k^{LM} = -J^T g_k \quad (10.12)$$

1912 其中 I 矩阵表示单位阵,LM 算法的原理是在近似 Hessian 矩阵上的对角线
 1913 上加上一个阻尼值，从而近似 Hessian 矩阵的奇异问题。通过观察公式10.12，
 1914 LM 方法是一种结构化的方法，其中包含两个部分： $J^T J$ 与 λI ，可以通过
 1915 调整 λ 的大小，调节 $J^T J$ 与 λ 之间的比重。当 λ 特别大时， $J^T J$ 可以忽
 1916 略，即公式10.12可以变成梯度下降方法公式：

$$d_k^{LM} = -\frac{1}{\lambda} J^T g_k \quad (10.13)$$

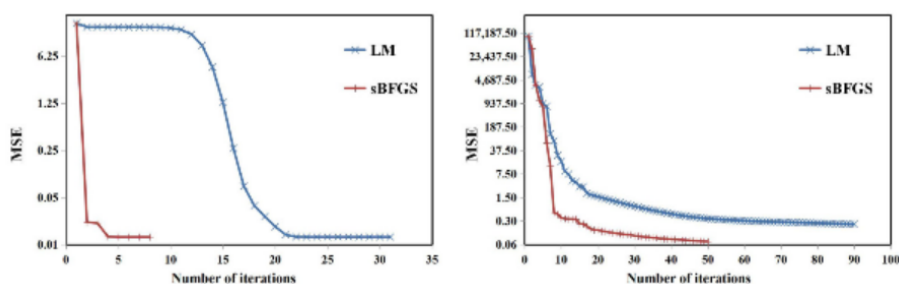


图 45: 两种方法收敛结果。其中红色稀疏结构化的 BFGS 方法的收敛结果, 蓝色为 LM 法收敛结果。

1917 当 λ 特别小时候, λI 可以忽略, 求解方法即变成 GN 方法。因此可以把 LM
 1918 方法看做是最速梯度下降法和高斯牛顿法的结合。使用置信区间来更新 λ ,
 1919 如果当前次迭代可以得到比前一次迭代更小的优化结果, 此时公式 10.12 就
 1920 需要减小 λ 的大小 ($\lambda = \lambda \cdot t, t = \max(\rho, \frac{1}{3})$); 假如不满足条件时, 需要增
 1921 大 λ 值 ($\lambda = \lambda \cdot v, v = 2v$)。

1922 进入的优化环节的初值情况是无法控制的, 与高斯牛顿的初值敏感特
 1923 点一样, LM 方法对于初值较差的情况, LM 同样是被动的调整阻尼值, 不
 1924 能根据当前迭代的实际情况来完成阻尼的更新, 因此 LM 方法常常在需要
 1925 较多次的迭代才能完成。从图 45 的实验结果可以看出, LM 迭代方向并不是
 1926 稳定的, 会不会在迭代初期经历近似平原区域, 还要看光束平差模型中优化
 1927 的输入数据。

1928 10.2.3 稀疏结构化 BFGS 方法

1929 在上述篇幅分析了常见的解算方法: 高斯牛顿方法和 LM 方法, 清楚的
 1930 把这两种方法的优势与缺陷, 及其成因展现在读者眼中。针对上述问题, 稀
 1931 疏结构化 BFGS (sBFGS) [20] 通过探索海森矩阵中高阶信息的保留方式,
 1932 提出一个保持更好下降方向的解算方法。

1933 **BFGS 正定保持特性** BFGS 是一种常用的拟牛顿方法, 由于该方法经过
 1934 秩 2 修正, 因此在迭代求解过程中, 对于输入的正定矩阵, 在满足情况的条

1935 件下，能保持正定性。

$$A_{k+1} = A_k - \frac{A_k r_k r_k^T A_k}{r_k^T A_k r_k} + \frac{z_k z_k^T}{z_k^T r_k} \quad (z_k^T r_k > 0) \quad (10.14)$$

1936 其中， $r_k = X_{k+1} - X_k$ 是两次迭代之间的残差向量。另外， z_k 是雅克比矩
1937 阵的差， $z_k = (J_{k+1}^T - J_k^T) / \|r_{k+1}\|$ 。

1938 **结构化的 BFGS 方法** 在上面的分析中我们明白高斯牛顿方法无法保持矩
1939 阵正定，现在我们为其加入一个正定矩阵，使得在任何时候都能够保持正
1940 定，这样可以解算任意情况的初值。如下式所示：

$$B_{k+1} = J_{k+1}^T J_{k+1} + A_{k+1} \quad (10.15)$$

1941 在本方法中通过 BFGS 算法来更新，具体的证明过程可以参考优化书籍，不
1942 再赘述。

1943 由于 BFGS 对于正定性质的保持能力，假设为正定矩阵，使用 BFGS
1944 算法更新时，只要，并且那么可以保障的正定性，进而为正定矩阵，保证了
1945 下降方向始终是正确的。

$$A_{k+1} = \begin{cases} A_k - \frac{A_k s_k s_k^T A_k}{s_k^T A_k s_k} + \frac{z_k z_k^T}{z_k^T s_k} & \text{if } (z_k^T s_k > v) \\ A_k & \text{otherwise} \end{cases} \quad (10.16)$$

1946 因此上面的式子中，使用 A_{k+1} 近似 Hessian 矩阵中的含有高阶导的部分，
1947 我们通过 BFGS 方法来更好的模拟 Hessian 矩阵。通过公式 10.16 获得的
1948 的 A_{k+1} 是一个密集矩阵，由于 BA 问题涉及到的矩阵的维数非常大，因此
1949 对于密集矩阵的求解将会非常困难。因此本文根据 Hessian 矩阵的低阶导数
1950 信息的稀疏化高阶导数矩阵，因此可以实现稀疏化的 A_{k+1} ，记做 A_{k+1}^s 。

$$B_k^s = \begin{cases} J_{k+1}^T J_{k+1} + A_{k+1}^s & \text{if } (z_k^T s_k > v) \\ J_{k+1}^T J_{k+1} + \|r_k\| I & \text{otherwise} \end{cases} \quad (10.17)$$

1951 如公式 10.18 所示，可以通过上面的稀疏化的 B_k^s 矩阵，在求解下降方
1952 向以后，可以顺利的对光束法平差问题进行解算。与 LM 方法相比，本方法
1953 在迭代点远离真值时，依然能保持超线性收敛速度，随着迭代的进行，当前
1954 迭代点逐渐接近真值，本方法中的高阶导数矩阵的作用就比较微弱，本方法
1955 以高斯牛顿方法进行收敛，同样保持超线性收敛速度。

$$B_k^s d_k^{sBFGS} = -J^T g_k \quad (10.18)$$

¹⁹⁵⁶ 其中 d_k^{sBFGS} 表示 sBFGS 方法下降方向。

盐粒-SLAM 中的几何与学习方法

1957 11 优化中的位姿描述

1958 11.1 四元数

1959 11.2 李群与李代数

1960 12 优化中的特征描述

1961 12.1 点特征-逆深度

1962 12.2 线特征

1963 众所周知，线特征比点特征可以提供更多的约束条件，并且在某些场合
1964 下比点特征要鲁棒。但是如果我们想要把线特征加入到 SLAM 系统中面临
1965 的是参数化和优化求导。这篇文章介绍了贺博的 PL-VIO 中线特征的参数
1966 化和求雅克比。

