

计算机视觉总结

SUMMARY OF COMPUTER VISION

(第1版)
LVSHUAILIN

OPEN SOURCE
BEIJING

VERSION 1

- 一. 数据结构与算法-LeetCode Hot 100
- 二. PYTHON: 1) NUMPY; 2) PANDAS; 3) PYTHON多进程; 4) PYTHON分布式; 5) PYTHON界面;
- 三. 深度学习: TensorFlow 2.0; PYTORCH;
- 四. 图像配准
- 五. 强化学习
- 六. OTHERS: 1) Model INFERENCE by EXE; 2) GIT; 3) DOCKER

LVSHUAILIN

2020年2月

目 录

| | |
|--|----|
| 第1章 LeetCode Hot 100..... | 1 |
| 1.1 两数之和 | 1 |
| 1.1.1 知识点(unordered_map) | 1 |
| 1.1.2 解题代码 | 3 |
| 1.2 两数相加 | 4 |
| 1.2.1 知识点(linked list)..... | 4 |
| 1.2.2 解题代码 | 9 |
| 1.3 无重复字符的最长子串 | 12 |
| 1.3.1 知识点(double pointer algorithm和unordered_set)..... | 12 |
| 1.3.2 解题思路 | 13 |
| 1.3.3 解题代码 | 13 |
| 1.4 寻找两个有序数组的中位数 | 14 |
| 1.4.1 知识点(二分查找算法)..... | 14 |
| 1.4.2 解题思路 | 16 |
| 1.4.3 解题代码 | 16 |
| 第2章 Bayesian-Python..... | 18 |
| 2.1 Probabilities..... | 18 |
| 2.1.1 Gaussian Distribution..... | 18 |
| 2.1.2 Coin-Flipping Problem | 19 |
| 2.1.2.1 The general model | 19 |
| 2.1.2.2 Choosing the likelihood | 20 |
| 2.1.2.3 Choosing the prior | 21 |
| 2.1.2.4 Getting the posterior | 22 |

| | | |
|---------|--|----|
| 2.1.2.5 | Model notation and visualization | 24 |
| 2.1.3 | Inference engines | 24 |
| 2.1.3.1 | Non-Markovian methods | 25 |
| 2.1.4 | PyMC3 introduction | 34 |
| 第3章 | 深度学习 | 35 |
| 3.1 | Pytorch | 35 |

第1章 LeetCode Hot 100

Goals to Achieve

1. unordered_map.

§ 1.1 两数之和

HOT100 1.1 问题描述

给定一个整数数组**nums** 和一个目标值**target**, 请你在该数组中找出和为目标值的那两个整数, 并返回他们的数组下标. 你可以假设每种输入只会对应一个答案. 但是, 你不能重复利用这个数组中同样的元素.

示例: 给定nums = [2, 7, 11, 15], target = 9; 因为nums[0] + nums[1] = 2 + 7 = 9;

所以返回[0, 1]

<https://leetcode-cn.com/problems/two-sum>

1.1.1 知识点(unordered_map)

unordered_map内部是一个关联容器, 采用hash 表结构, 有快速检索的功能.

哈希表是通过key关键字直接访问对应value值的数据结构. 特点是键和值一一对应, 查找时间复杂度**O(1)**.

Example_1: unordered_map插入, 迭代遍历.

unordered_map example_1 code

```
1 #include <iostream>
2 #include <unordered_map>
3 #include <string>
4 using namespace std;
5 int main()
6 {
```

```

7 unordered_map<string, double> umap;
8 umap["PI"] = 3.14;
9 umap.insert(make_pair("a", 2.1));
10
11 // find in umap
12 string key = "PI";
13 if (umap.find(key) == umap.end())
14     cout << "cannot find PI" << endl;
15 else
16     cout << "find " << umap.find(key)->first << " = " << umap.find(key)->second << endl;
17
18 // iterator of umap
19 cout << "entire unorded_map is: " << endl;
20 unordered_map<string, double>::iterator itr;
21 for (itr = umap.begin(); itr != umap.end(); ++itr)
22     cout << " ( " << itr->first << ", " << itr->second << " ) " << endl;
23 system("pause");
24 return 0;
25 }

```

output:

```

find PI = 3
all elements are:
(PI,3.14)
(a,2.1)

```

Example_2: 利用unordered_map输出一段文字中重复单词的个数

unordered_map example_2 code

```

1 #include <iostream>
2 #include <unordered_map>
3 #include <string>
4 #include <sstream>
5
6 using namespace std;
7
8 void printWordFreq(const string& str)
9 {
10     unordered_map<string, int> wordFreq;
11     string word;
12     stringstream ss(str);

```

```
13 while (ss >> word)
14     wordFreq[word]++;
15
16 cout << "all elements are:" << endl;
17 for (auto u : wordFreq)
18     cout << " (" << u.first << ", " << u.second << ") " << endl;
19 }
20
21 int main()
22 {
23     string str = "studies very very hard";
24     printWordFreq(str);
25     return 0;
26 }
```

output:

```
all elements are:
(studies, 1)
(very, 2)
(hard, 1)
```

1.1.2 解题代码

```
1 #include <iostream>
2 #include <unordered_map>
3 #include <vector>
4 using namespace std;
5
6 vector<int> twoSum(vector<int>& nums, int target)
7 {
8     unordered_map<int, int> map;
9     vector<int> result={};
10    int n = (int)nums.size();
11    for(int i = 0; i < n; ++i) {
12        auto p = map.find(target-nums[i]);
13        if(p != map.end()) {
14            result.push_back(p->second);
15            result.push_back(i);
16        }
```

```
17     map[nums[i]] = i;
18     }
19     return result;
20 }
21
22 int main()
23 {
24     vector<int> nums = {2,7,11,15};
25     vector<int> result;
26     result = twoSum(nums,9);
27     cout<<" ["<<result[0]<<" , "<<result[1]<<" ] "<<endl;
28     return 0;
29 }
```

§ 1.2 两数相加

HOT100 1.2 问题描述

给出两个非空的链表用来表示两个非负的整数。其中，它们各自的位数是按照逆序的方式存储的，并且它们的每个节点只能存储一位数字。如果，我们将这两个数相加起来，则会返回一个新的链表来表示它们的和。

您可以假设除了数字0之外，这两个数都不会以0开头。

示例: 输入(2 → 4 → 3) + (5 → 6 → 4), 输出: 7 → 0 → 8, 原因: 342 + 465 = 807

<https://leetcode-cn.com/problems/add-two-numbers>

1.2.1 知识点(linked list)

这里用c++ 链表来解决

Example_1: 创建链表并初始化

linked list example_1 code

```
1 #include <iostream>
2
3 using namespace std;
4
5 class Node{
6 public:
```



```
7     int data;
8     Node* next;
9 };
10
11 int main()
12 {
13     Node* head = nullptr;
14     Node* second = nullptr;
15     Node* third = nullptr;
16
17     head = new Node();
18     head->data = 1;
19
20     second = new Node();
21     second->data = 2;
22
23     third = new Node();
24     third->data = 3;
25
26     cout << head->data << " " << second->data << " " << third->data << endl;
27
28     delete head;
29     delete second;
30     delete third;
31     return 0;
32 }
```

output:

