

目录

第一章 MATLAB	2
1.1 内置函数	2
1.2 MATLAB 包推荐	2
1.3 MATLAB 数据结构	2
1.3.1 table	2
1.3.2 struct	3
1.3.3 cell	3
1.3.4 handle	3
1.4 如何遍历当前文件夹及其子文件夹中的全部文件?	3
1.5 聊一聊自然排序	5
1.6 如何给 struct 排序?	5
1.7 如何隔行取数据?	6
1.8 如何在遍历数组的同时删除被遍历过的元素?	6
1.9 再谈向量化操作	7
1.10 文件读取	8
1.11 如何将两个维度不一致的矩阵串联起来?	8
1.12 数组/矩阵去重	9
1.13 颜色的处理	9
1.14 谈谈 MATLAB 中的图形对象	9
1.14.1 图形对象	9
1.14.2 格式和注释	11
1.14.3 子图	11
1.14.4 最佳实践	11
1.14.5 谈一谈 colorbar	12
1.14.6 如何给 bar 加图例等说明性对象	12
1.14.7 如何导出大小和图形窗口大小一致的 pdf 文件?	13
1.15 函数最佳实践	13
1.16 关于路径地一切	13
1.17 如何在矩阵中插入一行/一列	14
1.18 图形处理	14
1.18.1 基本法	14
1.18.2 图形插值	14
1.19 优雅的数据流处理	15
1.20 如何更新绘图?	15
1.21 如何给不断更新的绘图制作 gif?	15

1.22 更优雅的颜色表	16
1.23 关于箱型图	16
第二章 Python	17
2.1 如何展开一个嵌套的序列?	17
2.2 如何遍历当前文件夹及其子文件夹中的全部文件?	19
2.3 如何在遍历 list 时删除元素?	19
2.4 字符串的处理	20
2.4.1 字符串自然排序	20
2.5 关于路径的一切	20
2.5.1 实例: 遍历文件夹内的 PDF 文件并提取第一页合并	21
2.5.2 pathlib	22
2.5.3 实例: 添加父文件夹名到文件名	22
2.6 图像处理	23
2.6.1 实例: 给图片添加文字	23
2.7 最佳实践	23
2.7.1 代码的模块化	23
2.8 matplotlib 最佳实践	24
2.8.1 生成矢量图	24
2.9 并行执行 shell 命令	24
2.10 协助 ffmpeg 进行批量字幕压制	25
第三章 C 和 C++	27
3.1 C 语言的动态数组	27
第四章 算法	29
4.1 排序	29
4.1.1 冒泡排序	29
4.1.2 插入排序	30
4.1.3 归并排序	31
4.1.4 选择排序	32
4.1.5 快速排序	33
第五章 Git	36
5.1 如何给 Git 仓库添加一个空文件夹?	36

导言

这份文档主要用来存放一些实际工作中碰到的实用的代码片段，可能包含 MATLAB、Python、C/C++ 和一些 \TeX 的小知识。个人笔记，个人娱乐。

如果有人想编译这份手册或想学习一下实现，请务必读以下说明。

字体设置，为了避免侵权，尽可能使用开源字体^①。

- Source Han Sans: <https://github.com/adobe-fonts/source-han-sans/tree/release>
- Source Han Serif: <https://github.com/adobe-fonts/source-han-serif/tree/release>
- Source Code Pro: <https://github.com/adobe-fonts/source-code-pro>
- PT Sans Narrow: <https://fonts.google.com/specimen/PT+Sans+Narrow>
- TeX Gyre: 有问题前往<https://www.ctan.org>获取，一般来说 \TeX 发行版自带
- 等宽字体：大多数等宽字体都是程序员使用的，开源居多，颇易获取。我常用 DejaVu Sans Mono, Fira Code 和 Source Code Pro 三种。

```
%% 字体设置
\usepackage{fontspec}
\setmainfont{Adobe Garamond Pro} % TeX Gyre Pagella
\setsansfont{TeX Gyre Heros}
\setmonofont{Source Code Pro} % Consolas, DejaVu Sans Mono
\setCJKmainfont[BoldFont={Source Han Sans SC}, ItalicFont={KaiTi}]{Source Han Serif
↪ SC}
\setCJKmonofont{FangSong}
\setCJKsansfont{Source Han Sans SC}

%% 数学字体
\usepackage{unicode-math}
\setmathfont[math-style = ISO, bold-style = ISO]{TeX Gyre Pagella Math}

%% url 样式
\newfontfamily\urlfont{PT Sans Narrow}
```

编译环境设置，代码高亮环境由 minted 宏包提供（需要 Python 环境）。

代码测试环境，各种代码的运行环境为 MATLAB 2017b、Anaconda、Visual Studio 2017 community、 \TeX （各宏包均为最新）。

如果你觉得本文档里的代码有用，请不要直接复制文档里面的代码（直接复制会复制到换行产生的符号及空格，可能导致代码出现难以预计的错误），请到[github 项目](#)的 code 文件夹找对应的文件。

^① 西文主字体 Adobe Garamond Pro、楷体、仿宋暂时没有找到理想的替代方案

第一章 MATLAB

1.1 内置函数

在 MATLAB 中工作时，很多操作其实都是有内置函数的，对 MATLAB 不熟悉就用不成，然后就“曲线救国”了，效率不高且不说，关键自己实现很费脑！这里收集一些数值计算工作中极其常用的函数。

`meshgrid`

1.2 MATLAB 包推荐

MATLAB 相比 Python 最大的弱点就是没有包管理器，在此收集一些常用的包。

1. MBeautifier, 代码格式化包, 代码必须要看起来美, 哪怕写得很烂, <https://github.com/davidvarga/MBeautifier>
2. 遗传算法工具箱, 比内置的逻辑清晰, 可能是我脑回路奇怪, <http://codem.group.shef.ac.uk/index.php/ga-toolbox>
3. 作者撰写了大量自然排序的工具箱, 爱不释手啊, <https://www.mathworks.com/matlabcentral/fileexchange/?term=profileid:3102170>

1.3 MATLAB 数据结构

1.3.1 table

`table` 是个很有意思的数据结构, 如果不指定行索引, 就跟 `struct` 类似, 指定行索引, 就像 Excel 电子表格了, 极为方便。使用时需要声明。

任何数据都能拿来创建 `table`, 但是行数必须相同。`table` 含三个方法, `Properties`, `Row`, `Variables`, 分别用来查看表的属性、查看行名称、拼接所有数据并以矩阵形式返回。

```
LastName = {'Sanchez'; 'Johnson'; 'Li'; 'Diaz'; 'Brown'};
Age = [38; 43; 38; 40; 49];
Smoker = logical([1; 0; 1; 0; 1]);
Height = [71; 69; 64; 67; 64];
Weight = [176; 163; 131; 133; 119];
BloodPressure = [124, 93; 109, 77; 125, 83; 117, 75; 122, 80];

T = table(Age, Smoker, Height, Weight, BloodPressure);
```

表的列索引是按变量名称默认生成的, 当然也可以主动指定。使用 `table` 传递参数或者修改名为 `VariableNames` 的属性。同样, 指定行名称的属性为 `RowNames`, 指定名称属性时使用 `cell` 或者 `string` 数据结构。使用 `join`, `innerjoin`, `outerjoin` 来合并不同的表。

1.3.2 struct

非常经典的整合不同数据的数据结构。使用时不需要声明，直接用 `.` 运算符就能原地创建 `struct`。

1.3.3 cell

使用 `[{}, {}]` 合并两个 `cell`，使用 `{{}, {}}` 将会创建嵌套的 `cell`。

1.3.4 handle

函数句柄是个相当有用的东西，某种程度上，相当于 MATLAB 中的指针。

1.4 如何遍历当前文件夹及其子文件夹中的全部文件？

假设现在我们有这样一个文件夹 A，它含有一些文件和子文件夹 B、C、D.....，这些子文件夹又包含若干层子文件夹。我们需要将这个父文件夹（A）及其子文件夹（B、C、D.....）和孙文件夹中的所有文件名和其路径取出来。