1967 下面我们就来详细介绍一下空间直线的两种参数化方法：Plücker 参数
1968 化方法，直线正交表示方法。为什么需要两种参数化方法呢？因为空间中的
1969 直线有 4 个自由度，而 Plücker 参数化方法需要使用 6 个参数表示直线，这
1970 样就会导致过参数化，过参数化在优化的时候就需要采用带约束的优化，不
1971 太方便。于是引入了可以用 4 个参数更新直线的正交表示来方便优化。这
1972 两种参数化方法可以很方便的相互转换，所以我们可以 SLAM 系统中同
1973 时使用这两种参数化形式，在初始化和进行空间变换的时候使用 Plücker 坐
1974 标，在优化的时候使用正交表示。下面我们就来详细了解这两种参数化方法
1975 和优化的雅克比求导。

1976 1、Plücker 参数化方法

1977 在 3D 空间中，直线 [公式] 的 Plücker 坐标表示为 [公式]，其中 [公
1978 式] 是直线的方向向量，[公式] 是由直线和光心构成的平面 [公式] 的法向量。
1979 Plücker 坐标是过参数化的，比如 [公式] 和 [公式] 之间存在约束关系 ([公
1980 式])。于是 Plücker 坐标不能直接使用无约束的优化。不过这种使用法向量
1981 和方向向量的表现形式在初始化直线和进行空间变换的时候很方便。所以
1982 在 SLAM 系统中我们可以使用 Plücker 坐标来初始化和变换，具体的初始
1983 化和变换方法如下。

1984 图 1 线特征的 Plücker 坐标 (a) Plücker 线坐标系 (b) 新观测到的线的
1985 初始化当给定从世界坐标系 [公式] 到相机 [公式] 的变换矩阵 [公式]，可以

1986 通过下面的形式将世界坐标系和相机坐标系下的 Plücker 线坐标进行变换：

1987 [公式]

1988 上式中 [公式] 表示的是三维向量的反对称矩阵，[公式] 表示的是从世界
1989 坐标 [公式] 到相机坐标 [公式] 的 Plücker 线坐标的变换矩阵。

1990 [公式] 当我们从不同的两帧相机观测到一个新的线路标的时候，Plücker
1991 坐标初始化方式也很简单。如图 1b 中所示，直线 [公式] 在两帧图像 [公式]
1992 和 [公式] 中表示为两个线段 [公式] 和 [公式]。线段 [公式] 在归一化平面可
1993 以被两个端点表示，[公式] 和 [公式]。再加上坐标原点 [公式]，这三个点
1994 可以确定一个平面：[公式]

1995 [公式]

1996 其中：[公式]

1997 然后对同一个空间直线在两帧相机平面的投影和相机光心可以确定两
1998 个平面 [公式] 和 [公式]，通过两个平面可以确定唯一的空间直线。有了这
1999 两个平面我们可以得到对偶 Plücker 矩阵 [公式]

2000 这里对偶矩阵的内容可以看《多视图几何》[公式]

2001 然后从对偶矩阵中可以提取出 Plücker 坐标 [公式]。法向量 [公式] 和
2002 方向向量 [公式] 都不要是单位向量。

2003 2、正交表示法

2004 由于 3D 空间的直线只有 4 个自由度，使用 Plücker 坐标这种过参数化
2005 的表示形式在优化的时候不方便，所以我们引入四个参数的正交表示 [公式]
2006 。Plücker 坐标和正交表示之间可以很方便的互相转换，我们之后会分别介
2007 绍如何从 Plücker 坐标到正交表示和正交表示怎么变换成 Plücker 坐标。首
2008 先我们需要知道了直线的 Plücker 坐标 [公式]，然后对 [公式] 进行 QR 分
2009 解，得到：

2010 [公式]

2011 分解得到的第一项是正交矩阵 [公式]，是一个旋转矩阵。所表示的是
2012 相机坐标系到直线坐标系的旋转。其中直线坐标系的定义如下：用直线的方
2013 向向量以及直线和光心组成平面的法向量作为坐标的两个轴，再用他们叉
2014 乘得到的向量作为第三个轴，所以

2015 [公式]

2016 其中 [公式] 代表的是相机坐标系到直线坐标系在 [公式]，[公式] 和 [公
2017 式] 轴的旋转角。

2018 到此我们得到了正交表示的第一项，第二项需要做一些小变换。由于将

2019 [公式] 结合之后只有一个自由度，所以我们可以用三角函数矩阵参数化：

2020 [公式]

2021 上式中的 [公式] 是旋转角。由于坐标原点（相机光心）到 3D 直线的距
2022 离是 [公式]，所以 [公式] 包含了距离信息 [公式] 的。根据 [公式] 和 [公式]
2023 的定义可以看出，4 个自由度包括旋转的 3 个自由度和距离的一个自由度。
2024 在优化的时候，我们使用 [公式] 作为空间直线更新的最小表示。

2025 正交表示到 Plücker 坐标之间的变换可以通过下面的方式计算出来：

2026 [公式]

2027 其中 [公式] 代表的是 [公式] 矩阵的第 i 列。虽然 [公式] 和 [公式] 有一
2028 个尺度的差，但是它们表示的是同一条空间直线。

2029 3、直线的观测模型和误差

2030 图 2 空间直线投影到像素平面要想知道线特征的观测模型，我们需要
2031 知道线特征从归一化平面到像素平面的投影内参矩阵 [公式]。如图 2，点 C
2032 和 D 是直线 [公式] 上两点，点 c 和 d 是它们在像素平面上的投影。[公式]
2033 ,K 是相机的内参矩阵。[公式]，[公式]。那么有

2034 [公式]

2035 上式表明，直线的线投影只和法向量有关和方向向量无关。

2036 关于投影的误差，我们不可以直接从两幅图像的线段中得到，因为同
2037 一条直线在不同图像线段的长度和大小都是不一样的。衡量线的投影误差
2038 必须从空间中重投影回当前的图像中才能定义误差。在给定世界坐标系下
2039 的空间直线 [公式] 和正交表示 [公式]，我们首先使用 [公式] 将直线变换到
2040 相机 [公式] 坐标下。然后再将直线投影到成像平面上得到投影线段 [公式]
2041 ，然后我们就得到了线的投影误差。我们将线的投影误差定义为图像中观测
2042 线段的端点到从空间重投影回像素平面的预测直线的距离。

2043 [公式]

2044 其中 [公式] 和 [公式] 是图像中观测到的线段端点，[公式] 是重投影的
2045 预测的直线。

2046 4、误差雅克比推导

2047 如果要优化的话，需要知道误差的雅克比矩阵：

2048 线特征在 VIO 下根据链式求导法则：

2049 [公式]

2050 其中第一项 [公式]，因为

2051 [公式]

2052 所以：
 2053 [公式]
 2054 第二项 [公式]，因为
 2055 [公式]
 2056 所以：
 2057 [公式]
 2058 最后一项矩阵包含两个部分，一个是相机坐标系下线特征对的旋转和
 2059 平移的误差导数，第二个是直线对正交表示的四个参数增量的导数
 2060 第一部分中，
 2061 [公式]
 2062 在 VIO 中，如果要计算线特征的重投影误差，需要将在世界坐标系 [公
 2063 式] 下的线特征变换到 IMU 坐标系 [公式] 下，再用外参数 [公式] 变换到相
 2064 机坐标系 [公式] 下。所以
 2065 [公式]
 2066 其中 [公式]
 2067 [公式] 线特征 [公式] 只优化状态变量中的位移和旋转，所以只需要对
 2068 位移和旋转求导，其他都是零。下面我们来具体分析旋转和位移的求导。首
 2069 先是线特征对旋转的求导：
 2070 [公式]
 2071 然后是线特征对位移的求导：
 2072 [公式]
 2073 第二部分中 [公式]，先解释第一个 [公式]
 2074 [公式]
 2075 所以 [公式]
 2076 然后后面的 [公式] 有两种思路，先介绍第一种：
 2077 [公式]
 2078 其中 [公式] 对 [公式] 和 [公式] 求导，因为 [公式]，所以
 2079 [公式]
 2080 [公式]
 2081 然后是 [公式] 对 [公式] 和 [公式] 对 [公式] 的求导，
 2082 因为 [公式]，所以
 2083 [公式]
 2084 所以，可得