1 2 3

Example_2: 打印链表中的所有元素

linked list example_2 code

```
1 #include <iostream>
2
3 using namespace std;
4
5 class Node{
6 public:
7     int data;
8     Node* next;
```

```
9 };
10
11 void PrintLinkedList(Node* head)
12 {
13     Node* temp = head;
14     while (temp != nullptr) {
15         cout << temp->data << " ";
16         temp = temp->next;
17     }
18     cout << endl;
19 }
20
21 int main()
22 {
23     Node* head = nullptr;
24     Node* second = nullptr;
25     Node* third = nullptr;
26
27     head = new Node();
28     second = new Node();
29     third = new Node();
30
31     head->data = 1;
32     head->next = second;
33
34     second->data = 2;
35     second->next = third;
36
37     third->data = 3;
38     third->next = nullptr;
39
40
41     PrintLinkedList(head);
42
43     delete head;
44     delete second;
45     delete third;
46     return 0;
47 }
```

output:

1 2 3

Example_3: 链表插入节点

linked list example_3 code

```
1 #include <iostream>
2
3 using namespace std;
4
5 class Node{
6 public:
7     int data;
8     Node* next;
9 };
10
11 // 在链表前面插入节点
12 void push(Node** head_ref, int newData)
13 {
14     Node* newNode = new Node();
15     newNode->data = newData;
16     newNode->next = (*head_ref);
17     (*head_ref) = newNode;
18 }
19
20 // 在节点后面插入节点
21 void insertAfter(Node** prev_node, int newData)
22 {
23     if ((*prev_node) == nullptr) {
24         cout << "the previous node cannot be nullptr" << endl;
25         return;
26     }
27
28     Node* newNode = new Node();
29     newNode->data = newData;
30     newNode->next = (*prev_node)->next;
31     (*prev_node)->next = newNode;
32 }
33
34 // 在尾节点后插入节点
35 void append(Node** head_ref, int newData)
```

```
36 {
37     Node* newNode = new Node();
38     newNode->data = newData;
39     newNode->next = nullptr;
40     if ((*head_ref) == nullptr) {
41         (*head_ref) = newNode;
42         return;
43     }
44
45     Node* move = (*head_ref);
46     while (move->next != nullptr) {
47         move = move->next;
48     }
49     move->next = newNode;
50 }
51
52 // 打印链表
53 void PrintLinkedList(Node* head)
54 {
55     Node* temp = head;
56     while (temp != nullptr) {
57         cout << temp->data << " ";
58         temp = temp->next;
59     }
60     cout << endl;
61 }
62
63 void destroyLinkedList(Node** head_ref) {
64     Node* move = (*head_ref);
65     Node* next = nullptr;
66     while (move != nullptr) {
67         next = move->next;
68         delete move;
69         move = next;
70     }
71     (*head_ref) = nullptr;
72 }
73
74 int main()
75 {
```

```
76 Node* head = nullptr;
77
78 append(&head, 6);
79
80 push(&head, 7);
81
82 push(&head, 1);
83
84 append(&head, 4);
85
86 insertAfter(&(amp;head->next), 8);
87
88 cout << "linked list is: ";
89 PrintLinkedList(head);
90 destroyLinkedList(&head);
91 return 0;
92 }
```

output:

linked list is: 1 7 8 6 4

1.2.2 解题代码

```
1 #include <iostream>
2 using namespace std;
3
4 struct ListNode {
5     int val;
6     ListNode *next;
7     ListNode(int x) : val(x), next(NULL) {}
8 };
9
10 ListNode* addTwoNumbers(ListNode* l1, ListNode* l2) {
11     int len1 = 1;//记录的长度l1
12     int len2 = 1;//记录的长度l2
13     ListNode* p = l1;
14     ListNode* q = l2;
15     while (p->next != NULL)//获取的长度l1
16     {
```

```
17     len1++;
18     p = p->next;
19 }
20 while (q->next != NULL)//获取的长度l2
21 {
22     len2++;
23     q = q->next;
24 }
25 if (len1 > len2)//较长，在末尾补零l1l2
26 {
27     for (int i = 1; i <= len1 - len2; i++)
28     {
29         q->next = new ListNode(0);
30         q = q->next;
31     }
32 }
33 else//较长，在末尾补零l2l1
34 {
35     for (int i = 1; i <= len2 - len1; i++)
36     {
37         p->next = new ListNode(0);
38         p = p->next;
39     }
40 }
41 p = l1;
42 q = l2;
43 bool count = false;//记录进位
44 ListNode* l3 = new ListNode(-1);//存放结果的链表
45 ListNode* w = l3;//的移动指针l3
46 int i = 0;//记录相加结果
47 while (p != NULL && q != NULL)
48 {
49     i = count + p->val + q->val;
50     w->next = new ListNode(i % 10);
51     count = i >= 10 ? true : false;
52     w = w->next;
53     p = p->next;
54     q = q->next;
55 }
56 if (count)//若最后还有进位
```

```
57     {
58         w->next = new ListNode(1);
59         w = w->next;
60     }
61     return l3->next;
62 }
63
64 void printLinkedList(ListNode* head)
65 {
66     ListNode* move = head;
67     while (move != nullptr) {
68         cout << move->val << " ";
69         move = move->next;
70     }
71 }
72
73 int main()
74 {
75     #if 1
76         ListNode* l1 = new ListNode(2);
77         ListNode* l1_1 = new ListNode(4);
78         ListNode* l1_2 = new ListNode(3);
79
80         l1->next = l1_1;
81         l1_1->next = l1_2;
82
83
84         ListNode* l2 = new ListNode(5);
85         ListNode* l2_1 = new ListNode(6);
86         ListNode* l2_2 = new ListNode(4);
87         l2->next = l2_1;
88         l2_1->next = l2_2;
89     #endif
90
91     #if 0
92         ListNode* l1 = new ListNode(5);
93         ListNode* l2 = new ListNode(5);
94     #endif
95
96     ListNode* result = addTwoNumbers(l1, l2);
```

```
97     printLinkedList(result);  
98     return 0;  
99 }
```

output:

7 0 8

§ 1.3 无重复字符的最长子串

HOT100 1.3 问题描述

给定一个字符串, 请你找出其中不含有重复字符的最长子串的长度.

示例1:

输入: “abcabcbb”

输出: 3

解释: 因为无重复字符的最长子串是“abc”, 所以其长度为3.

示例2:

输入: “bbbbbb”

输出: 1

解释: 因为无重复字符的最长子串是“b”, 所以其长度为1.

示例3:

输入: “pwwkew”

输出: 3

解释: 因为无重复字符的最长子串是“wke”, 所以其长度为3. 请注意, 你的答案必须是子串的长度, “pwke”是一个子序列, 不是子串.

<https://leetcode-cn.com/problems/longest-substring-without-repeating-characters>

1.3.1 知识点(double pointer algorithm和unordered_set)

c++提供两种关联型数据结构, 1) 树型结构, 如: map, set; 2) hash结构, 如: unordered_map, unordered_set. map和set是有序的, 其他两个是无序的.

1.3.2 解题思路

<https://cloud.tencent.com/developer/article/1377650>

这道题主要用到思路是: 滑动窗口

什么是滑动窗口?

其实就是一个队列, 比如例题中的`abcabcbb`, 进入这个队列(窗口)为`abc`满足题目要求, 当再进入`a`, 队列变成了`abca`, 这时候不满足要求. 所以, 我们要移动这个队列!

如何移动?

我们只要把队列的左边的元素移出就行了, 直到满足题目要求!

一直维持这样的队列, 找出队列出现最长的长度时候, 求出解!