如果你用的是 MATLAB 2016b 及更新的版本，那真的太棒了！`dir()` 函数已经支持遍历搜索了。尝试敲入：

```
dir_data = dir('**/*');  
dir_data([dir_data.isdir]) = []; % 去除所有 . 和 .. 文件夹
```

这将会返回一个包含文件信息的 `struct`，现在你可以任意操作这些 `struct` 了，随意拼接路径。解放大脑，哦也！方便归方便，但是，一来肯定有大多数人使用的是 MATLAB 2016b 之前的版本，二来，解放大脑意味着我失去了一次独立思考的机会。

思考

对于实现方法^①，多层次的遍历，我第一时间想到的是递归。然后就是数据的存储了，`dir()` 函数返回的是一个 `struct`，这个数据结构储存有文件的信息，我们要充分利用这个数据结构。所以现在思路是，写一个递归函数，这个函数返回包含所有文件信息的 `struct`。

这个函数应对先处理父文件夹，获取文件和子文件夹，然后储存文件信息，同时去除子文件夹中的 `'.'` 和 `'..'` 这两个特殊文件夹。我们对获取的子文件夹再次调用该函数，并储存文件信息。如此，利用递归获取子子孙孙无穷尽文件夹的信息^②，最后函数返回存储有所有文件信息的 `struct`。现在，你可以对这个结构体做你想做的事情。

解

MATLAB 2016b 以上的版本我们可以用函数返回 `struct`，这个数据结构包含 `[folder, name, date, bytes, isdir, datenum]` 六个字段的信息，我们可以按自己意愿使用 `folder` 和 `name` 拼接出文件的完整路径。

```
% get all file name in current dir and sub dir, Compatible with MATLAB R2016b and newer  
function file_list = get_all_file_name_R2016b_newer(path)  
  
dir_data = dir(path);
```

^① 思路来源：[How to get all files under a specific directory in MATLAB?](#)

^② 其实这并不可能，因为递归是有栈高度限制的，调用函数压入栈，返回函数弹出栈，如果文件夹层次太深，一直压栈就会到达栈溢出警告的极限，例如 Python 的栈往往是 100 层，我想 MATLAB 的栈也大致如此，不会太高

```

file_list = dir_data(~[dir_data.isdir]); % file name of current dir

% get sub dir information
sub_dir = dir_data([dir_data.isdir]); % struct
dot_dir = ismember({sub_dir.name}, {'.', '..'}); % logical
sub_dir = sub_dir(~dot_dir); % struct, remove specific folder

% recursion
for i = 1:length(sub_dir)
    next_dir = fullfile(sub_dir(i).folder, sub_dir(i).name); % str
    file_list = [file_list; get_all_file_name_R2016b_newer(next_dir)]; % struct
end

end

```

MATLAB 2016a 及之前的版本 dir struct 信息并不包含 folder, 如果返回 struct, 将只有文件的 [name, date, bytes, isdir, datenum] 五个字段的信息, 所以我们并不能根据函数返回的 struct 拼接出文件完整路径, 我们需要自己将路径拼接成一个 cell, 然后使用函数返回 cell。

```

% get all file name in current dir and sub dir, Compatible with MATLAB R2016a and older
function file_list = get_all_file_name_R2016a_older(path)

% file name of current dir
dir_data = dir(path);
file_list_struct = dir_data(~[dir_data.isdir]);
file_list = fullfile(path, {file_list_struct.name})';

% get sub dir information
sub_dir = dir_data([dir_data.isdir]); % struct
dot_dir = ismember({sub_dir.name}, {'.', '..'}); % logical
sub_dir = sub_dir(~dot_dir); % struct, remove specific folder

% recursion
for i = 1:length(sub_dir)
    next_dir = fullfile(path, sub_dir(i).name); % str
    file_list = [file_list; get_all_file_name_R2016a_older(next_dir)]; % struct
end

end

```

总结

`dir()` 函数遍历整个 F 盘共 2 万余文件文件大约需要 1.555823s。我们实现的递归函数遍历 F 盘文件大约需要 3.703009s。慢是慢了点，但我们成功运用了递归解决问题，不是吗？

1.5 聊一聊自然排序

通常，我们会遇到处理一系列文件名有规律的文件的情况，比如：a1.txt、a2.txt a100.txt。但是，当读取文件名到一个 `cell` 里后，我们发现文件名往往是乱序排列的，甚至当你使用 `sort` 函数后，排序也不会改变。搜索了一下，在 Mathworks File Exchange 网站找到了一个自然排序的函数 `natsort`^③，感谢作者 Stephen Cobeldick。效果如下：

	1	2
1	a0.txt	
2	a1.txt	
3	a10.txt	
4	a11.txt	
5	a12.txt	
6	a13.txt	
7	a14.txt	
8	a15.txt	
9	a16.txt	
10	a17.txt	
11	a18.txt	
12	a19.txt	
13	a2.txt	
14	a20.txt	
15	a3.txt	

图 1.1: 乱序的文件名

```
{'a0.txt' }
{'a1.txt' }
{'a2.txt' }
{'a3.txt' }
{'a4.txt' }
{'a5.txt' }
{'a6.txt' }
{'a7.txt' }
{'a8.txt' }
{'a9.txt' }
{'a10.txt'}
{'a11.txt'}
{'a12.txt'}
{'a13.txt'}
{'a14.txt'}
{'a15.txt'}
{'a16.txt'}
{'a17.txt'}
{'a18.txt'}
{'a19.txt'}
{'a20.txt' }
```

图 1.2: 排序后自然顺序的文件名

另：作者还提供两个额外的函数，用于文件名自然排序的 `natsortfiles`^④，用于按行自然排序的 `natsortrows`^⑤。

这三个函数都能传入额外的参数，其中最有用的就是传入正则表达式，用以匹配文件名来对结构更复杂字符串进行排序。正则表达式比较复杂，多次匹配不成功都很正常，可以使用该作者^⑥提供的小工具反复检查。

有时候还需要给 `struct` 排序，而显然 `natsort` 函数只能处理 `cell` 数据结构。使用函数的返回值配合一些索引的技巧可以轻易做到给复杂结构体排序。

```
st_file_name = dir('test/*');
st_file_name([st_file_name.isdir]) = []; % delete . and .. folder
% natsort filename
[~, Index] = natsort({st_file_name.name});
st_file_name = st_file_name(Index);
```

1.6 如何给 struct 排序？

`struct` 是一种非常常见的数据结构，在对其进行操作的过程中免不了要排序，那么我们怎样以一个 `field`(域) 为基准给整个 `struct` 排序呢？

^③ <https://cn.mathworks.com/matlabcentral/fileexchange/34464-customizable-natural-order-sort>

^④ <https://www.mathworks.com/matlabcentral/fileexchange/47434-natural-order-filename-sort>

^⑤ <https://www.mathworks.com/matlabcentral/fileexchange/47433-natural-order-row-sort>

^⑥ <https://www.mathworks.com/matlabcentral/fileexchange/48930>

基本思想^⑦是先把 `fields`、`data` 从 `struct` 里面剥离出来，同时将 `data` 转化为 `cell` 数据类型以供排序，使用 `sortrows` 以指定的 `fields` 列为标尺进行排序，最后将结果转化为 `struct` 数据类型。如果比较懒的话，直接使用现成的脚本^⑧吧！

```
st_file_name = dir('*.txt');
fields_file_name = fieldnames(st_file_name); % extract fields
cell_file_name = struct2cell(st_file_name)'; % extract data
cell_file_name = sortrows(cell_file_name, 4, 'descend'); % sort by 4th column
st_file_name = cell2struct(cell_file_name', fields_file_name);
```