2085 [公式]

2086 12.3 面特征

2087 众所周知，SLAM 中添加线和面这样的高级特征可以增加系统的鲁棒
2088 性，但是在使用面特征的时候有两个问题：面特征怎么去提取？提取出来的
2089 面特征怎么参数化和优化？在这篇文章中，我们会给出几种常用的面特征参
2090 数化方法。如果有错误或者疏漏的地方，欢迎大家批评指正。

2091 Hesse 形式对于面特征来说，大家都知道的一种参数化形式就是 Hesse
2092 形式，用一个平面的单位法向量和平面距离原点的距离来表示：

2093 [公式]

2094 其中 [公式] 表示平面上的一点。在使用 Hesse 形式的时候，我们也可
2095 以很方便的把平面表示在不同坐标系下进行变换。假设 [公式] 是世界坐标
2096 系下的平面方程，[公式] 是本地坐标系下平面方程，[公式] 是本地坐标到世
2097 界坐标的旋转，[公式] 是本地坐标到世界坐标的平移，则有

2098 [公式]

2099 那么我们是不是就可以愉快的使用这种方式呢。答，不可以。因为 SLAM
2100 中有一个优化的问题，优化的时候上面这种 Hesse 形式会导致过参数化
2101 (overparametrization) 的问题。

2102 过参数化问题什么是过参数化问题呢？三维空间中的平面只有 3 个自
2103 由度，但是上面的 Hesse 形式有 4 个参数，这种参数化之后参数的个数大于
2104 实际表示的自由度的表现形式就被成为过参数化 (overparametrization)。
2105 过参数化的形式会引发什么样的问题呢？对于高斯-牛顿优化问题来说，如
2106 果是过参数化的表现形式，优化过程中计算出的 Hessian 矩阵会不满秩，也
2107 就不能求逆。所以对于过参数化的问题，一种解决办法就是列文伯格-马夸
2108 尔特算法 (L-M) 提供的正则化 (regularization)。然而，由于代价函数在某
2109 些方向上是平面，因而存在无穷多的解，会使得算法的收敛速度变为梯度下
2110 降法的线性收敛速度，而不是高斯-牛顿法的二次收敛速度，因此会导致收
2111 敛的速度变慢。除此之外，过参数化还不可以使用信赖区域法。所以下面我
2112 们就来介绍几种平面参数化的方法。

2113 球坐标形式其实平面的 Hesse 参数化形式可以被分成两个部分，一个
2114 单位法向量部分和一个距离部分。平面的 Hesse 形式的过参数化就是因为
2115 单位法向量部分有 3 个参数，但是实际只有两个自由度导致的。那么我们就
2116 可以从这个法向量着手解决过参的问题。单位法向量可以被看成是一个单

2117 位圆球上的一点。那么我们就可以用两个角度 [公式] 来参数化这个点，从
2118 而表示出这个单位法向量。那么我们就可以把 Hesse 形式变换成如下的形
2119 式：[公式]

2120 图 1 半球坐标系 g_{2o} 中也实现了这种面特征参数化形式的节点。这
2121 种方式我们可以便捷的计算出 Hesse 形式对其的雅克比：

2122 [公式]

2123 但是这种形式的参数化虽然是最小参数化形式，但是在优化的时候会
2124 产生奇异性。当 [公式] 的时候，上面的雅克比的秩会从 3 降到 2。

2125 切平面形式那么上面的球坐标的形式会产生奇异性，还有没有更好的
2126 形式呢，也是有的。比如 GTSAM 里面表示面的方式用切平面的方式来更
2127 新单位法向量，同样也可以被用来优化单位法向量。和上面类似，在单位球
2128 平面上找到该法向量对应的点的切平面。然后就可以任意找到一对正交单
2129 位法向量，这两个法向量乘以单位就可以表示这个平面上的任意一点。

2130 [公式]

2131 法向量的优化参数就可以被表示为切平面上这两个正交向量的大小系
2132 数。然后优化的时候求雅克比的话也很方便。不过这种参数化形式一般只用
2133 来对法向量方向进行微调。

2134 图 2 切平面最近点 (Closest Point, CP) 这是黄国权黄老师团队在论
2135 文 LIPS: LiDAR-Inertial 3D Plane SLAM Patrick 中提出的一种平面参数
2136 化方法。这种方式是用在平面上距离当前表示坐标系原点最近的 3D 坐标点
2137 来表示平面。这种表示方式有两个优点：1. 是最小化参数形式；2. 具有真实
2138 的物理意义。

2139 [公式]

2140 但是需要说明的是，这种参数化方式在平面参数 d 为 0 的时候会有奇
2141 异性。所以一般情况需要将平面参数化在检测到平面的坐标系之下，这样
2142 就可以消除奇异性。

2143 图 3 最近点单位四元数这是 CMU 的 Michael Kaess 教授在论文 Si-
2144 multaneous Localization and Mapping with Infinite Planes 中提出的方法。
2145 对于单位四元数来说，这是一种采用了 4 个参数，但是只有 3 个自由度，同
2146 时只需要 3 个参数去更新的参数化形式。那么这种形式很巧合的和平面有
2147 4 个参数，但是只有 3 个自由度在数学上高度吻合。那么如果我们可以把平
2148 面的 4 个参数投影到单位四元数平面上，那么就可以用单位四元数去优化
2149 平面了。

2150 所以我们只需要把 Hesse 形式进行一个小小的变形就可以得到单位四
2151 元数化的平面参数化形式啦。

2152 [公式]

2153 对于这种参数化形式，主要是用来优化不同坐标系下的平面的变换的
2154 误差。只能整个平面的方程一起更新。同时这种参数化形式的更新和其他方
2155 式的参数化不太一样，四元数的参数化更新和旋转矩阵的参数更新是一样
2156 的，是对应到李代数去更新的。所以这种参数化形式有个令人诟病的地方就
2157 是：这种参数化形式的几何意义并不明确。

2158 利用四元数去优化平面方程构建误差一般采用以下形式：

2159 [公式]

2160 退化二次曲面最后给大家介绍一种特殊的表示方法。这是 g2o 的作者
2161 最近在论文 Unified Representation and Registration of Heterogeneous Sets
2162 of Geometric Primitives 中提出的一种可以同时用一种参数化形式同时表示
2163 点线面的参数化形式。大神这么做的意义是想把所有特征都统一起来，也可
2164 以方便的构建出不同特征之间的约束（比如，点到面，点到线，线到面）。本
2165 人认为这应该也是未来多特征 SLAM 的一个大的研究方向吧。一个二次曲
2166 面可以通过下面这种方式参数化，其中 [公式] 是二次曲面上的点，[公式] 是
2167 二次曲面的原点，[公式] 是一个 3×3 的对称矩阵。

2168 [公式]

2169 对 [公式] 做一个特征值分解可以得到 [公式]。其中 [公式] 的每一列都
2170 是二次曲面坐标系的一个轴，[公式] 是包含了 [公式] 的特征值的一个对角
2171 阵。

2172 这种表示方法是把一个点看成是一个半径为 0 的一个球；一条直线看
2173 成是一个半径为 0 的一个圆筒；一个平面看成是两个匹配上的平面。所以
2174 它们对应的分解得到的 [公式] 的形式分别如下表所示。（1 代表该特征值和
2175 其他特征值相等，而 0 代表特征值为 0）

2176 所以对于任意一个点线面的特征实体 [公式]，可以用下面 3 个部分来
2177 参数化：

2178 [公式]

2179 那么两个特征之间的变换也可以从旋转和平移中简单的得到。

2180 [公式]

2181 一些其他形式最后，还有很多种方法没有被提到，比如 VINS on wheels
2182 中的一个 2 自由度的四元数加距离，Probabilistic plane fitting in 3d and an