时间复杂度: $O(n)$

1.3.3 解题代码

```
1 #include <iostream>
2 #include <string>
3 #include <unordered_set>
4 #include <algorithm> // max, min
5
6 using namespace std;
7
8 int lengthOfLongestSubstring(string s) {
9     if (s.size() == 0) return 0;
10    unordered_set<char> lookup;
11    int maxStr = 0;
12    int left = 0;
13    for (int i = 0; i < s.size(); i++) {
14        while (lookup.find(s[i]) != lookup.end()) {
15            lookup.erase(s[left]);
16            left++;
17        }
18        maxStr = max(maxStr, i - left + 1);
19        lookup.insert(s[i]);
20    }
21    return maxStr;
22 }
23
24 int main()
```

```
25 {  
26     string str = "abcbabcb";  
27     cout << lengthOfLongestSubstring(str) << endl;  
28     return 0;  
29 }
```

output:

3

§ 1.4 寻找两个有序数组的中位数

HOT100 1.4 问题描述

给定两个大小为 m 和 n 的有序数组 $nums1$ 和 $nums2$.

请你找出这两个有序数组的中位数, 并且要求算法的时间复杂度为 $O(\log(m + n))$.

你可以假设 $nums1$ 和 $nums2$ 不会同时为空.

示例1:

$nums1 = [1, 3]$

$nums2 = [2]$

则中位数是2.0

示例2:

$nums1 = [1, 2]$

$nums2 = [3, 4]$

则中位数是 $(2 + 3)/2 = 2.5$

<https://leetcode-cn.com/problems/median-of-two-sorted-arrays>

1.4.1 知识点(二分查找算法)

用二分查找算法, 也叫做折半查找算法.

Example_1: 二分查找

二分查找算法*example_1 code*

```
1 // 二分查找— 折半查找  
2 int search(int arr[], int key, int left, int right)  
3 {  
4     while (left <= right)
```

```
5     {
6         int mid = left + (right - left) / 2;
7         if (key < arr[mid])
8             right = mid - 1;
9         else if (key > arr[mid])
10            left = mid + 1;
11        else
12            return mid;
13    }
14    return -1;
15 }
16
17 int main()
18 {
19     int arr[] = { 0,2 ,3,4};
20     int value = 3;
21
22     // left Index of the array
23     int left = 0;
24
25     // right Index of the array
26     int right = sizeof(arr) / sizeof(arr[0]) - 1;
27
28     cout << "left: " << left << ", right: " << right << endl;
29
30     int ret = search(arr, value, left, right);
31     if (ret == -1)
32         printf("cannot find the value");
33     else
34         printf("found the value, the index is: %d\n", ret);
35     system("pause");
36     return 0;
37 }
```

output:

left: 0, right: 3

found the value, the index is: 2

1.4.2 解题思路

这道题让我们求两个有序数组的中位数，而且限制了时间复杂度为 $O(\log(m+n))$ ，看到这个时间复杂度，自然而然的想到了应该使用二分查找法来求解。那么回顾一下中位数的定义，如果某个有序数组长度是奇数，那么其中位数就是最中间那个，如果是偶数，那么就是最中间两个数字的平均值。这里对于两个有序数组也是一样的，假设两个有序数组的长度分别为 m 和 n ，由于两个数组长度之和 $m+n$ 的奇偶不确定，因此需要分情况来讨论，对于奇数的情况，直接找到最中间的数即可，偶数的话需要求最中间两个数的平均值。为了简化代码，不分情况讨论，我们使用一个小trick，我们分别找第 $(m+n+1)/2$ 个，和第 $(m+n+2)/2$ 个，然后求其平均值即可，这对奇偶数均适用。假如 $m+n$ 为奇数的话，那么其实 $(m+n+1)/2$ 和 $(m+n+2)/2$ 的值相等，相当于两个相同的数字相加再除以2，还是其本身。

这里我们需要定义一个函数来在两个有序数组中找到第 K 个元素，下面重点来看如何实现找到第 K 个元素。首先，为了避免产生新的数组从而增加时间复杂度，我们使用两个变量 i 和 j 分别来标记数组 $nums1$ 和 $nums2$ 的起始位置。然后来处理一些边界问题，比如当某一个数组的起始位置大于等于其数组长度时，说明其所有数字均已经被淘汰了，相当于一个空数组了，那么实际上就变成了在另一个数组中找数字，直接就可以找出来了。还有就是如果 $K=1$ 的话，那么我们只要比较 $nums1$ 和 $nums2$ 的起始位置 i 和 j 上的数字就可以了。难点就在于一般的情况怎么处理？因为我们需要在两个有序数组中找到第 K 个元素，为了加快搜索的速度，我们要使用二分法，对 K 二分，意思是我们需要分别在 $nums1$ 和 $nums2$ 中查找第 $K/2$ 个元素，注意这里由于两个数组的长度不定，所以有可能某个数组没有第 $K/2$ 个数字，所以我们需要先检查一下，数组中到底存不存在第 $K/2$ 个数字，如果存在就取出来，否则就赋值上一个整型最大值。如果某个数组没有第 $K/2$ 个数字，那么我们就淘汰另一个数字的前 $K/2$ 个数字即可。有没有可能两个数组都不存在第 $K/2$ 个数字呢，这道题里是不可能的，因为我们的 K 不是任意给的，而是给的 $m+n$ 的中间值，所以必定至少会有一个数组是存在第 $K/2$ 个数字的。最后就是二分法的核心啦，比较这两个数组的第 $K/2$ 小的数字 $midVal1$ 和 $midVal2$ 的大小，如果第一个数组的第 $K/2$ 个数字小的话，那么说明我们要找的数字肯定不在 $nums1$ 中的前 $K/2$ 个数字，所以我们可以将其淘汰，将 $nums1$ 的起始位置向后移动 $K/2$ 个，并且此时的 K 也自减去 $K/2$ ，调用递归。反之，我们淘汰 $nums2$ 中的前 $K/2$ 个数字，并将 $nums2$ 的起始位置向后移动 $K/2$ 个，并且此时的 K 也自减去 $K/2$ ，调用递归即可。

1.4.3 解题代码

```
1  /*分清 起始位置和第几个元素*/
2  #include <vector>
3  #include <iostream>
4  #include <algorithm>
5  using namespace std;
6
7  int findKthNumber(vector<int>& nums1, int i, vector<int>& nums2, int j, int k) {
8      if (i >= nums1.size()) return nums2[j + k - 1];
9      if (j >= nums2.size()) return nums1[i + k - 1];
10     //if(k == 1) return (double(nums1[i] + nums2[j]));wrong
```

```
11     if (k == 1) return min(nums1[i], nums2[j]);
12     //查找有没有k个元素的位置/2  $i + k/2 - 1$ 
13     int midVal1 = (i + k / 2 - 1 < nums1.size()) ? nums1[i + k / 2 - 1] : INT_MAX;
14     int midVal2 = (j + k / 2 - 1 < nums2.size()) ? nums2[j + k / 2 - 1] : INT_MAX;
15     if (midVal1 < midVal2)
16         return findKthNumber(nums1, i + k / 2, nums2, j, k - k / 2);
17     else
18         return findKthNumber(nums1, i, nums2, j + k / 2, k - k / 2);
19 }
20 class Solution {
21 public:
22     double findMedianSortedArrays(vector<int>& nums1, vector<int>& nums2) {
23         int m = nums1.size(), n = nums2.size();
24         int left = (m + n + 1) / 2, right = (m + n + 2) / 2;
25         return (findKthNumber(nums1, 0, nums2, 0, left) + findKthNumber(nums1, 0, nums2, 0, right)) / 2.0;
26     }
27 };
```

output:

Null

第2章 Bayesian-Python

学习目标与要求

1. .
2. .
3. .
4. .