2019 年 6 月 11 日更新：最近发现一种更简单、更优雅的方法^⑨，故上述的方法可以作废了。新方法的基本思路是：利用 `sort` 函数获取排序的索引，而不是结果；然后直接利用该索引对 `struct` 中的数据进行索引，完成排序。

```
FileList = dir(fullfile(Folder, '*..*'));
 [~, Index] = natsort({FileList.name});
FileList = FileList(Index);
```

1.7 如何隔行取数据？

闭上眼睛，想象现在有这样一个数组 `[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]`，我们要隔一列取一个数据，或者隔两列取一个数据。得益于 MATLAB 的向量化编程，我们可以很方便的做到，

```
mat_a = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10];
mat_b = mat_a(:, 1:2:length(mat_a));
```

如果你用循环，那么你的代码就不优雅，另，向量化操作比循环快，大型数组优势明显。以上。

1.8 如何在遍历数组的同时删除被遍历过的元素？

闭上眼睛，想象现在有这样一个数组 `[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]`，我们需要边遍历元素边删除元素。实现方法和 Python 章节方法一致。

```
mat_a = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10];

while ~isempty(mat_a)
    fprintf("The element being traversed is %d\n", mat_a(1));
    mat_a(1) = [];
    disp(mat_a);
end
```

^⑦ <https://blogs.mathworks.com/pick/2010/09/17/sorting-structure-arrays-based-on-fields/>

^⑧ <https://www.mathworks.com/matlabcentral/fileexchange/28573-nestedsortstruct>

^⑨ <https://www.mathworks.com/matlabcentral/answers/389300-how-to-nat-sort-rows-of-a-struct>

1.9 再谈向量化操作

今天又碰到一个数组操作的问题，同样，如果用一般的方法来解决，代码是很冗长的，向量化操作再次助我一臂之力。

有一个 2 列的数组 `all_data`，第一列有正有负，我们称第一列为 x ，第二列为 y 。现在需要索引 $x > 0$ 时对应的 $[x, y]$ 为一个新的数组 `a`。并且需要从 `a` 中返回 $y = \min(y)$ 时所对应的数组 $[x_0, y_0]$ 。

```
... ..
-1.44319267634370e-06 9.80637785912817e-06
-1.68967863180042e-07 9.73806551956721e-06
-6.45218837777561e-07 9.75074561060079e-06
6.28923217787410e-07 9.75037059307950e-06
1.54045071473931e-07 9.73772289244816e-06
1.42591401762642e-06 9.80510552313044e-06
... ..
```

先谈向量化获取数组 `a`，利用逻辑索引，保证 x 全大于 0，并取出 1、2 两列；然后利用 `find` 获取 $y = \min(y)$ 的行索引；最后利用索引轻松找到需要的数据。可能看起来比较难理解，但是此时再在外面套文件操作的循环等循环操作是不是清晰多了。

```
a = all_data(all_data(:, 1) > 0, 1:2);
```

```
[r, ~] = find(a(:, 2) == min(a(:, 2)));
what_is_i_need = a(r, c);
```

```
clear
clc

file_name_struct = dir('./0518*.txt');
file_name = {file_name_struct(:).name};
file_name = natsort(file_name);
what_is_i_need = [];

for file_num = 1:length(file_name)
    all_data = load(file_name{file_num});
    a = all_data(all_data(:, 1) > 0, 1:2);

    [r, ~] = find(a(:, 2) == min(a(:, 2))); %
    what_is_i_need = [what_is_i_need; a(r, 1), a(r, 2)];
end

plot(what_is_i_need(:, 2))
```

再来记录一个问题：循环操作里面有一个的 2 列数组 `all_data`，每次循环取第一列中与 0.5 最接近的数据和对应的列，所以该数组大小会不断变，设其维度为 $n \times 2$ 。如果 $n < 2$ ，我们将数据置为 0 并保存到

一个新数组里面去；如果 $n \geq 2$ ，保存其最小值和最大值到新数组里面去。同样，利用向量化操作最大程度减少代码量。

```
pos = [];
for i = 1:length(time)
    % [x, phi]
    all_data = load(strcat(mph_file, '\', num2str(i), '.txt'));
    all_data = all_data(abs(all_data(:, 2)-0.5) < 0.01, :);
    [r, c] = size(all_data);
    if r == 0 || r == 1
        x_min = 0;
        x_max = 0;
        phi = 0;
        pos = [pos; time(i), x_min, phi; time(i), x_max, phi];
    else
        x_min = all_data(all_data(:, 1) == min(all_data(:, 1)));
        x_max = all_data(all_data(:, 1) == max(all_data(:, 1)));
        [r, ~] = find(all_data == x_min);
        phi_min = all_data(r, 2);
        [r, ~] = find(all_data == x_max);
        phi_max = all_data(r, 2);
        pos = [pos; time(i), x_min, phi_min; time(i), x_max, phi_max];
    end
end
```

1.10 文件读取

`csvread` 适合读取纯 Comma-Separated Values 文件，`load` 适合读取带注释的 Comma-Separated Values 文件（示例如下），其在读取过程中会自动忽略 csv 文件的注释。

% x	y	IsoLevel
-1.348651530446577E-5	1.798983884698175E-5	0.5
-1.4987361701783775E-5	2.4219367655756994E-5	0.5
-1.494145158530118E-5	2.3443649068772022E-5	0.5
...

1.11 如何将两个维度不一致的矩阵串联起来？

现有矩阵 $a = [1]$ 和矩阵 $b = [5; 9; 4; 4]$ ，将其横向拼接成一个矩阵 c ，

```
c = [1, 5;
     1, 9;
     1, 4;
     1, 4]
```

思路很简单，因为 `a` 的维度不够，所以将其扩维，然后拼接。

```
a = [1];  
b = [5; 9; 4; 4];  
[r, ~] = size(b);  
c = [repmat(a, r, 1), b];
```

1.12 数组/矩阵去重

假设现在有一个数组或矩阵，里面包含许多重复的项目，我们需要合并这些重复的项目。MATLAB 提供一个名为 `unique` 的函数可以用来去重。很可惜，你用过就会发现，去重后数据的先后顺序乱了（MATLAB 给排序了）。借用排序算法的一个概念，这叫不稳定的去重。

其实我们可以在去重的同时记录数据第一次出现的位置，然后根据这个位置把去重后的数据再次排序即可。

```
...  
tmp = load(file_name_str(num).name);  
[data, pos] = unique(tmp(:, 2), 'first');  
data = sortrows([pos, data]);
```

最后，忙活了半天，发现 `unique` 函数有一个 `'stable'` 参数实现了稳定的去重，fuck you! 这件事情告诉我们，要认真读文档。当然，还是有收获的，起码学会了一个新的函数 `sortrows` 的用法。

1.13 颜色的处理

Hex 是一种常用的十六进制颜色码，MATLAB 并不能识别，从 MathWorks File Exchange 找了一个 `rgb2hex` and `hex2rgb` 函数^⑩，另，MathWorks File Exchange 真特么是个宝库，缺什么找什么，一找一个准。

有时候需要将 `double` 类型的矩阵转换成 `rgb` 色值的三维矩阵，没有 MATLAB 内置函数能做到，在 MathWorks File Exchange 上找了个 `double2rgb` 函数^⑪可以很方便的做到这一点。

1.14 谈谈 MATLAB 中的图形对象

可视化是一项很费时费力的工作，MATLAB 可视化成本更高，由于参数混杂，很难进行快速调整。我脑子不好使，先记录一下 MATLAB 基本图像元素的构成，更详细介绍请查阅 MATLAB 帮助文档^⑫。