2183 application to robotic mapping 和 Probabilistic RGB-D odometry based on
2184 points, lines and planes under depth uncertainty 中分别对 Hesse 形式的四
2185 个参数除以一个参数导致可以使用 3 个参数最小表示平面。就不在此赘述
2186 了。因为它们都有一定的限制条件，并不大众化。

盐粒-SLAM 中的几何与学习

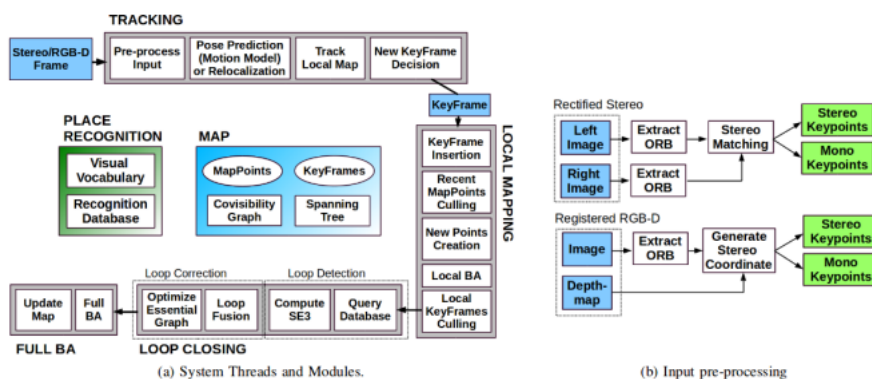


图 46: ORB-SLAM2 系统框图。[24]

2187 13 经典的纯几何 SLAM

2188 13.1 ORBSLAM (V1 and V2)

2189 ORB-SLAM2 [24] 是由 Raúl Mur-Artal 博士等人提出的一个完备的基
 2190 于特征点的纯视觉实时 SLAM 系统, 适用于单目、双目和 RGB-D 相机, 发
 2191 表于 2017 年。作者在 2015 年第一次将单目版本 (ORB-SLAM [23]) 发表
 2192 在 IEEE Transactions on Robotics (ToR) 期刊上, 这是一个基于纯特征点
 2193 的单目实时 SLAM 框架。总体来说, 该 ORB-SLAM 的两篇文章是基于特
 2194 征的 SLAM 系统的大成之作, 系统本身非常的完整, 包含了稀疏地图, 回
 2195 环和优化环节。ORB-SLAM2 将传感器从单目扩展到双目和 RGB-D, 成为
 2196 了迄今为止影响力最大的特征法 SLAM 系统之一, 为 SLAM 在工业界和学
 2197 术界的发展作出巨大贡献。

2198 在特征匹配、跟踪和回环检测等所有环节全部采用统一的 ORB 描述
 2199 子, 通过地图复用, 有效提高了系统的鲁棒性。该系统分为四个模块, 分别
 2200 是前端、局部建图和回环检测以及闭环, 系统使用 3 个线程。。下面我们着
 2201 重介绍 ORB-SLAM2 中重要的几个模块。

2202 13.1.1 单目系统的初始化

2203 系统有两个初始化模块, 双目和 RGB-D 共用一个模块, 单目有自己的
 2204 一个初始化模块, 在追踪开始后, 根据输入图像的样式, 进去相应的初始化

2205 程序。相比于双目和 RGB-D，单目不具有尺度，因此更加的复杂，这里我
2206 们介绍单目的初始化过程。首先是对连续的两帧 (F_i 和 F_{i+1}) 提取特征点，
2207 如果他们满足特征匹配数量的约束。就进一步计算两帧之间的位姿。如果，
2208 不是连续的两帧那么就重新开始。

2209 为了提升初始化的精度，作者同时估计两帧图像之间的单应矩阵 (Ho-
2210 mograph) 和基础矩阵 (Fundamental)，单应矩阵可以用来处理平面之间/较
2211 小视差角的旋转与平移，基础矩阵则能处理更普遍的特征点。对于两帧图像
2212 之间的同名点，系统同时估计上述的两个矩阵 H_{cr} 和 $F(cr)$ ，其中 c 和 r 分
2213 别表示当前帧和参考帧。对于同名点随机选择八个点，并建立很多组。然后
2214 通过公式 13.1 将这两个矩阵分别将一帧图像中同名点转换到另一帧图像上，
2215 来计算同名点在每帧图像上的像素距离。

$$S_M = \sum_i (\rho_M(d_{cr}^2(x_c^i, x_r^i, M)) + \rho_M(d_{rc}^2(x_c^i, x_r^i, M))) \quad (13.1)$$

2216 这里 i 表示同名点的数量， M 为 H_{cr} 或 $F(cr)$ 矩阵中的一个。最后，利用
2217 公式 13.2 来比较两个误差的大小，并选择一个矩阵。

$$R_H = \frac{S_H}{S_H + S_F} \quad (13.2)$$

2218 在 ORB-SLAM 中，如果 $R_H > 0.045$ 选择单应矩阵，否则选择基础矩阵来
2219 初始化两帧图像之间的位姿。在完成位姿之后，系统利用该信息和同名点信
2220 息进行路标点的三角化 (Triangulation) 工作，创建稀疏地图。

2221 13.1.2 姿态估计与优化

2222 对于相机姿态的估计是 SLAM 系统的一个重要任务，这个环节为 AR/VR，
2223 机器人和无人驾驶等应用提供支持。在 ORB-SLAM2 中，利用相机基线来
2224 实现 RGB-D 向双目相机的转化。

- 2225 • 系统在 2D 特征提取中为了使特征分布更加均匀，将图像分块提取特
2226 征。对于双目和 RGB-D 相机，系统可以直接获得特征点的 3D 信息，
2227 而单目相机需要在后面进行三角化操作。
- 2228 • 完成当前帧的特征提取之后，系统通过追踪参考帧或使用运动模型的
2229 方式获得当前帧的初始 6D 位姿。两个方法都是先给当前相机的位姿
2230 赋一个值，然后再用重投影误差去优化。然后使用局部地图再优化一
2231 遍，局部建图模块是利用前端输入的关键帧插入局部地图进行局部优
2232 化。

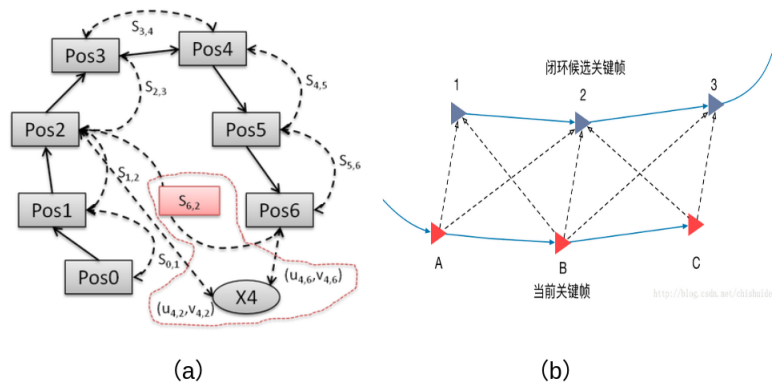


图 47: 闭环检测检测

2233 • 如果检测出来是关键帧，另一个线程基于共视图的优化。

2234 ORB-SLAM 的这种相机位姿估计策略非常鲁棒，尤其是当检测的特征点充
2235 足时，可以实现高精度的位姿估计。

2236 13.1.3 基于词袋的闭环检测

2237 对于基于 Frame-to-Frame 和 Map-to-Frame 的 SLAM 系统，随着追踪
2238 时间的增长，系统也在不断积累着误差，这种现象被称之为系统漂移 (drift)。
2239 回环检测 (Loop Closure) 是一种消除漂移的有效的解决方案。闭环检测是
2240 在历史轨迹中查找是否存在与新关键帧具有相同观测的和约束的关键帧，筛
2241 选过程主要包括位置识别和几何验证。