§ 2.1 Probabilities

A common and useful conceptualization in statistics is to think that data was generated from some probability distribution with unobserved parameters.

2.1.1 Gaussian Distribution

$$pdf(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp \frac{-(x-\mu)^2}{2\sigma^2}$$

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 from scipy import stats
4 import seaborn as sns
5
6 mu_params = [-1, 0, 1]
7 sd_params = [0.5, 1, 1.5]
8 x = np.linspace(-7, 7, 100)
9 f, ax = plt.subplots(len(mu_params), len(sd_params), sharex=True, sharey=True)
10 for i in range(3):
11     for j in range(3):
```

```

12     mu = mu_params[i]
13     sd = sd_params[j]
14     y = stats.norm(mu, sd).pdf(x)
15     ax[i, j].plot(x,y)
16     ax[i, j].plot(0, 0, label="$\\mu$ = { :3.2f} \\n$\\sigma$={ :3.2f} ".format(mu, sd), alpha=0)
17     ax[i, j].legend(fontsize=12)
18     ax[2,1].set_xlabel(' $x$ ', fontsize=16)
19     ax[1,0].set_ylabel(' $pdf(x)$ ', fontsize=16)
20     plt.tight_layout()

```

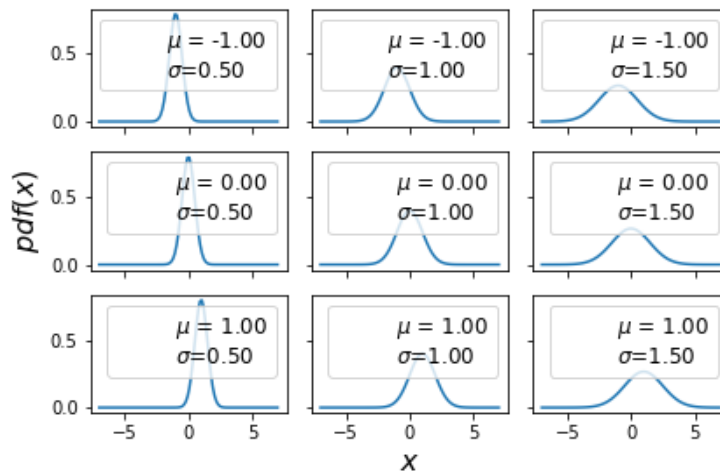


图 2.1 Gaussian Distribution

2.1.2 Coin-Flipping Problem

We will answer this question in a Bayesian setting. We will need data and a probabilistic model.

Data: we will assume that we have already tossed a coin a number of times and we have recorded the number of observed head, so the data part is done.

Model will be discussed soon.

2.1.2.1 The general model

The first thing we will do is generalize the concept of bias. We will say that a coin with a bias of 1 will always land heads, one with a bias of 0 will always land tails, and one with a bias of 0.5 will land half of the time heads and half of the time tails. To represent the bias, we will use the parameter θ , and to represent the total number of heads for an N number of tosses, we will use the variable y . According to Bayes' theorem we have the following formula:

$$p(\theta|y) \propto p(y|\theta)p(\theta)$$

Notice that we need to specify which prior $p(\theta)$ and likelihood $p(y|\theta)$ we will use. Let's start with the likelihood.

2.1.2.2 Choosing the likelihood

Let's assume that a coin toss does not affect other tosses, that is, we are assuming coin tosses are independent of each other. Let's also assume that only two outcomes are possible, heads or tails. Given these assumptions, a good candidate for the likelihood is the binomial distribution (二项分布):

$$p(y|\theta) = \frac{N!}{N!(N-y)!} \theta^y (1-\theta)^{N-y}$$

This is a discrete distribution returning the probability of getting y heads (or in general, success) out of N coin tosses (or in general, trials or experiments) given a fixed value of θ . The following code generates 9 binomial distributions; each subplot has its own legend indicating the corresponding parameters:

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from scipy import stats
4 import seaborn as sns
5
6 n_params = [1,2,4]
7 p_params = [0.25, 0.5, 0.75]
8 x = np.arange(0, max(n_params) + 1)
9 f, ax = plt.subplots(len(n_params), len(p_params), sharex=True, sharey=True)
10 for i in range(3):
11     for j in range(3):
12         n = n_params[i]
13         p = p_params[j]
14         y = stats.binom(n=n, p=p).pmf(x)
15         ax[i,j].vlines(x, 0, y, colors='b', lw=5)
16         ax[i,j].set_ylim(0,1)
17         ax[i,j].plot(0,0,label="n = {:3.2f}\np = {:3.2f}".format(n,p), alpha=0)
18         ax[i,j].legend(fontsize=12)
19 ax[2,1].set_xlabel('$\\theta$', fontsize=14)
20 ax[1,0].set_ylabel('$p(y|\\theta)$', fontsize=14)
21 ax[0,0].set_xticks(x)
22 plt.savefig('binomial_distribution.png')
```

The binomial distribution is also a reasonable choice for the likelihood. Intuitively, we can see that θ indicates how likely it is that we will obtain a head when tossing a coin, and we have observed that event y times. Following the same line of reasoning we get that $1 - \theta$ is the chance of getting a tail, and that event has occurred $N - y$ times.

OK, so if we know θ , the binomial distribution will tell us the expected distribution of head. The only problem is that we do not know θ ! But do not despair; in Bayesian statistics, every time we do not know the value of a parameter, we put a prior on it, so let's move on and choose a prior.

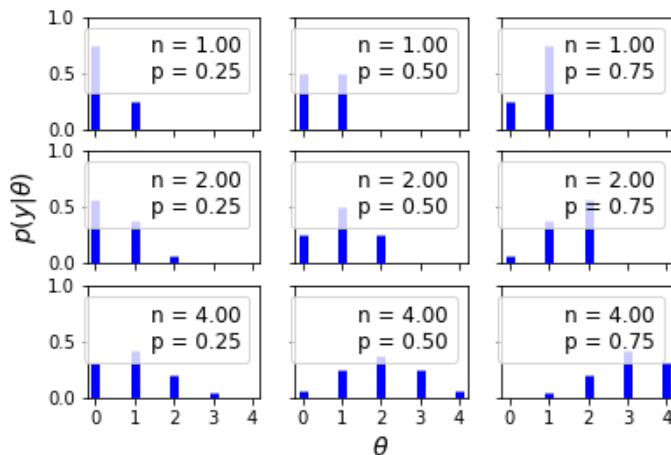


图 2.2 binomial_distribution

2.1.2.3 Choosing the prior

As a prior we will use a beta distribution (贝塔分布), which is a very common distribution in Bayesian statistics and looks like this:

$$p(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

The first term is a normalization constant that ensures the distribution integrates to 1.