1.14.1 图形对象

MATLAB 的图形系统中常用的对象分为：顶层对象（Root, Figure, Axes 等）、图表对象（如 Bar, Contour, Surface, Line 等）、插图对象（如 Colorbar, Legend），注释对象（如 Arrow, TextBox 等）和原始对象，函数对象，组对象，标尺对象等几种不常用对象。

我们重点关注两类对象，顶层对象 `Root`, `Figure`, `Axes` 和插图对象 `Colorbar`, `Legend`。图形对象相关的函数如图 1.3 所示，用于查找、复制和删除图形对象。这些操作推荐使用面向对象编程的方式进行，代码风格如下。

^⑩ <https://ww2.mathworks.cn/matlabcentral/fileexchange/46289-rgb2hex-and-hex2rgb>

^⑪ <https://www.mathworks.com/matlabcentral/fileexchange/30264-double2rgb>

^⑫ <https://ww2.mathworks.cn/help/matlab/graphics-objects.html>

<code>gca</code>	当前坐标区或图
<code>gcf</code>	当前图窗的句柄
<code>gcbf</code>	包含正在执行其回调的对象的图窗句柄
<code>gcbo</code>	正在执行其回调的对象的句柄
<code>gco</code>	当前对象的句柄
<code>groot</code>	图形根对象
<code>ancestor</code>	图形对象的父级
<code>allchild</code>	查找指定对象的所有子级
<code>findall</code>	查找所有图形对象
<code>findobj</code>	查找具有特定属性的图形对象
<code>findfigs</code>	查找可见的屏幕外图窗
<code>gobjects</code>	初始化图形对象的数组
<code>isgraphics</code>	对有效的图形对象句柄为 True
<code>ishandle</code>	测试是否有效的图形或 Java 对象句柄
<code>copyobj</code>	复制图形对象及其子级
<code>delete</code>	删除文件或对象

图 1.3: MATLAB 中的图形对象的标识

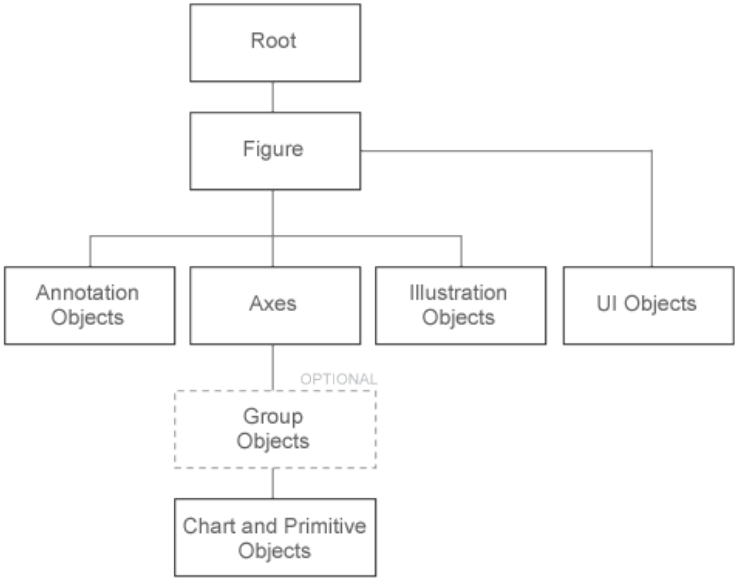


图 1.4: MATLAB 中的图形对象

```
r = groot;  
fig = figure;  
ax = gca;  
c = colorbar;  
lgs = legend('a','b','c');
```

`figure` 控制图窗窗口的外观和行为, `axes` 控制绘图区域外观和行为的对象, 图表对象 (如 `Bar`, `Line` 等) 才控制展示数据的图形的外观和行为, 这几个关系一定要拎得清。如, 你要修改图形窗口的大小和位置, 那就改 `figure` 对象的属性; 如果你要修改坐标轴的线条或者坐标轴标签, 就修改 `axes` 的属性; 最后, 改图形内容属性, 比如 `plot` 线条的粗细, 你才需要修改图表对象的属性。MATLAB 图形系统的设计层级关系非常明晰, 理清关系才能高效进行“码”操作 (笑)。

1.14.2 格式和注释

大部分调整图像格式和注释的函数都是直接修改图形对象的属性, 从代码可维护性的角度来看, 是没必要用这些函数来修改图形的属性的。但是有一些很奇怪的函数调整的不是图像对象的属性, 比如 `caxis` 函数用于修改 `colorbar` 的色值范围, 但是 `colorbar` 对象根本没有这个属性, 这就卧槽了! 有必要单独列出来记一下。

格式和注释分如下几类^⑬: 标题和标签、坐标区外观、颜色图、三维场景控制。

看, 是吧, 大部分内容直接修改图形对象的属性就可以了, 但这里面偏偏有特例, 可要注意了!

1.14.3 子图

`subplot` 函数可以画子图, 有两种调用方式,

```
subplot(m,n,p)  
subplot('Position',pos)
```

第一种方式是创建一个 $m \times n$ 的网格, 并在第 p 个网格上绘图; 第二种方式是在指定 `pos` 上绘图, 座标属性为 `[left bottom width height]`。可以通过查阅图像对象来获取图形属性。

值得注意的是 `subplot` 这个函数创建的是一个 `axes` 对象, 所有 `axes` 对象可修改的属性对 `subplot` 都起作用, 哦也。

1.14.4 最佳实践

用一个图像对象的关键是搞清楚这个对象属于图像系统中的哪一类主要对象, 然后去阅读该类对象的文档, 切记, 阅读帮助文档是最有效的解决方案。

新版本 MATLAB 推荐使用调用对象的方式修改图形对象属性, 老版本用 `get`, `set` 函数分别查阅和修改对象属性, 或者使用参数传递的方法修改图形对象的属性, 没有调用对象的方式优雅。下面的代码展示两种方式的区别。

```
p = plot(1:10);  
p.Color = 'r';  
set(p, 'Color', 'red');
```

^⑬ <https://www2.mathworks.cn/help/matlab/formatting-and-annotation.html>

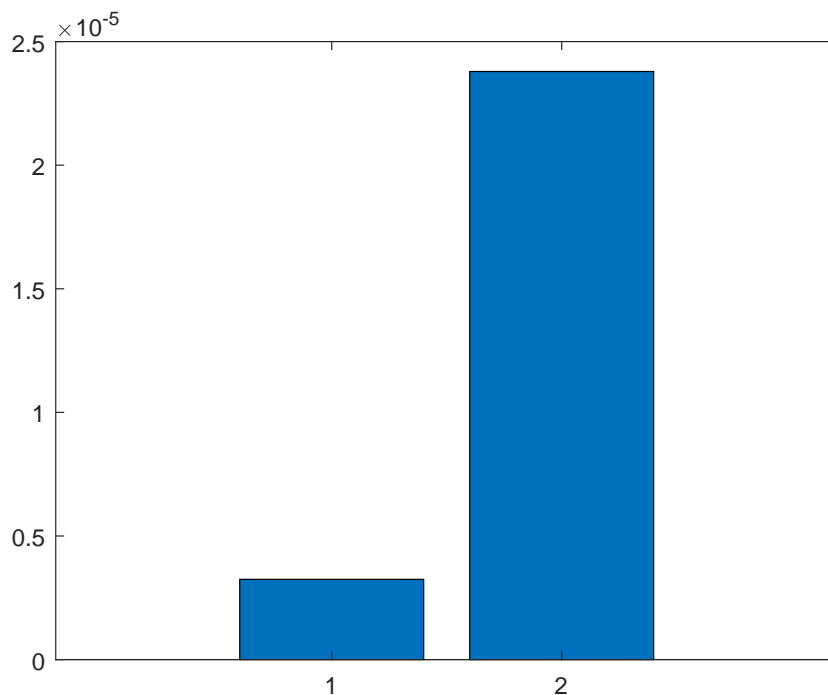


图 1.5: 默认样式的 bar 图

绘图最佳实践的代码风格如下所示，依次建立 figure, axes 和图表对象（如 Bar, Line 等），然后使用面向对象编程的方式修改其属性，有些图形属性有多层，此时切勿烦躁，及时使用 `get()` 函数查阅对象属性。坚持使用面向对象编程，代码风格会更加清晰，不宜混乱。实际上，传参和 `set` 这两种方式最大的劣势就是不能自动补全，而且写字符串很别扭。

```
fig = figure;
ax = axes;
fig_bar = bar(cate, [all_data{:, 2}]);
ax.YLabel.String = 'Variance';
```