2242 在 ORBSLAM2 中，闭环检测线程通过词袋法 (DBoW2) 加速闭环匹
2243 配帧的筛选。这种方法的目的是用一种更加抽象、更加高效的特征集合来描
2244 述一幅图像。通过 DBoW2，作者训练了一个较大的特征字典。对于每一帧
2245 图像中提取的特征，只要在字典树中逐层查找，最后都能找到与之对应的单
2246 词。最后可以使用一个分布或直方图来准确的描述这个图像。从图47 (a) 中
2247 可以看到 Pose2 和 Pose6 之间的共同观测。并通过共同观测计算两帧图像
2248 之间的旋转矩阵和平移矩阵，对于单目相机，还要通过 Sim3 优化尺度。然
2249 后在上一步求得的 Sim3 和对应同名点的基础上，纠正了当前帧以及所有有
2250 同视关系的关键帧的位姿，以及路标点。

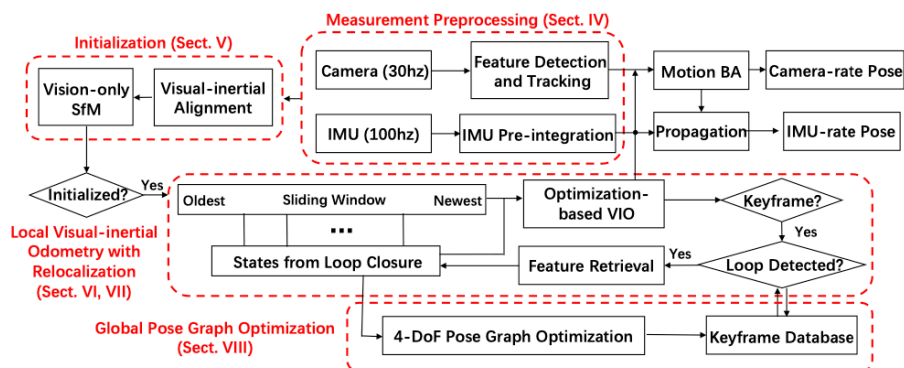


图 48: VINS-Mono 的系统框架。[cite]

2251 13.2 VINS 系列

2252 VINS 系列主要指港科沈绍杰教授团队所提出的 VIO 相关论文和开源
 2253 系统，其中更以 VINS-Mono 为代表，这篇文章的引用量迅速上升至 1000
 2254 (Google Scholar, 06.2021)，成为了 VIO 领域短期之内避不开的重要论文。
 2255 这里我们一起来回顾 VINS-Mono 所提出来的经典框架。

2256 13.2.1 系统概述

2257 如图48所示，系统首先分别以 30hz 和 100hz 的帧率输入单目图片和
 2258 IMU 信息，使用 KLT 光流算法跟踪每帧图像提取的 Harris 角点，IMU 部
 2259 分的处理方式是预积分。系统运行起来后，第一步是初始化程序。无论是
 2260 VINS 还是后面的 ORB-SLAM3，系统都是先进行纯视觉的初始化，视觉初
 2261 始化成功之后再继续进行视觉和 IMU 的共同优化。在追踪过程中，该系统通过
 2262 使用优化模型，构建包含特征点约束、IMU 预积分约束和闭环检测约束的
 2263 优化方程，来求解滑窗内所有帧的位置、速度、旋转和 IMU 的 bias。需要注
 2264 意的是，闭环检测中用到的图像帧是滑窗中选出来的关键帧，通过 DBoW
 2265 方法监测到闭环后，对涉及闭环的关键帧进行全局优化。

2266 13.2.2 预积分和视觉约束

2267 前面在介绍 ORB-SLAM 系列的时候，介绍了特征提取

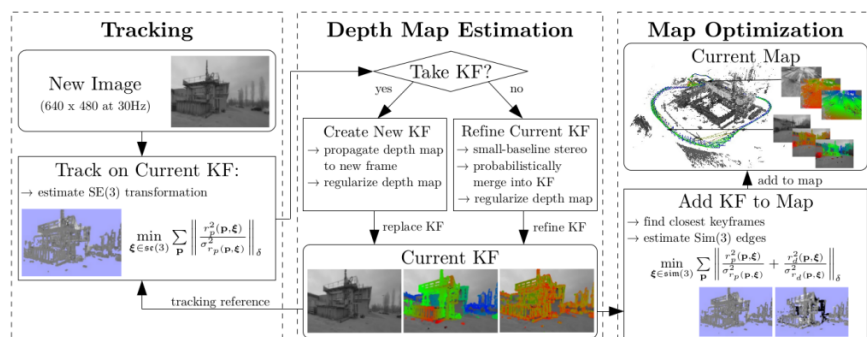


图 49: Pipeline of LSD-SLAM [5]

2268 13.3 LSD-SLAM

2269 LSD-SLAM (large-scale direct monocular SLAM) [5] 是发表在 ECCV2014
 2270 的直接法的 SLAM。与 ORB-SLAM 等所用的特征法不同，直接法视觉里程
 2271 计 (VO) 是直接利用图像像素点的灰度信息来构图与定位，克服了特征点
 2272 提取方法的局限性，可以使用图像上的所有信息。该方法在特征点稀少的环境
 2273 下仍能达到很高的定位精度与鲁棒性，而且提供了更多的环境几何信息，
 2274 这在机器人和增强现实应用中都是非常有意义的。

2275 这篇文章的第一作者 Jakob Engel 对直接法的发展作出了很大的贡献。
 2276 本篇文章提出的方法能够构建大尺度的，全局一致性的环境地图。除了能够
 2277 基于直接图像配准得到高度准确的姿态估计外，还能够将三维环境地图实
 2278 时重构为关键帧的姿态图和对应的半稠密的深度图。这些都是通过对大量
 2279 像素点对之间的基线立体配准结果滤波后得到的。算法提出了计算尺度漂
 2280 移的公式，即便是当图像序列的场景尺度变化较大时也能够适用。

2281 13.3.1 整体流程

2282 算法有三个主要组成部分，分别为图像跟踪、深度图估计和地图优化。
 2283 如图 49 所示。

- 2284 • 图像跟踪：连续跟踪从相机获取到的新“图像帧”。也就是说用前一帧
 2285 图像帧作为初始姿态，估算出当前参考关键帧和新图像帧之间的刚体
 2286 变换 $\xi \in se(3)$ 。
- 2287 • 深度图估计：使用被跟踪的“图像帧”，要么对当前关键帧深度更新，

2288 要么替换当前关键帧。深度更新是基于像素小基线立体配准的滤波方
2289 式，同时耦合对深度地图的正则化。如果相机移动足够远，就初始化
2290 新的关键帧，并把现存相近的关键帧图像点投影到新建立的关键帧上。

2291 • 地图优化模块：一旦关键帧被当前的图像替代，它的深度信息将不会
2292 再被进一步优化，而是通过地图优化模块插入到全局地图中。为了检
2293 测闭环和尺度漂移，采用尺度感知的直接图像配准方法来估计当前帧
2294 与现有邻近关键帧之间的相似性变换。