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from scipy import stats
4 import seaborn as sns
5
6 params = [0.5, 1, 2, 3]
7 x = np.linspace(0, 1, 100)
8 f, ax = plt.subplots(len(params), len(params), sharex=True, sharey=True)
9 for i in range(4):
10     for j in range(4):
11         a = params[i]
12         b = params[j]
13         y = stats.beta(a, b).pdf(x)
14         ax[i,j].plot(x, y)
15         ax[i,j].plot(0,0,label="$\\alpha$ = {:.2f} \n $\\beta$ = {:.2f}".format(a,b), alpha=0)
16         ax[i,j].legend(fontsize=12)
17 ax[3,0].set_xlabel('$\\theta$', fontsize=14)
18 ax[0,0].set_ylabel('$p(\\theta)$', fontsize=14)

```

```
plt.savefig('beta_distribution.png')
```

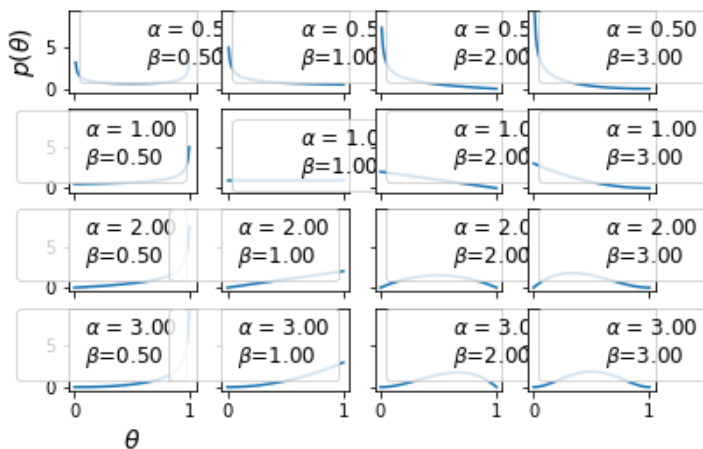


图 2.3 beta_distribution

Why are we using the beta distribution for our model? 1) One reason is that the beta distribution is restricted to be between 0 and 1, in the same way our parameter θ is. 2) Another reason is its versatility (通用性). As we can see in the preceding figure, the distribution adopts several shapes, including a uniform distribution, Gaussian-like distributions, U-like distributions, and so on. 3) A third reason is that the beta distribution is the conjugate prior (共轭先验) of the binomial distribution (which we are using as the likelihood). A conjugate prior of a likelihood is a prior that, when used in combination with the given likelihood, returns a posterior with the same functional form as the prior. There are other pairs of conjugate priors, for example, the Gaussian distribution is the conjugate prior of itself.

For many years, Bayesian analysis was restricted to the use of conjugate priors. Conjugacy ensures mathematical tractability of the posterior, which is important given that common problem in Bayesian statistics is to end up with a posterior we cannot solve analytically. This was a deal breaker before the development of suitable computational methods to solve any possible posterior.

However, modern computational methods to solve Bayesian problems whether we choose conjugate priors or not.

2.1.2.4 Getting the posterior

The Bayes' theorem says that the posterior is proportional to the likelihood times the prior:

$$p(\theta|y) \propto p(y|\theta)p(\theta)$$

which turns out to be

$$p(\theta|y) \propto \frac{N!}{N!(N-y)!} \theta^y (1-\theta)^{N-y} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

To our practical concerns we can drop all the terms that do not depend on θ and our results will still be valid. So we can write the following:

$$p(\theta|y) \propto \theta^y (1-\theta)^{N-y} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

Reordering it, we have

$$p(\theta|y) \propto \theta^{\alpha-1+y}(1-\theta)^{\beta-1+N-y}$$

We will see that this expression has the same functional form of a beta distribution (except for the normalization) with $\alpha_{\text{posterior}} = \alpha_{\text{prior}} + y$ and $\beta_{\text{posterior}} = \beta_{\text{prior}} + N - y$, which means that the posterior for our problem is the beta distribution:

$$p(\theta|y) = \text{Beta}(\alpha_{\text{prior}} + y, \beta_{\text{prior}} + N - y)$$

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from scipy import stats
4 import seaborn as sns
5
6 theta_real = 0.35
7 trials = [0, 1, 2, 3, 4, 8, 16, 32, 50, 150]
8 data = [0, 1, 1, 1, 1, 4, 6, 9, 13, 48]
9
10 beta_params = [(1,1), (0.5,0.5), (20,20)]
11 dist = stats.beta
12 x = np.linspace(0, 1, 100)
13
14 for idx, N in enumerate(trials):
15     if idx == 0:
16         plt.subplot(4, 3, 2)
17     else:
18         plt.subplot(4, 3, idx+3)
19     y = data[idx]
20     for (a_prior, b_prior), c in zip(beta_params, ('b', 'r', 'g')):
21         p_theta_given_y = dist.pdf(x, a_prior+y, b_prior+N-y)
22         plt.fill_between(x, 0, p_theta_given_y, color=c, alpha=0.6)
23
24     plt.axvline(theta_real, ymax=0.3, color='k')
25     plt.plot(0,0,label="{ :d} experiments\n{ :d} heads".format(N,y),alpha=0)
26     plt.xlim(0,1)
27     plt.ylim(0,12)
28     plt.xlabel(r'$\theta$')
29     plt.legend()
30     plt.gca().axes.get_yaxis().set_visible(False)
31 plt.tight_layout()
32 plt.savefig("posterior.png")

```

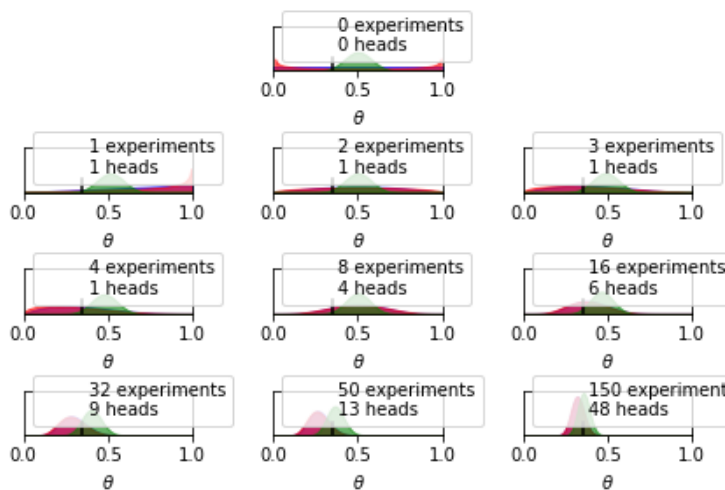


图 2.4 posterior

2.1.2.5 Model notation and visualization

A common notation to succinctly represent probabilistic is as follows:

$$\checkmark \theta \sim \text{Beta}(\alpha, \beta)$$

$$\checkmark y \sim \text{Bin}(n = 1, p = \theta)$$

2.1.3 Inference engines

There are several methods to compute the posterior even when it is not possible to solve it analytically. Some of the methods are:

1. Non-Markovian methods:

1.1 Grid computing

1.2 Quadratic approximation

1.3 Variation methods

2. Markovian methods:

2.1 Metropolis-Hastings

2.2 Hamiltonian Monte Carlo/No Y-Turn Sampler

Nowadays, Bayesian analysis is performed mainly by using Markov Chain Monte Carlo (MCMC) methods, with variational methods gaining momentum for bigger datasets. We do not need to really understand these methods to perform Bayesian analysis, that's the whole point of probabilistic programming languages, but knowing at least how they work at a conceptual level is often very useful.

2.1.3.1 Non-Markovian methods

Let's start our discussion of inference engines with the non-Markovian methods. These methods are in general faster than Markovian ones.

Grid computing

Grid computing is a brute-force approach. Even if you are not able to compute the whole posterior, you may be able to compute the prior and the likelihood for a given number of points. Let's assume we want to compute the posterior for a single parameter model. The grid approximation is as follows:

1. Define a reasonable interval for the parameter (the prior should give you a hint).
2. Place a grid of points (generally equidistant) on that interval.
3. For each point in the grid we multiply the likelihood and the prior.

Optionally, we may normalize the computed values (divide the result at each point by the sum of all points).