1.14.5 谈一谈 colorbar

colorbar 这个对象属于插图对象，要修改的属性常常是位置，阅读官方文档吧^⑭！

怎样用代码修改 colorbar 的颜色和色值范围呢？这个问题困扰了我很久，一直觉得就是一行代码的事，果然不错，见官方文档^⑮！

`caxis(target, [min max])` 函数接受两个对象，目标对象可以是一个 axes 或者 figure 对象。

```
caxis(ax, [1, 100]) % 修改 ax 对象的色值范围为 1 - 100
```

1.14.6 如何给 bar 加图例等说明性对象

如 1.5 所示一个普通的 bar 图，如果我们需要加一些修饰，比如图例，坐标轴标题。。。该怎么办？

^⑭ <https://www.mathworks.cn/help/matlab/ref/matlab.graphics.illustration.colorbar-properties.html>

^⑮ <https://www.mathworks.cn/help/matlab/ref/caxis.html>

可以通过传递参数的方式修改坐标轴标题，形如 `bar(x, y)`。

```
c = categorical({'apples','pears','oranges'});  
prices = [1.23 0.99 2.3];  
bar(c,prices)
```

当然，还可以修改 `XTickLabel` 的属性值达到修改坐标轴标题的目的。

```
fig = figure;  
ax = axes;  
ax.XTickLabel = {'Small', 'Big'};
```

修改制定条目的颜色，使条形图使用 `CData` 属性中定义的颜色，然后更改矩阵中的对应行即可更改特定条目的颜色。这个属性可以玩出很多花样，毕竟修改条目颜色就是修改矩阵数值这么简单。

```
b = bar(rand(10,1));  
b.FaceColor = 'flat';  
b.CData(2,:) = [.5 0 .5];
```

1.14.7 如何导出大小和图形窗口大小一致的 pdf 文件？

有时候我们想把绘制的图形导出成一个 pdf 文件，但是发现 MATLAB 的默认设置使用的是 A4 纸，图形总是无法铺满的。我们可以将图形和图窗的大小设置为一样，然后调整几个 `paper` 参数为图窗的对应参数，再导出就可以得到与图窗大小一致的 pdf 或其他图片文件了。

```
fig.Units = 'points';  
fig.PaperUnits = 'points';  
fig.PaperPositionMode = 'Auto';  
fig.PaperPosition = [0, 0, fig.Position(3:4)];  
fig.PaperSize = fig.Position(3:4);  
print(fig, '-dpdf', 'test');
```

1.15 函数最佳实践

这里收集几个关于函数的最佳实践，首先是函数的返回值问题。很多时候，MATLAB 函数的返回值不止一个，通常让变量一一返回，在外部使用变量一一接收。最佳做法是将需要返回的变量装入 `struct` 中，这样在外部仅使用一个变量就能接收所有返回值。

1.16 关于路径地一切

处理数据肯定要读取文件，那么就免不了要处理路径，这里记录一下关于路径处理的最佳实践。

1. 拼接路径时最好使用内置函数 `fullfile`，而不是手动拼接字符串；
2. 使用 `isfile`, `isfolder` 对输入的文件进行异常检查；
3. 使用 `fileparts` 函数分割路径，

4. 如果需要手动插入路径分割符, 务必使用 `filesep` 函数, 其不同系统上有不同表现形式, 用户不需要为跨平台而进行特殊处理;

5. `dir` 函数是人类的好朋友;

```
file = 'H:\user4\matlab\myfile.txt';  
[filepath, name, ext] = fileparts(file);
```

1.17 如何在矩阵中插入一行/一列

看了一圈, 只能通过拼接矩阵来实现, 并且 MATLAB 自己没有这样的函数。自己写了一个。

```
% 按行插入  
function mat = row_insert(mat_row, pos, mat_add)  
[~, c_row] = size(mat_row);  
[~, c_add] = size(mat_add);  
if c_row == c_add  
    mat = [mat_row(1:pos-1, :); mat_add; mat_row(pos:end, :)];  
else  
    error(' 要插入的矩阵列要与原始矩阵一致')  
end  
end  
  
% 按列插入  
function mat = col_insert(mat_row, pos, mat_add)  
[r_row, ~] = size(mat_row);  
[r_add, ~] = size(mat_add);  
if r_row == r_add  
    mat = [mat_row(:, 1:pos-1), mat_add, mat_row(:, pos:end)];  
else  
    error(' 要插入的矩阵行要与原始矩阵一致')  
end  
end
```

1.18 图形处理

1.18.1 基本法

1.18.2 图形插值

图形插值就是让模糊的图形(像素较少)变成比较清楚的图形(像素变多), 稍微科学一点的说法是, 图像插值就是利用已知邻近像素点的灰度值(或 rgb 值)来产生未知像素点的灰度值(或 rgb 值), 以便由原始图像再生出具有更高分辨率的图像。图像处理的学问太多, 涉及到的数学也比较深, 这里只简单记录应用。

常见的插值算法有最近邻插值算法, 双线性插值算法, 三次卷积等。具体内容可以去学习一些 DIP(Digital image processing) 的公开课。如:

<http://inside.mines.edu/~whoff/courses/EENG510/lectures/>

<http://cs.nju.edu.cn/liyf/dip15/dip.htm>

常用的是 MATLAB 内置的 `interp2` 函数,

Github 上找了个现成的工具箱[®], 快速实现 bicubic 和 bilinear 两种插值算法, 效果很不错。

1.19 优雅的数据流处理

1.20 如何更新绘图?

有时候我们会跟踪一个变量的变化, 需要把变量不断地 `plot` 出来。如果是跟踪变量的“轨迹”, 即不抹去上一时刻的图形, 可以非常简单的使用 `pause` 函数实现; 如果需要实实在在更新变量的绘图, 可以使用 `refreshdata` 更新绘图数据。

```
...
l = plot(X, ObjV, 'b*');
l.XDataSource = 'X';
l.YDataSource = 'ObjV';

while something
    X.update();
    ObjV.update();
    pause(0.1);
    refreshdata
end
```

1.21 如何给不断更新的绘图制作 gif?