2295 13.3.2 光度误差与深度误差

2296 关键帧是 LSD-SLAM 系统中重要的媒介。启动 LSD-SLAM 系统时，
2297 只需要初始化首帧关键帧即可，而关键帧深度信息初步设定为一个方差很
2298 大的随机变量。在算法运行最开始的几秒钟，一旦摄像头运动了足够的平移
2299 量，LSD-SLAM 算法就会“锁定”到某个特定的深度配置，经过几个关键
2300 帧的传递之后，就会收敛到正确的深度配置。

$$E(\xi) = \sum_i (I_{ref}(p_i) - I(w(p_i, D_{ref}(p_i), \xi)))^2 \quad (13.3)$$

2301 公式 13.3 是当前帧 I 与参考帧 I_{ref} 之间的光度误差。其中 $w(p_i, D_{ref}(p_i), \xi)$
2302 是将这些像素重投影到参考帧 I_{ref} 。一旦新图像帧被选择成为关键帧，把上
2303 一帧关键帧的兴趣点投影到新创建的关键帧上得到这一帧的兴趣点。之后，
2304 深度图被平均至逆深度为 1，就可以使用相似变换 $\text{sim}(3)$ 来计算关键帧之
2305 间的边，因为相似变换 $\text{sim}(3)$ 可以较好地结合关键帧之间的尺度缩放差异。

2306 除了包含光度测量残差 r_p ，该论文还引入深度残差 (depth residual) r_d
2307 来惩罚关键帧之间的逆深度偏差，能直接够估计出帧间的相似变换。LSD-
2308 SLAM 论文中使用的误差函数 $E(\xi_{ji})$ ，

$$E(\xi_{ji}) = \sum \left\| \frac{r_p^2(p, \xi_{ji})}{\sigma_{r_p(p, \xi_{ji})}^2} + \frac{r_d^2(p, \xi_{ji})}{\sigma_{r_d(p, \xi_{ji})}^2} \right\| \quad (13.4)$$

2309 其中 r_d 是深度残差， $\sigma_{r_p}^{-2}$ 和 $\sigma_{r_d}^{-2}$ 分别是光度残差和深度残差的逆方差。

2310 13.4 SVO

2311 全称 fast Semi-direct monocular Visual Odometry (半直接视觉里程
2312 计)，是苏黎世大学机器人感知组的克里斯蒂安·弗斯特 (Christian Forster

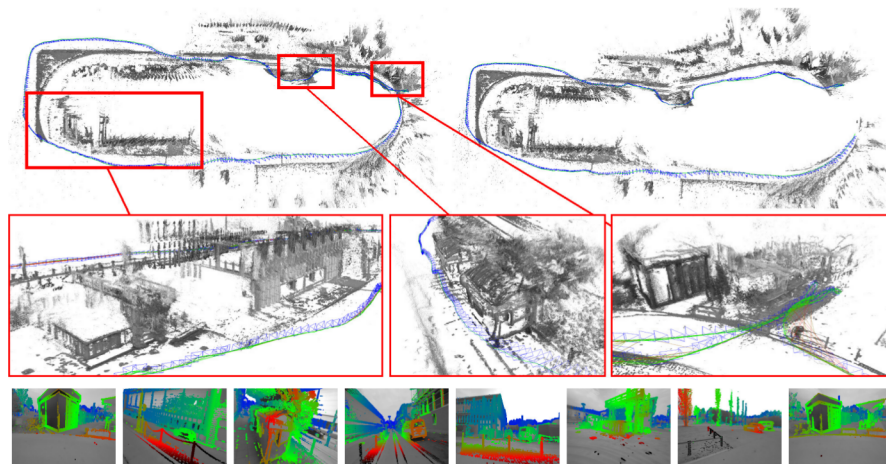


图 50: LSD-SLAM 中的闭环检测模块。左图是闭环检测的效果，右图是没有闭环检测。[5]

2313 等人于 2014 年 ICRA 会议上发表的工作。2016 年扩展了多相机和 IMU 之
 2314 后，写成期刊论文，称为 SVO 2.0。与基于特征法的 ORB-SLAM 和直接法
 2315 的 LSD-SLAM 不同，SVO 提取了稀疏点，但是使用直接法来做图像之间
 2316 的对齐，故称它为半直接法。高翔博士在 LDSO [8] 中同样使用特征点和光
 2317 度，特征的主要作用是实现闭环检测，只不过现在半稠密的概念已经不如以
 2318 往流行了。SVO 结合了直接法和特征点法，但是特征点只是从关键帧中提
 2319 取，以保证系统的运行速度。SVO 算法分成两部分：位姿估计，深度估计。
 2320 下文详细介绍这两个模块。

2321 13.4.1 位姿估计

2322 对稀疏的特征块使用 direct method 配准，获取相机位姿；通过获取的
 2323 位姿预测参考帧中的特征块在当前帧中的位置，由于深度估计的不准导致
 2324 获取的位姿也存在偏差，从而使得预测的特征块位置不准。由于预测的特
 2325 征块位置和真实位置很近，所以可以使用牛顿迭代法对这个特征块的预测
 2326 位置进行优化。特征块的预测位置得到优化，说明之前使用直接法预测的
 2327 有问题。利用这个优化后的特征块预测位置，再次使用直接法，对相机位姿
 2328 (pose) 以及特征点位置 (structure) 进行优化。而特征点值从关键帧中提取，
 2329 文章使用重投影误差来做 pose refinement。

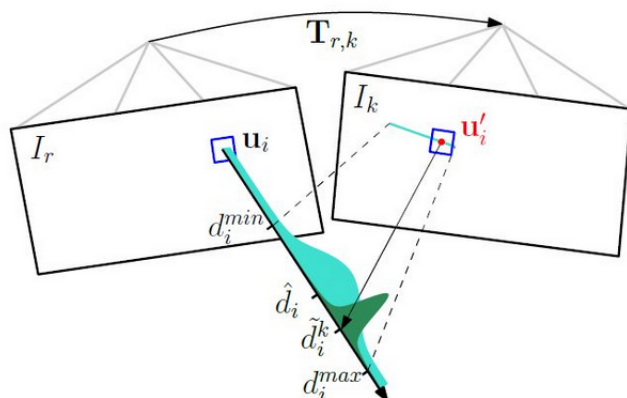


图 51: 深度的高斯-均匀分布。

2330 13.4.2 深度估计

2331 除了 tracking 线程，svo 还有一个 mapping 线程。2d 特征的深度估计
 2332 是在这个线程中完成的。SVO 中使用高斯-均匀分布对逆深度信息建模。如
 2333 果当前帧被选为关键帧，就提取关键帧上的新特征点，这些特征点并未有对
 2334 应的深度信息，因此需要获得计算深度。如果进来一个普通帧，就用普通帧
 2335 的信息，更新所有种子点的概率分布。于是你从这条极线的一个端点走到另
 2336 一个端点，把每处的图像块都和参考的去比较，就可以找到正确的匹配。如
 2337 果某个种子点的深度分布已经收敛，就把它放到地图中，供追踪线程使用。

2338 14 深度 SLAM

2339 深度 SLAM 用来描述融合几何与深度学习的 SLAM 系统。

2340 14.1 Structure-SLAM

2341 该文章提出一个单目 SLAM 系统，使用深度学习网络进行表面法向量
 2342 (surface normal) 估计，然后使用曼哈顿世界 (Manhattan World) 模型，来
 2343 估计相机的旋转矩阵，然后利用点线特征估计平移向量。获得旋转、平移
 2344 矩阵和局部地图之后，通过地图对相机姿态进行一次优化。论文使用 TUM
 2345 RGB-D 和 ICL NUIM 数据集进行测试。