It is easy to see that a larger number of points (or equivalently a reduced size of the grid) will result in a better approximation. In fact, if we take an infinite number of points we will get the exact posterior. The grid approach does not scale well for many parameters (also referred as dimensions); as you increase the number of parameters the volume of the posterior gets relatively smaller compared with the sampled volume. In other words, we will spend most of the time computing values with an almost null contribution to be posterior, making this approach unfeasible for many statistical and data problems.

The following code implements the grid approach to solve the coin-flipping problem.

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from scipy import stats
4 import seaborn as sns
5
6 def posterior_grid_approx(grid_points=100, heads=6, tosses=9):
7     """
8     A grid implementation for the coin-flip problem
9     """
10    grid = np.linspace(0, 1, grid_points)
11    prior = np.repeat(5, grid_points)
12    likelihood = stats.binom.pmf(heads, tosses, grid)
13    unstd_posterior = likelihood * prior
14    posterior = unstd_posterior / unstd_posterior.sum()
15    return grid, posterior
16
17 # Assuming we made 4 tosses and we observe only 1 head we have the following:
18
19 points = 15
20 h, n = 1, 4

```

```

21 grid, posterior = posterior_grid_approx(points, h, n)
22 plt.plot(grid, posterior, 'o-', label='heads = {} \ ntosses = {}'.format(h, n))
23 plt.xlabel(r'$\theta$')
24 plt.legend(loc=0)
25 plt.savefig('grid_approach.png')

```

$$\sqrt{\theta} \sim \text{Beta}(\alpha, \beta)$$

$$\sqrt{y} \sim \text{Bin}(n = 1, p = \theta)$$

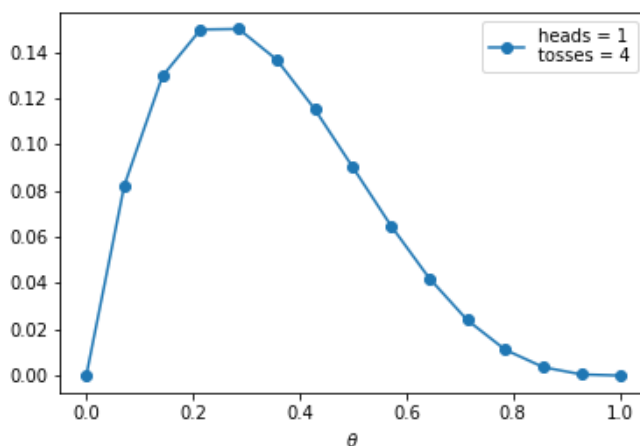


图 2.5 grid_approach to solve the coin flipping

Quadratic method The quadratic approximation, also known as the Laplace method or the normal approximation, consists of approximating the posterior with a Gaussian distribution. This method often works because in general the region close to the mode of the posterior distribution is more or less normal, and in fact in many cases is actually a Gaussian distribution. This method consists of two steps. First, find the mode of the posterior distribution. This method consists of two steps. First, find the mode of the posterior distribution. This can be done using optimization methods; that is, methods to find the maximum or minimum of a function, and there are many off-the-shelf methods for this purpose. This will be the mean of the approximating Gaussian. Then we can estimate the curvature of the function near the mode. Based on this curvature, the standard deviation of the approximating Gaussian can be computed. We are going to apply this method once we have introduced PyMC3.

Variational methods Most of modern Bayesian statistics is done using Markovian methods (see the next section), but for some problems those methods can be too slow and they do not necessarily parallelize well. The naive approach is to simply run n chains in parallel and then combine the results, but for many problems this is not a really good solution. Finding effective ways of parallelizing them is an active research area.

Variational methods could be a better choice for large datasets (think big data) and/or for likelihoods that are too expensive to compute. In addition, these methods are useful for quick approximations to the posterior and as starting points for MCMC methods.

The general idea of variational methods is to approximate the posterior with a simpler distribution; this may sound similar to the Laplace approximation, but the similarities vanish when we check the details of the method. The main drawback of variational methods is that we must come up with a specific algorithm for each model, so it is not really a universal inference engine, but a model-specific one.

Of course, lots of people have tried to automatize variational methods. A recently proposed method is the automatic differentiation variational inference (ADVI) (see <http://arxiv.org/abs/1603.00788>). At the conceptual level, ADVI works in the following way:

1. Transform the parameters to make them live in the real line. For example, taking the logarithm of a parameter restricted to positive values we obtain an unbounded parameter on the interval $[-\infty, \infty]$.
2. Approximate the unbounded parameters with a Gaussian distribution. Notice that a Gaussian on the transformed parameter space is non-Gaussian on the original parameter space, hence this is not the same as the Laplace approximation.
3. Use an optimization method to make the Gaussian approximation as close as possible to the posterior. This is done by maximizing a quantity known as the Evidence LowerBound (ELBO). How we measure the similarity of two distributions and what ELBO is exactly, at this point, is a mathematical detail.

Markovian methods There is a family of related methods collectively known as MCMC methods. As with the grid computing approximation, we need to be able to compute the likelihood and prior for a given point and we want to approximate the whole posterior distribution. MCMC methods outperform the grid approximation because they are designed to spend more time in higher probability regions than in lower ones. In fact, a MCMC method will visit different regions of the parameter space in accordance with their relative probabilities. If region A is twice as likely as region B, then we are going to get twice the samples from A as from B. Hence, even if we are not capable of computing the whole posterior analytically, we could use MCMC methods to take samples from it, and the larger the sample size the better the results.

What is in a name? Well, sometimes not much, sometimes a lot. To understand what MCMC methods are we are going to split the method into the two MC parts, the Monte Carlo part and the Markov Chain part.

Monte Carlo The use of random numbers explains the Monte Carlo part of the name. Monte Carlo methods are a very broad family of algorithms that use random sampling to compute or simulate a given process. Monte Carlo is a very famous casino located in the Principality of Monaco. One of the developers of the Monte Carlo method, Stanislaw Ulam, had an uncle who used to gamble there. The key idea Stan had was that while many problems are difficult to solve or even formulate in an exact way, they can be effectively studied by taking samples from them, or by simulating them. In fact, as the story goes, the motivation was to answer questions about the probability of getting a particular hand in a solitary game. One way to solve this problem was to follow the analytical combinatorial problem. Another way, Stan argued, was to play several games of solitaire and just count how many of the hands we play match the particular hand we are interested in! Maybe this sounds obvious, or at least pretty reasonable; for example, you may have used re-sampling methods to solve your statistical problems. But remember this mental experiment was

performed about 70 years ago, a time when the first practical computers began to be developed. The first application of the method was to solve a problem of nuclear physics, a problem really hard to tackle using the conventional tools at that time. Nowadays, even personal computers are powerful enough to solve many interesting problems using the Monte Carlo approach and hence these methods are applied to a wide variety of problems in science, engineering, industry, arts, and so on.

A classic pedagogical example of using a Monte Carlo method to compute a quantity of interest is the numerical estimation of π . In practice there are better methods for this particular computation, but its pedagogical value still remains. We can estimate the value of π with the following procedure:

1. Throw N points at random into a square of side $2R$.
2. Draw a circle of radius R inscribed in the square and count the number of points that are inside that circle.
3. Estimate $\hat{\pi}$ as the ration $\frac{4 \times \text{inside}}{N}$.

A couple of notes: We know a point is inside a circle if the following relation is true: $\sqrt{(x^2 + y^2)} \leq R$

The area of the square is $(2R)^2$ and the area of the circle is πR^2 . Thus we know that the ratio of the area of the square to be the area of the circle is $\frac{4}{\pi}$, and the area of the circle and square are proportional to the number of points inside the circle and the total N points, respectively.

Using a few line of Python we can run this simple Monte Carlo simulation and compute π and also the relative error of our estimate compared to the rule value of π :

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from scipy import stats
4 import seaborn as sns
5
6 N = 10000
7 x,y = np.random.uniform(-1, 1, size=(2, N))
8 inside = (x**2 + y**2) <= 1
9 pi = inside.sum() * 4 / N
10 error = abs((pi - np.pi) / pi) * 100
11
12 outside = np.invert(inside)
13
14 plt.plot(x[inside], y[inside], 'b.')
15 plt.plot(x[outside], y[outside], 'r.')
16 plt.plot(0, 0, label='$\hat{\pi}$ = {:.4f}\nerror = {:.3f}'.format(pi, error), alpha=0)
17 plt.axis('square')
18 plt.legend(frameon=True, framealpha=0.9, fontsize=16)
19
20 plt.savefig('monte_carlo.png')
```