上一节讲了更新绘图, 所以肯定会有导出为动画的需求。顺手写了一个导出 gif 的函数, 传入一个图形对象, 控制参数和文件名即可在脚本当前所在目录生成一个 gif 文件。该函数比较简便, 主要利用到 `imwrite` 函数。

```
function ex_gif(fig, n, filename)
% Input - figure object
%       - control parameter
%       - file name
frame = getframe(fig);
im = frame2im(frame);
[ind, map] = rgb2ind(im, 256);
% Write to the GIF File
if n == 1
    imwrite(ind, map, filename, 'gif', 'Loopcount', inf);
else
    imwrite(ind, map, filename, 'gif', 'WriteMode', 'append');
```

[®] <https://github.com/FlorentBrunet/image-interpolation-matlab>

```
end  
end
```

1.22 更优雅的颜色表

仿 matplotlib^{①7}。

1.23 关于箱型图

统计学上对数据进行分析时，往往需要同时分析数据的诸多方面，如同时显示数据的极值、四分位矩等，箱型图能非常方便展示这些数据，还能揭示数据间的离散程度、异常值、分布差异等。MATLAB 中的箱型图函数是 `boxplot`，不过有一件事情很奇怪，该函数的返回值不是图表对象，而是绘图的数值。该箱型图的本质其实是很多 `line` 对象构成的普通线图，使用 `handle` 函数能很方便地查看 `boxplot` 函数返回的对象数组。所以想要修改默认的箱型图样式，只需要逐个修改这些 `line` 对象，可惜这方法太蠢了，后面再讨论更优雅的办法吧！

```
load carsmall  
h = boxplot(MPG, Origin);  
h = handle(h)
```

`boxplot` 有 7 个返回值，分别是：

1. Upper Whisker
2. lower Whisker
3. Upper Adjacent value
4. Lower Adjacent value
5. Box
6. Median
7. Outliers

如果要在同一张箱型图上绘制多个不相容的数据，需要进行分组。先把所有需要统计的数据整理为一行或者一列，同时整理一个相同大小的分组数组，用 1, 2, 3... 或者字符类型的数组进行分组标识。下面时一个例子。

```
A = [16, 20, 15, 17, 22, 19, 17]';  
B = [22, 15, 16, 16, 16, 18]';  
C = [23, 9, 15, 18, 13, 27, 17, 14, 16, 15, 21, 19, 17]';  
  
group = [ones(size(A)); ...  
2 * ones(size(B)); ...  
3 * ones(size(C))];  
  
figure();  
boxplot([A; B; C], group)
```

^{①7} <https://www.mathworks.com/matlabcentral/fileexchange/62729-matplotlib-2-0-colormaps-perceptually-uniform-and-beautiful>

第二章 Python

2.1 如何展开一个嵌套的序列？

我们现在有这样一个序列 `items = [1, 2, [3, 4, [5, 6, [9, 8], 7], 8]]`，我们想逐级展开这个序列，然后将所有元素装入一个序列。

如果这个序列层级较少，我们可以用多层 `for` 循环来遍历这个序列。一旦这个序列超过 3 层，过多的循环会让你很头疼。同样，这种多层级的问题我们可以用递归来解决。构建一个函数，这个函数能处理第一层的元素，由于第二层是 `list`，它是一个可迭代对象，我们只需要判断第二层是不是可迭代对象，同时忽略 `str`, `bytes` 对象^①。只要内层是可迭代的，我们就开始递归，对其应用该函数。

```
from collections import Iterable

#def unfold(items, unfolded=None, ignore_types=(str, bytes)):
#    unfolded = list() if unfolded is None else unfolded
#
#    for item in items:
#        if isinstance(item, Iterable) and not isinstance(item, ignore_types):
#            unfold(item, unfolded=unfolded)
#        else:
#            unfolded.append(item)
#
#    return unfolded

# 更优雅的版本
def unfold(items, ignore_types=(str, bytes)):
    unfolded = []
    for item in items:
        if isinstance(item, Iterable) and not isinstance(item, ignore_types):
            unfolded.extend(unfold(item))
        else:
            unfolded.append(item)

    return unfolded
```

^① `str`, `bytes` 也是可迭代对象，我们要避免其展开成单个字符。

由于存在递归，所以函数会被调用很多次，每次调用所得的数据都需要保留，如何在多次的调用之间共享保留数据呢？我采用一个默认参数来实现^②，首次调用时不给默认参数新值，这会产生一个空的 `list`，当对内层对象调用时，将上一次产生的数据赋值给这个参数。输出结果：

```
>>> items1 = ['Paula', ['Thomas', 'Lewis', ['siyu', 'ziyan', ['jianyuan']]]]
>>> items2 = [1, 2, [3, 4, [5, 6, [9, 8], 7], 8]]
>>> items3 = [[1, 2], 3, (4, [5, 6])]
>>> print(unfold(items1))
>>> print(unfold(items2))
>>> print(unfold(items3))
['Paula', 'Thomas', 'Lewis', 'siyu', 'ziyan', 'jianyuan']
[1, 2, 3, 4, 5, 6, 9, 8, 7, 8]
[1, 2, 3, 4, 5, 6]
```

但这样做有两个显而易见的坏处，一是当我们的嵌套序列有无限多层，递归会栈溢出；二是序列整个被读取到内存中了，当序列元素非常多，比如 1 亿，内存会被撑死。坏处一我们不去管他，大多数情况下是适用的，坏处二可以很容易的利用 `generator` 来解决^③。

```
from collections import Iterable

def unfold(items, ignore_types=(str, bytes)):

    for item in items:
        if isinstance(item, Iterable) and not isinstance(item, ignore_types):
            yield from unfold(item)
        else:
            yield item
```

使用 `generator` 一来能防止内存爆炸，二来不需要在函数的多次调用间传递数据，代码更清晰明朗。需要注意，`generator` 是惰性序列，边调用边计算，我们需要使用 `for` 迭代出每一个元素或者直接用 `list()` 获取全部元素。

```
items1 = ['Paula', ['Thomas', 'Lewis', ['siyu', 'ziyan', ['jianyuan']]]]
items2 = [1, 2, [3, 4, [5, 6, [9, 8], 7], 8]]
items3 = [[1, 2], 3, (4, [5, 6])]
print(list(unfold(items1)))
print(list(unfold(items2)))
print(list(unfold(items3)))
```

^② 前几天没有回想起 `list` 有个 `extend` 方法，显然用 `extend` 方法来实现更加优雅

^③ 思路来源 http://python3-cookbook.readthedocs.io/zh_CN/latest/c04/p14_flattening_nested_sequence.html

2.2 如何遍历当前文件夹及其子文件夹中的全部文件？

前面用 MATLAB 实现了一个，现在用 Python 来实现。第一种方法是利用递归来实现，思路同样是先找文件，然后找子文件夹，最后对子文件夹递归；第二种方法是利用 `os.walk` 模块，并将其做成 generator，这样在应对大量的文件时会有优势。推荐第二种方法，一来 `os` 模块考虑了很多我们忽略了的细节^④，二来 generator 是一个优雅的设计，用 Python 就应该好好学用 generator。

```
from os import listdir, walk
from os.path import isfile, isdir, join

def get_all_file_name(path):

    file_list = [join(path, f) for f in listdir(path) if isfile(join(path, f))]
    sub_dir = [join(path, d) for d in listdir(path) if isdir(join(path, d))]

    for i in range(len(sub_dir)):
        next_dir = sub_dir[i]
        file_list.extend(get_all_file_name(next_dir))

    return file_list

def get_all_file_name_generator(path):

    for folder, subdirs, files in walk(path):
        for f in files:
            yield join(folder, f)
```

2.3 如何在遍历 list 时删除元素？

存在一个 `list_a = [1, 2, 3, 4, 5, 6, 7, 8, 9]`，现在需要逐一操作内部元素，并在操作结束之后删除它。使用 `while` 判断 `list` 是否为空，不为空则 `pop` 第一个元素，在循环下依次操作每一个元素。

```
list_a = [1, 2, 3, 4, 5, 6, 7, 8, 9]

while list_a:
    temp = list_a.pop(0)
    print(temp)
```

^④ 比如，如果递归版本的函数遍历的根目录是一个磁盘，这个磁盘上的特殊的文件夹“System Volume Information”又是禁止被访问的，这时就会抛出一个 `PermissionError`。笨一点的解决办法是从子目录的 `list` 中删除这个目录，好一点的办法就是用 `os` 模块了。