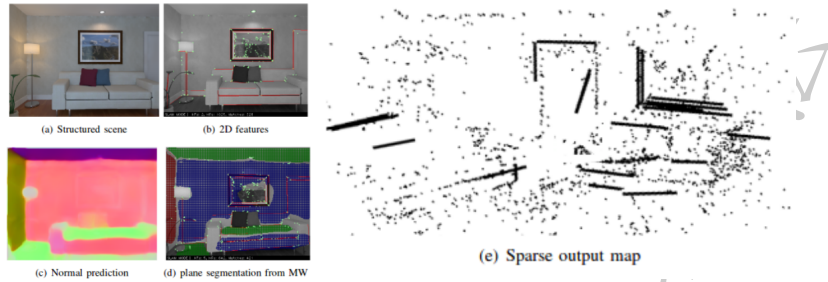


图 52: 效果图

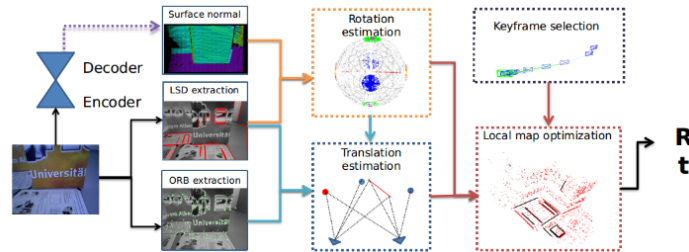


Fig. 2. Proposed SLAM framework (StructureSLAM). In the front-end, the encoder-decoder network predicts dense surface normals. In parallel, point and line features are extracted from the RGB image. In the back-end, first the scene structure in the form of normals and lines is used to estimate the global rotation of the camera. Then, the remaining 3-DoF for the translation are obtained using point and line features. The initial pose estimate is validated and refined using the local map. Keyframes are selected based on the availability of point and line features.

图 53: Structure-SLAM 框架 [19]

2346 2. 基于曼哈顿世界假设的旋转矩阵估计

2347 **什么是曼哈顿世界假设** 曼哈顿假设：如下图所示，环境中存在垂直/正交
 2348 的信息，如地板、天花板、墙面等。通过判断相机和曼哈顿世界之间的旋转
 2349 矩阵，进而可以获得相机与相机之间的相对旋转矩阵。

2350 **旋转矩阵估计** 首先是进行 normal 估计，第四列是使用深度图计算出来的
 2351 平面分割结果。二三列是网络基于 RGB 图估计出来的结果。

2352 将估计的法向量结果投影在一个单位球 (unit sphere) 上，如下如所示，
 2353 一部分法向量会聚集在一起 (平行平面)，由于估计误差，也会有分散在球
 2354 面上的法向量。该文章使用聚类方法-mean shift，对球面上的法向量进行聚
 2355 类。

2356 基于 sphere mean shift 的处理方式如下：

2357 开始假设最开始的 $R_{c_0, M}$ 是单位矩阵，然后将法向量投影在 x, y, z

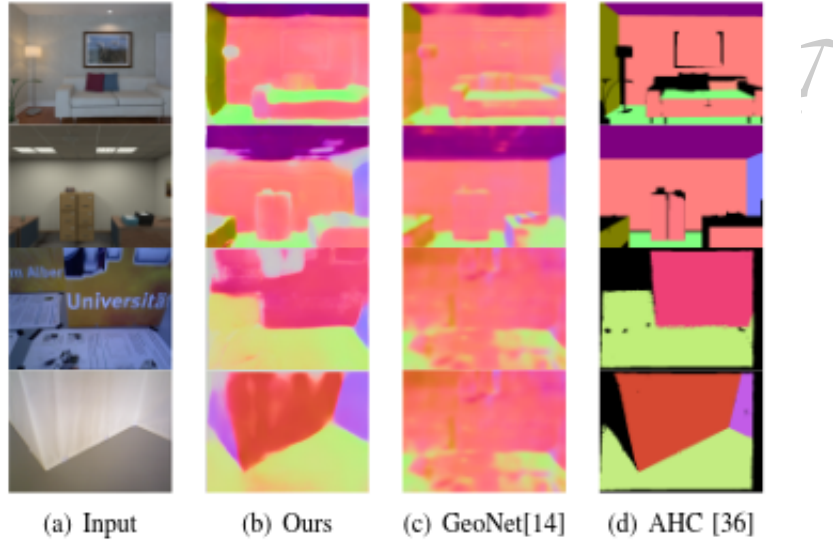


Fig. 4. Results of normal prediction on ICL-NUIM (top) and TUM-RGBD (bottom) scenes for different approaches.

图 54: Caption

2358 轴的切平面上，然后在 2D 且平面上进行 mean-shift 聚类：

$$s'_n = \frac{\sum_{in} e^{-c\|m'_{in}\|^2} m'_{in}}{\sum_{in} e^{-c\|m'_{in}\|^2}} \quad (14.1)$$

2359 聚类结果之后，我们把聚类结果投影回 3D 空间 s_n ，形成两个坐标系
2360 之间的夹角。

$$\hat{r}_n = Q_n s_n \quad (14.2)$$

2361 这里 $Q_n = [r_{\text{mod}(n,3)}, r_{\text{mod}(n+1,3)}, r_{\text{mod}(n+2,3)}]$ ，其中 $\text{mod}()$ 是一个取模
2362 操作，经过一个迭代过程，我们可以获得 $R_{C_1 M} = [\hat{r}_1, \hat{r}_2, \hat{r}_3]^T$ 。这样就更新
2363 出来了 $R_{C_{k+1}, M}$ ，也就可以获得 $R_{C_{k+1}, w} = R_{C_{k+1}, M} * R_{M, w}$

2364 Note 我们获得 $R_{C_{k+1}, w}$ ，但是它不依赖与前面的图像帧或关键帧，而是
2365 一个不变的量 $R_{M, w}$ 。因此这就是一个很有趣的概念：**drift-free rotation**
2366 **estimation.**

2367 **平移向量估计** 由于这里只计算 3 自由度的平移向量，因此能更好的适应
2368 low-textured 的场景。

$$e_{k,j}^p = p_k - \pi(R_{k,j}P_j + t_{k,j}) \quad (14.3)$$

2369

$$e_{k,j}^l = l\pi(R_{k,j}P_x + t_{k,j}) \quad (14.4)$$

2370 这是使用点线特征来完成平移向量的计算。

参考文献

- 2371
- 2372 [1] Michael Calonder, Vincent Lepetit, Christoph Strecha, and Pascal Fua.
2373 Brief: Binary robust independent elementary features. In *European*
2374 *conference on computer vision*, pages 778–792. Springer, 2010.
- 2375 [2] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM
2376 Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source
2377 library for visual, visual-inertial and multi-map slam. *arXiv preprint*
2378 *arXiv:2007.11898*, 2020.
- 2379 [3] David Eigen and Rob Fergus. Predicting depth, surface normals and
2380 semantic labels with a common multi-scale convolutional architecture.
2381 In *Proceedings of the IEEE international conference on computer vision*,
2382 pages 2650–2658, 2015.
- 2383 [4] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map predic-
2384 tion from a single image using a multi-scale deep network. In *Advances*
2385 *in neural information processing systems*, pages 2366–2374, 2014.
- 2386 [5] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-
2387 scale direct monocular slam. In *European conference on computer vi-*
2388 *sion*, pages 834–849. Springer, 2014.
- 2389 [6] Chen Feng, Yuichi Taguchi, and Vineet R Kamat. Fast plane extraction
2390 in organized point clouds using agglomerative hierarchical clustering.
2391 In *2014 IEEE International Conference on Robotics and Automation*
2392 *(ICRA)*, pages 6218–6225. IEEE, 2014.
- 2393 [7] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich,
2394 and Dacheng Tao. Deep ordinal regression network for monocular depth
2395 estimation. In *Proceedings of the IEEE Conference on Computer Vision*
2396 *and Pattern Recognition*, pages 2002–2011, 2018.
- 2397 [8] Xiang Gao, Rui Wang, Nikolaus Demmel, and Daniel Cremers. Ldso:
2398 Direct sparse odometry with loop closure. In *2018 IEEE/RSJ Inter-*
2399 *national Conference on Intelligent Robots and Systems (IROS)*, pages
2400 2198–2204. IEEE, 2018.