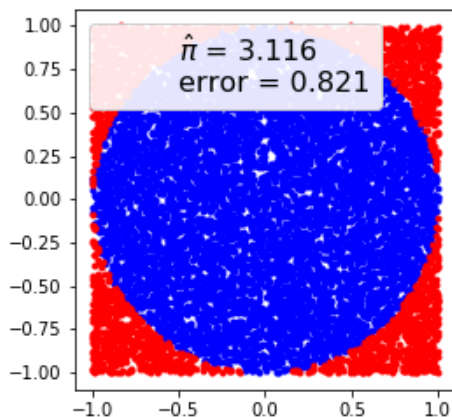



图 2.6 monte_carlo

In the preceding code we can see that the outside variable is only used to get the plot; we do not need it for computing $\frac{4 \times \text{inside}}{N}$. Another clarification: because our computation is restricted to the unit circle we can omit computing the square root from the computation of the inside variable.

Markov chain A Markov chain is a mathematical object consisting of a sequence of states and a set of probabilities describing the transitions among those states. A chain is Markovian if the probability of moving to other states depends only on the current state. Given such a chain, we can perform a random walk by choosing a starting point and moving to other states following the transition probabilities. If we somehow find a Markov chain with transitions proportional to the distribution we want to sample from (the posterior distribution in Bayesian analysis), sampling becomes just a matter of moving between states in this chain. So, how do we find this chain if we do not know the posterior in the first place? Well, there is something known as **detailed balance condition**. Intuitively, this condition says that we should move in a reversible way (a reversible process is a common approximation in physics). That is, the probability of being in state i and moving to state j should be the same as the probability of being in state j and moving towards state i .

In summary, all this means that if we manage to create a Markov Chain satisfying detailed balance we can sample from that chain with the guarantee that we will get samples from the correct distribution. This is a truly remarkable result! The most popular method that guarantees detailed balance is the **Metropolis-Hasting algorithm**.

Metropolis-Hastings To conceptually understand this method, we are going to use the following analogy. Suppose we are interested in finding the volume of water a lake contains and which part of the lake has the deepest point. The water is really muddy so we can't estimate the depth just by looking to the bottom, and the lake is really big, so a grid approximation does not seem like a very good idea. In order to develop a sampling strategy, we seek help from two of our best friends, Markovia and Monty. After some discussion they come up with the following algorithm that requires a boat; nothing fancy, we can even use a wooden raft, and a very long stick. This is cheaper than a sonar and we have already spent all our money on the boat, anyway!

1. Initialize the measuring by choosing a random place in the lake and move the boat there.

2. Use the stick to measure the depth of the lake.
3. Move the boat to some other point and take a new measurement.
4. Compare the two measures in the following way:
 - 4.1. If the new spot is deeper than the old one, write down in your notebook the depth of the new spot and repeat from 2.
 - 4.2 If the spot is shallower than the old one, we have two options: to accept or reject. Accepting means to write down the depth of the new spot and repeat from 2. Rejecting means to go back to the old spot and write down (again) the value for the depth of the old spot.

How do we decide to accept or reject a new spot? Well, the trick is to apply the Metropolis-Hastings criteria. This means to accept the new spot with a probability that is proportional to the ratio of the depth of the new and old spots.

If we follow this iterative procedure, we will get not only the total volume of the lake and the deepest point, but we will also get an approximation of the entire curvature of the bottom of the lake. As you may have already guessed, in this analogy the curvature of the bottom of the lake is the posterior distribution and the deepest point is the mode. According to our friend Markovia, the larger the number of iterations the better the approximation.

Indeed, theory guarantees that under certain general circumstances, we are going to get the exact answer if we get an infinite number of samples. Luckily for us, in practice and for many, many problems, we can get a very accurate approximation using a relatively small number of samples.

Let's look at the method now in a little bit more formal way. For some distributions, like the Gaussian, we have very efficient algorithms to get samples from, but for some other distributions such as many of the posterior distributions, we are going to find this is not the case. Metropolis-Hastings enables us to obtain samples from any distribution with probability $p(x)$ given that we can compute at least a value proportional to it. This is very useful since in a lot of problems like Bayesian statistics the hard part is to compute the normalization factor, the denominator of the Bayes' theorem. The Metropolis-Hastings algorithm has the following steps:

1. Choose an initial value for our parameter x_i . This can be done randomly or by using some educated guess.
2. We choose a new parameter value x_{i+1} , sampling from an easy-to-sample distribution such as a Gaussian or uniform distribution $Q(x_{i+1}|x_i)$. We can think of this step as perturbing the state x_i somehow.
3. We compute the probability of accepting a new parameter value by using the Metropolis-Hastings criteria $p_\alpha(x_{i+1}|x_i) = \min \left(1, \frac{p(x_{i+1})q(x_i|x_{i+1})}{p(x_i)q(x_{i+1}|x_i)} \right)$.
4. If the probability computed on 3 is larger than the value taken from a uniform distribution on the interval $[0, 1]$ we accept the new state, otherwise we stay in the old state.
5. We iterate from 2 until we have enough samples. Later we will see what enough means.

A couple of things to take into account:

1. If the proposal distribution $Q(x_{i+1}|x_i)$ is symmetric we get $p_\alpha(x_{i+1}|x_i) = \min \left(1, \frac{p(x_{i+1})}{p(x_i)} \right)$, often referred to as **Metropolis criteria** (we drop the **Hastings** apart).

2. Steps 3 and 4 imply that we will always accept or move to a most probable state, to a most probable parameter value. Less probable parameter values are accepted probabilistically given the ratio between the probability of the

new parameter value x_{i+1} and the old parameter value x_i . This criteria for accepting steps gives us a more efficient sampling approach compared to the grid approximation, while ensuring a correct sampling.

3. The target distribution (the posterior distribution in Bayesian statistics) is approximated by saving the sampled (or visited) parameter values. We save a sampled value x_{i+1} if we accept moving to a new state x_{i+1} . If we reject moving to x_{i+1} , we save the value of x_i .

At the end of the process we will have a list of values sometimes refereed to as a **sample chain or trace**. If everything was done the right way these samples will be an approximation of the posterior. The most frequent values in our trace will be the most probable values according to the posterior. An advantage of this procedure is that analyzing the posterior is simple. We have effectively transformed integrals (of the posterior) into just summing values in our vector of sampled values. The following code illustrates a very basic implementation of the Metropolis algorithm. Is not meant to solve any real problem, only to show it is possible to sample from a function if we know how to compute its value at a given point. Notice also that the following implementation has nothing Bayesian in it; there is no prior and we do not even have data! Remember that the MCMC methods are very general algorithms that can be applied to a broad array of problems. For example, in a (non-Bayesian) molecular model, instead of `func.pdf(x)` we would have a function computing the energy of the system for the state x .