2.4 字符串的处理

2.4.1 字符串自然排序

自然排序库, <https://github.com/SethMMorton/natsort>。

```
>>> a = ['2 ft 7 in', '1 ft 5 in', '10 ft 2 in', '2 ft 11 in', '7 ft 6 in']
>>> sorted(a)
['1 ft 5 in', '10 ft 2 in', '2 ft 11 in', '2 ft 7 in', '7 ft 6 in']

>>> from natsort import natsorted
>>> a = ['2 ft 7 in', '1 ft 5 in', '10 ft 2 in', '2 ft 11 in', '7 ft 6 in']
>>> natsorted(a)
['1 ft 5 in', '2 ft 7 in', '2 ft 11 in', '7 ft 6 in', '10 ft 2 in']
```

2.5 关于路径的一切

无论你用 python 干什么, 只要你操作文件, 就会涉及路径。接下来记录一些关于路径操作的心得, 避免以后踩坑。

脚本所在目录

一般情况, 如果你在当前目录启动 python, 你可以使用 `'.'` 或者 `os.getcwd()` 获取脚本所在目录 (也就是当前目录)。但是, 如果 python 不从当前目录启动, 此法就很有风险 (逃)。所有要坚定不移地使用如下这种优雅的方式获取脚本所在的绝对路径。

```
import os

current_path = os.path.split(os.path.abspath(__file__))[0]
```

获取指定文件类型的文件的路径

听起来很绕口, 就是说如果你只想获取一个目录里面的 pdf 文件, 该怎么办? 遇到问题一定要不假思索, 直接查阅 python 库 (人生苦短)。walk 函数可以很方便的遍历目录, splitext 可以把扩展名分开, 再加个判断, 并注意合并成绝对路径哦, 最后不要忘记使用优雅的生成器, 哦耶!

```
import os

def get_designated_file_name(path, ext):
    for folder, subdirs, files in os.walk(path):
        for file in files:
            if os.path.splitext(file)[1] == ext:
                yield os.path.join(folder, file)

current_path = os.path.split(os.path.abspath(__file__))[0]
file_name = list(get_designated_file_name(current_path, '.pdf'))
```

2.5.1 实例：遍历文件夹内的 PDF 文件并提取第一页合并

结合路径处理来个实例，我们要处理某个路径（本例是脚本运行目录下的）下的所有 pdf 文件，遍历所有文件并把第一页提取出来合并成一个单独的 pdf 文件。我们需要注意如下几点：遍历 pdf 文件时跳过上次输出的文件；避免意外，全部转化为绝对路径处理。

```
from PyPDF2 import PdfFileReader, PdfFileWriter, PdfFileMerger
import os

def get_designated_file_name(path, ext, out_file_name):
    for folder, subdirs, files in os.walk(path):
        for file in files:
            if os.path.splitext(file)[1] == ext:
                if file == out_file_name:
                    print(" 跳过 %s" % out_file_name)
                    continue
                else:
                    yield os.path.join(folder, file)

current_path = os.path.split(os.path.abspath(__file__))[0]
out_file_name = 'out.pdf'
out_file_path = os.path.join(current_path, out_file_name)
file_name = list(get_designated_file_name(current_path, '.pdf', out_file_name))

# ===== #

# pdf_out = PdfFileWriter()
# if len(file_name) != 0:
#     for f in file_name:
#         print(" 正在处理 %s" % f)
#         pdf_input = PdfFileReader(open(f, 'rb'))
#         pdf_out.addPage(pdf_input.getPage(1))
#     pdf_out.write(open(out_file_path, 'wb'))
# else:
#     print(" 当前文件夹没有 pdf 文件！")

# ===== #

pdf_out = PdfFileMerger()
if len(file_name) != 0:
    for f in file_name:
        print(" 正在处理 %s" % f)
        pdf_out.append(f, pages=(1, 2))
```

```
pdf_out.write(out_file_path)
else:
    print(" 当前文件夹没有 pdf 文件！")
```

2.5.2 pathlib

os 库就一个特点，用起来非常过程式，而且比较繁琐，完全没有爽快感。

2.5.3 实例：添加父文件夹名到文件名

首先提取父文件夹名到 `list`，然后遍历每个文件夹内的文件，对文件名的字符串进行处理后，使用 `os.rename` 重命名。

```
import os
import shutil

par_dirs = [name for name in os.listdir('.') if os.path.isdir(name)]

def rename(par_dir):
    print("Rename dir [{}] now!".format(par_dir))
    # get all mp4 file
    video_files = [
        video for video in os.listdir(par_dir) if video.endswith('.mp4')
    ]

    for video_file in video_files:
        newname = video_file.replace('_Downloadly.ir', '') # delete something
        if not video_file.startswith(par_dir):
            newname = par_dir + ' - ' + video_file # add par dir name before filename
        else:
            print("Pass add dir name")
        newname = os.path.join(par_dir, newname)
        oldname = os.path.join(par_dir, video_file)
        os.rename(oldname, newname)
        shutil.move(newname, '.')

for par_dir in par_dirs:
    rename(par_dir)
```


2.6 图像处理

一般都用 Pillow 库，面向对象，简单实用。

2.6.1 实例：给图片添加文字

我在给 Fluent 后处理的时候，由于数据量十分庞大，边计算边保存 case 以便后处理是不可能的，只能边计算边输出图片，然后后处理成动画。但是 Fluent 输出的图片没有当前时间，只能手动网上加了。核心功能由 ImageDraw 这个类提供。

```
from PIL import Image, ImageDraw, ImageFont
import os

im_width = 960
im_height = 720

# list all file and dir in current dir and extract png filename
dir_info = os.listdir('.')
png_files = [
    png for png in dir_info if os.path.isfile(png) and png.endswith('png')
]
print(len(png_files))

list_texts = [str(num) + ' s' for num in range(0, 1210, 10)]
print(len(list_texts))

def write_text_in_png(png_file, text):
    im = Image.open(png_file)
    draw = ImageDraw.Draw(im)
    font = ImageFont.truetype('arial.ttf', 36)
    draw.text((im_width / 2, im_height / 20), text, fill='black', font=font)
    im.save(png_file)

for num in range(0, len(png_files)):
    print("processing ", num)
    write_text_in_png(png_files[num], list_texts[num])
```

2.7 最佳实践

2.7.1 代码的模块化

要「模块化」就要分割代码。模块其实是一个逻辑层面的定义。它是说：如果一个东西，它有定义良好的输入和输出，那么它就是一个模块。例如一个定义良好的函数是一个模块；因为它接受固定格式的输

入（函数参数），根据固定逻辑产出输出。

撰写模块化代码遵循三个原则：

1. 单一职责：设计函数时，最好让函数的职责足够简单，只有一个；
2. 输入决定（input dominated）：设计函数时，尽可能使函数的行为完全由其输入参数决定。特别地，尽可能不要让函数的行为受到全局变量、类的成员变量的取值影响。某种意义上，这和「可重入」的概念比较像；
3. 自我限制：不写超过 50 行的代码；不超过 5 行的通用功能也要拉出去做工具函数。

2.8 matplotlib 最佳实践

2.8.1 生成矢量图

2.9 并行执行 shell 命令

最近接管了一台高性能 Linux 集群，root 权限在手，终于能为所欲为了。这台服务器有 1 个管理节点，13 个计算节点。管理节点的 `/opt/` 是共享的，计算节点通过访问这个目录获取计算程序。所以其他非共享目录的程序是各个节点独有的，比如 `top`。如此，就没办法在管理节点通过 `shell` 程序获取全部节点的 `cpu` 使用情况。我始终还不是专业人士，技术还没有学到家，一时也不知如何方便获取各个节点 `cpu` 占用率等信息。所以平时想看节点有没有人在计算，就只能一个个节点登陆，然后查看。