- 2401 [9] Ruben Gomez-Ojeda, Francisco-Angel Moreno, David Zuñiga-Noël,
2402 Davide Scaramuzza, and Javier Gonzalez-Jimenez. Pl-slam: a stereo
2403 slam system through the combination of points and line segments.
2404 *IEEE Transactions on Robotics*, 35(3):734–746, 2019.
- 2405 [10] Yijia He, Ji Zhao, Yue Guo, Wenhao He, and Kui Yuan. Pl-vio:
2406 Tightly-coupled monocular visual-inertial odometry using point and
2407 line features. *Sensors*, 18(4):1159, 2018.
- 2408 [11] Dirk Holz, Stefan Holzer, Radu Bogdan Rusu, and Sven Behnke. Real-
2409 time plane segmentation using rgb-d cameras. In *Robot Soccer World*
2410 *Cup*, pages 306–317. Springer, 2011.
- 2411 [12] Pyojin Kim, Brian Coltin, and H Jin Kim. Linear rgb-d slam for planar
2412 environments. In *European Conference on Computer Vision (ECCV)*,
2413 pages 333–348, 2018.
- 2414 [13] Pyojin Kim, Brian Coltin, and H Jin Kim. Visual odometry with
2415 drift-free rotation estimation using indoor scene regularities. In *British*
2416 *Machine Vision Conference (BMVC)*, 2017.
- 2417 [14] Pyojin Kim, Brian Coltin, and H Jin Kim. Low-drift visual odometry
2418 in structured environments by decoupling rotational and translational
2419 motion. In *IEEE International Conference on Robotics and Automation*
2420 *(ICRA)*, pages 7247–7253. IEEE, 2018.
- 2421 [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic
2422 optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- 2423 [16] Georg Klein and David Murray. Parallel tracking and mapping for small
2424 ar workspaces. In *2007 6th IEEE and ACM international symposium*
2425 *on mixed and augmented reality*, pages 225–234. IEEE, 2007.
- 2426 [17] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico
2427 Tombari, and Nassir Navab. Deeper depth prediction with fully con-
2428 volutional residual networks. In *2016 Fourth international conference*
2429 *on 3D vision (3DV)*, pages 239–248. IEEE, 2016.

- 2430 [18] Xin Li, Yanyan Li, Evin Pnar Ornek, Jinlong Lin, and Federico
2431 Tombari. Co-planar parametrization for stereo-slam and visual-inertial
2432 odometry. *IEEE Robotics and Automation Letters*, 2020.
- 2433 [19] Yanyan Li, Nikolas Brasch, Yida Wang, Nassir Navab, and Federico
2434 Tombari. Structure-slam: Low-drift monocular slam in indoor envi-
2435 ronments. *IEEE Robotics and Automation Letters*, 5(4):6583–6590,
2436 2020.
- 2437 [20] Yanyan Li, Shiyue Fan, Yanbiao Sun, Qiang Wang, and Shanlin Sun.
2438 Bundle adjustment method using sparse bfgs solution. *Remote Sensing*
2439 *Letters*, 9(8):789–798, 2018.
- 2440 [21] Yanyan Li, Raza Yunus, Nikolas Brasch, Nassir Navab, and Federico
2441 Tombari. Rgb-d slam with structural regularities. *arXiv preprint*
2442 *arXiv:2010.07997*, 2020.
- 2443 [22] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolu-
2444 tional networks for semantic segmentation. In *Proceedings of the IEEE*
2445 *conference on computer vision and pattern recognition*, pages 3431–
2446 3440, 2015.
- 2447 [23] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos.
2448 Orb-slam: a versatile and accurate monocular slam system. *IEEE*
2449 *transactions on robotics*, 31(5):1147–1163, 2015.
- 2450 [24] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam
2451 system for monocular, stereo, and rgb-d cameras. *IEEE Transactions*
2452 *on Robotics*, 33(5):1255–1262, 2017.
- 2453 [25] Richard A Newcombe, Steven J Lovegrove, and Andrew J Davison.
2454 Dtam: Dense tracking and mapping in real-time. In *2011 international*
2455 *conference on computer vision*, pages 2320–2327. IEEE, 2011.
- 2456 [26] Albert Pumarola, Alexander Vakhitov, Antonio Agudo, Alberto Sanfe-
2457 liu, and Francesc Moreno-Noguer. Pl-slam: Real-time monocular visual
2458 slam with points and lines. In *2017 IEEE international conference on*
2459 *robotics and automation (ICRA)*, pages 4503–4508. IEEE, 2017.

- 2460 [27] Xiaojuan Qi, Renjie Liao, Zhengzhe Liu, Raquel Urtasun, and Jiaya
2461 Jia. Geonet: Geometric neural network for joint depth and surface
2462 normal estimation. In *Proceedings of the IEEE Conference on Computer
2463 Vision and Pattern Recognition*, pages 283–291, 2018.
- 2464 [28] Tong Qin, Peiliang Li, and Shaojie Shen. Vins-mono: A robust and
2465 versatile monocular visual-inertial state estimator. *IEEE Transactions
2466 on Robotics*, 34(4):1004–1020, 2018.
- 2467 [29] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov,
2468 and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear
2469 bottlenecks. In *Proceedings of the IEEE conference on computer vision
2470 and pattern recognition*, pages 4510–4520, 2018.
- 2471 [30] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva,
2472 and Thomas Funkhouser. Semantic scene completion from a single
2473 depth image. In *Proceedings of the IEEE Conference on Computer
2474 Vision and Pattern Recognition*, pages 1746–1754, 2017.
- 2475 [31] Jean-Philippe Tardif. Non-iterative approach for fast and accurate van-
2476 ishing point detection. In *2009 IEEE 12th International Conference on
2477 Computer Vision*, pages 1250–1257. IEEE.
- 2478 [32] Keisuke Tateno, Federico Tombari, Iro Laina, and Nassir Navab. Cnn-
2479 slam: Real-time dense monocular slam with learned depth prediction.
2480 In *Proceedings of the IEEE Conference on Computer Vision and Pattern
2481 Recognition*, pages 6243–6252, 2017.
- 2482 [33] Deepak Geetha Viswanathan. Features from accelerated segment test
2483 (fast).
- 2484 [34] L. von Stumberg, V. Usenko, and D. Cremers. Direct sparse visual-
2485 inertial odometry using dynamic marginalization. May 2018.
- 2486 [35] Xiaoyu Zhang, Wei Wang, Xianyu Qi, Ziwei Liao, and Ran Wei. Point-
2487 plane slam using supposed planes for indoor environments. *Sensors*,
2488 19(17):3795, 2019.

- 2489 [36] Yinda Zhang and Thomas Funkhouser. Deep depth completion of a
2490 single rgb-d image. In *Proceedings of the IEEE Conference on Computer
2491 Vision and Pattern Recognition*, pages 175–185, 2018.
- 2492 [37] Yinda Zhang, Shuran Song, Ersin Yumer, Manolis Savva, Joon-Young
2493 Lee, Hailin Jin, and Thomas Funkhouser. Physically-based rendering
2494 for indoor scene understanding using convolutional neural networks. In
2495 *Proceedings of the IEEE Conference on Computer Vision and Pattern
2496 Recognition*, pages 5287–5295, 2017.
- 2497 [38] Yi Zhou, Laurent Kneip, Cristian Rodriguez, and Hongdong Li. Divide
2498 and conquer: Efficient density-based tracking of 3d sensors in manhat-
2499 tan worlds. In *Asian Conference on Computer Vision (ACCV)*, pages
2500 3–19. Springer, 2016.
- 2501 [39] Danping Zou, Yuanxin Wu, Ling Pei, Haibin Ling, and Wenxian Yu.
2502 Structvio: visual-inertial odometry with structural regularity of man-
2503 made environments. *IEEE Transactions on Robotics*, 35(4):999–1013,
2504 2019.