The first argument of the metropolis function is a SciPy distribution; we are assuming we do not know how to directly get samples from this distribution.

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from scipy import stats
4 import seaborn as sns
5
6 import matplotlib.pyplot as plt
7 import numpy as np
8 from scipy import stats
9 import seaborn as sns
10
11 def metropolis(func, steps=10000):
12     """
13     A very simple Metropolis implementation.
14     """
15     samples = np.zeros(steps)
16     old_x = func.mean()
17     old_prob = func.pdf(old_x)
18
19     for i in range(steps):
20         new_x = old_x + np.random.normal(0, 0.5)
21         new_prob = func.pdf(new_x)
22         acceptance = new_prob / old_prob

```

```

23         if acceptance >= np.random.random():
24             samples[i] = new_x
25             old_x = new_x
26             old_prob = new_prob
27         else:
28             samples[i] = old_x
29     return samples
30
31 """
32 In the next example we have defined func as a beta function,
33 simply because is easy to change their parameters and get
34 different shapes. We are plotting the samples obtained by
35 metropolis as a histogram and also the True distribution as
36 a red line:
37 """
38
39 func = stats.beta(0.4, 2)
40 samples = metropolis(func=func)
41 x = np.linspace(0.01, .99, 100)
42 y = func.pdf(x)
43 plt.xlim(0, 1)
44 plt.plot(x, y, 'r-', lw=3, label='True distribution')
45 plt.hist(samples, bins=30, normed=True, label='Estimated distribution')
46 plt.xlabel(' $x$ ', fontsize=14)
47 plt.ylabel(' $p_{\text{pdf}}(x)$ ', fontsize=14)
48 plt.legend(fontsize=14)
49 plt.savefig('metropolis_algorithm.png')

```

Hamiltonian Monte Carlo/NUTS MCMC methods, including Metropolis-Hastings, come with the theoretical guarantee that if we take enough samples we will get an accurate approximation of the correct distribution. However, in practice it could take more time than we have to get enough samples. For that reason, alternatives to the general Metropolis-Hastings algorithm have been proposed. Many of those alternative methods such as the Metropolis-Hastings algorithm itself, were developed originally to solve problems in statistical mechanics, a branch of physics that studies properties of atomic and molecular systems. One such modification is known as **Hamiltonian Monte Carlo or Hybrid Monte Carlo (HMC)**. In simple terms a Hamiltonian is a description of the total energy of a physical system. The name Hybrid is also used because it was originally conceived as a hybridization of molecular mechanics, a widely used simulation technique for molecular systems, and Metropolis-Hastings. The HMC method is essentially the same as Metropolis-Hastings except that instead of proposing random displacements of our boat we do something more clever instead; we move the boat following the curvature of the lake's bottom. Why is this clever? Because in doing so we try to avoid one of the main problems of Metropolis-Hastings: the exploration is slow and samples tend

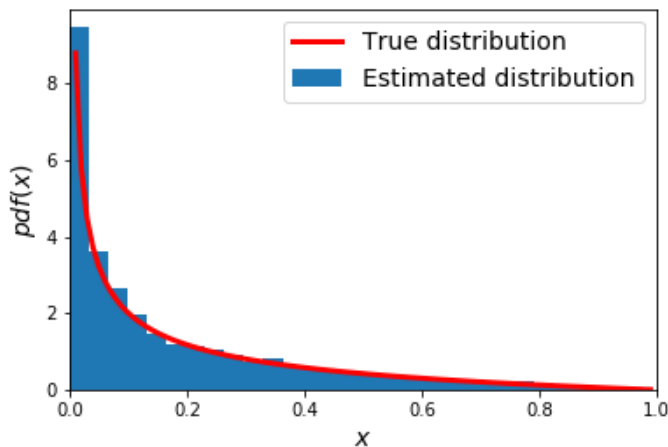


图 2.7 metropolis_algorithm

to be autocorrelated, since most of the proposed moves are rejected.

So, how can we try to understand this method without going into mathematical details? Imagine we are in the lake with our boat. In order to decide where to move next we let a ball roll at the bottom of the lake, starting from our current position. Remember that this method was brought to us by the same people that treat horses as spheres, so our ball is not only perfectly spherical, it also has no friction and thus is not slowed down by the water or mud. Well, we throw a ball and we let it roll for a short moment, and then we move the boat to where the ball is. Now we accept or reject this step using the Metropolis criteria just as we saw with the Metropolis-Hastings method. The whole procedure is repeated a number of times. This modified procedure has a higher chance of accepting new positions, even if they are far away relative to the previous position.

Out of our Gedanken experiment and back to the real world, the price we pay for this much cleverer Hamiltonian-based proposal is that we need to compute gradients of our function. A gradient is just a generalization of the concept of derivative to more than one dimension. We can use gradient information to simulate the ball moving in a curved space. So, we are faced with a trade-off; each HMC step is more expensive to compute than a Metropolis-Hastings one but the probability of accepting that step is much higher with HMC than with Metropolis. For many problems, this compromise turns in favor of the HMC method, especially for complex ones. Another drawback with HMC methods is that to have really good sampling we need to specify a couple of parameters. When done by hand it takes some trial and error and also requires experience from the user, making this procedure a less universal inference engine than we may want. Luckily for us, PyMC3 comes with a relatively new method known as **No-U-Turn Sampler (NUTS)**. This method has proven very useful in providing the sampling efficiency of HMC methods, but without the need to manually adjust any knob.

Other MCMC methods There are plenty of MCMC methods out there and indeed people keep proposing new methods, so if you think you can improve sampling methods there is a wide range of persons that will be interested in your ideas. Mentioning all of them and their advantages and drawbacks is completely out of the scope of this book. Nevertheless, there are a few worth mentioning because you may hear people talk about them, so it is nice

to at least have an idea of what are they talking about.

Another sampler that has been used extensively for molecular systems simulations is the **Replica Exchange** method, also known as **parallel tempering** or **Metropolis Coupled MCMC** (or MC3; maybe that's too many MCs). The basic idea of this method is to simulate different replicas in parallel. Each replica follows the Metropolis-Hastings algorithm. The only difference between replicas is that the value of a parameter called temperature (physics influence once more time!) controls the probability of accepting less probable positions. From time to time, the method attempts a swap between replicas. The swapping is also accepted/rejected according to the Metropolis-Hastings criteria, but this time taking into account both replicas' temperatures. The swapping between chains can be attempted between random chains but it is generally preferable to do it for neighboring replicas; that is, replicas with similar temperatures and hence a higher probability-of-acceptance ratio. The intuition for this method is that as we increase the temperature the probability of accepting the new proposed position increases, and decreases with lower and lower temperatures. Replicas at higher temperatures explore the system more freely; for these replicas the surface becomes effectively flatter and thus easier to explore. For a replica with infinite temperature, all states are equally likely. The exchange between replicas avoids replicas at low temperatures getting trapped in local minima. This method is well suited for exploring systems with multiple minima.

2.1.4 PyMC3 introduction

第3章 深度学习

Goals to Achieve

1. pytorch basics
2. pytorch projects

§ 3.1 Pytorch

打印模型结构

```
pip install torchsummary
```

```
summary(model, (3, 32, 32))
```

```
print(model)
```