最近在网上搜了一圈，发现有个技术叫 Parallel SSH。人生苦短，马上加上 Python 这个关键字搜索相应的库。在 Github 上找到了一个最近一直在更新，star 和 fork 数量也合格的库，`parallel-ssh`。例子给得很直观，毫不费力就能仿制一个程序。接下来考虑如何方便地获取 `cpu` 占用率，随便搜索一下，发现 `psutil` 是这方面最出名的库。将最难的部分，并行执行 `shell` 命令和获取 `cpu` 占用率交给库后，很轻易就能写出同时获取所有节点 `cpu` 占用率的脚本。

如下，是并行执行 `shell` 的脚本 `parallel-ex-shell.py`。需要注意处理的是，代码中获取 `cpu` 占用率脚本 `get_cpu_info.py` 的路径最好使用绝对路径完全定位脚本位置。

```
# -*- coding: utf-8 -*-
from pssh.clients import ParallelSSHClient

hosts = [
    "cu01", "cu02", "cu03", "cu04", "cu05", "cu06", "cu07", "cu08", "cu09",
    "cu10", "cu11", "cu12", "cu13"
]

print(hosts)

client = ParallelSSHClient(hosts)
output = client.run_command("python /home/zousiyu1/bin/get_cpu_info.py")

for host, host_output in output.items():
    for line in host_output.stdout:
        print(host, ":", line)
    for line in host_output.stderr:
        print(line)
```

这里获取的是某一段时间内 cpu 占用率的平均值。

```
# -*- coding: utf-8 -*-
import psutil
import sys
import os

cpu_percents = []

for num in range(10):
    cpu_percents.append(psutil.cpu_percent(interval=0.5))

print(round(sum(cpu_percents) / len(cpu_percents)))
```

输出如下：

```
cu01 : 100
cu02 : 5
cu03 : 0
cu04 : 0
cu05 : 0
cu06 : 59
cu07 : 98
cu08 : 0
cu09 : 0
cu10 : 92
cu11 : 0
cu12 : 90
cu13 : 0
```

开心！同样，仿造这两个库的例子还能写出其他东西。

2.10 协助 ffmpeg 进行批量字幕压制

之前看过很多视频教程（英文的），但不想丢掉，于是想把字幕压进去收藏起来，同时也上传到哔哩哔哩分享给网友。简单 Google 一下发现 ffmpeg 能很方便压制字幕，就一条命令^⑤。

```
ffmpeg -i video.avi -vf "ass=subtitle.ass" out.avi
```

在视频文件很多的情况下，我们借助 Python 来辅助批量压制字幕。代码如下：

```
import subprocess
import os
```

^⑤ <https://trac.ffmpeg.org/wiki/HowToBurnSubtitlesIntoVideo>

```
def burn_subtitle(video, subtitle):
    ffmpeg = "/home/zousiyu1/bin/ffmpeg/ffmpeg"
    exec_list = [ffmpeg, "-i"]

    # parameters
    exec_list.append(video) # video file
    exec_list.append("-vf") # subtitle
    vf_args = "ass=" + subtitle
    exec_list.append(vf_args)

    # output
    exec_list.append(video[0:-4] + "-batch" + ".mp4")
    print("\n", exec_list, "\n")

    # run
    subprocess.call(exec_list)

# main

videos = []
subtitles = []

# get files
for _file in os.listdir("."):
    if _file.endswith("batch.mp4"):
        continue
    elif _file.endswith(".mp4"):
        videos.append(_file)
    elif _file.endswith(".ass"):
        subtitles.append(_file)

videos = sorted(videos)
subtitles = sorted(subtitles)

print(videos)
print(subtitles)

for num in range(len(videos)):
    burn_subtitle(videos[num], subtitles[num])
```

第三章 C 和 C++

3.1 C 语言的动态数组

大多数时候为了方便（其实是我菜），会使用库较多（方便）的 C++，但是 C 语言在实际生产中使用率仍然很高，比如长期使用的 ANSYS Fluent 的 UDF 就不得不用 C 语言。下面是一个简易的动态数组的实现，来源^①。

```
#include <iostream>

typedef struct
{
    int *array;
    size_t used;
    size_t size;
} Array;

void initArray(Array *a, size_t initialSize)
{
    a->array = (int *)malloc(initialSize * sizeof(int));
    a->used = 0;
    a->size = initialSize;
}

void insertArray(Array *a, int element)
{
    // a->used is the number of used entries, because a->array[a->used++] updates a->used
    // ↳ only *after* the array has been accessed.
    // Therefore a->used can go up to a->size
    if (a->used == a->size)
    {
        a->size *= 2;
        a->array = (int *)realloc(a->array, a->size * sizeof(int));
    }
    a->array[a->used++] = element;
}
```

^① <https://stackoverflow.com/questions/3536153/c-dynamically-growing-array>

```
void freeArray(Array *a)
{
    free(a->array);
    a->array = NULL;
    a->used = a->size = 0;
}

int main()
{
    Array a;
    int i;

    initArray(&a, 5); // initially 5 elements
    for (i = 0; i < 100; i++)
        insertArray(&a, i); // automatically resizes as necessary
    printf("%d\n", a.array[9]); // print 10th element
    printf("%d\n", a.used); // print number of elements
    freeArray(&a);
    return 0;
}
```

第四章 算法

算法可视化的网站，<https://www.cs.usfca.edu/%7Egalles/visualization/Algorithms.html>，晕乎乎的时候看几遍动画就明白了。

4.1 排序

先来了解几个基本概念：

排序：将一组“无序”的记录序列调整为“有序”的记录序列。

稳定性：假定在待排序的记录序列中，存在多个具有相同的关键字的记录，若经过排序，这些记录的相对次序保持不变，则称这种排序算法是稳定的，否则称为不稳定的。

排序算法的分类：插入类、交换类、选择类、归并类和基数类。

基于比较的排序算法的最佳性能为 $O(n \log n)$ 。

4.1.1 冒泡排序

复杂度： $O(n^2)$

1、对数组 `array[n]` 进行从 0 `n-1` 项的扫描，每碰到相邻两项数值大的在前小的在后时，对二者进行交换。当扫描进行完成后，0 `n-1` 中最大的元素必然已经在 `array[n-1]` 就位，而所有数值较小，序号却靠后的元素，序号也减小了 1。

2、既然最大的元素已在 `array[n-1]` 的位置就位，接下来的扫描只需从 0 `n-2`。具体过程同 1。同样的，扫描结束后 0 `n-2` 中最大的元素（全数组第二大的元素）必然已经在 `array[n-2]` 就位，而所有数值较小，序号却靠后的元素，序号也减小了 1。

3、如此不断重复，直到最小的元素在 `array[0]` 的位置就位。

从上述描述中我们可以看到“冒泡排序”这个名字的由来：每一次扫描，都可以使得数值较小，序号却靠后的元素的序号减少 1，宏观来看这些元素就像是从数组底部向上慢慢上浮的泡泡。

```
#include <iostream>
#include <algorithm>
using namespace std;

template <typename T>
void bubble_sort(T arr[], int size)
{
    int i, j;
    for (i = 0; i < size - 1; i++)
        for (j = 0; j < size - 1 - i; j++)
            if (arr[j] > arr[j + 1])
                swap(arr[j], arr[j + 1]);
}
```

```

}

int main()
{
    int arr[] = {61, 17, 29, 22, 34, 60, 72, 21, 50, 1, 62};
    int size = (int)sizeof(arr) / sizeof(*arr);
    bubble_sort(arr, size);
    for (int i = 0; i < size; i++)
        cout << arr[i] << ' ';
    cout << endl;

    float arrf[] = {17.5, 19.1, 0.6, 1.9, 10.5, 12.4, 3.8, 19.7, 1.5, 25.4, 28.6, 4.4,
        ↪ 23.8, 5.4};
    size = (int)sizeof(arrf) / sizeof(*arrf);
    bubble_sort(arrf, size);
    for (int i = 0; i < size; i++)
        cout << arrf[i] << ' ';
    return 0;
}

```

4.1.2 插入排序

复杂度： $O(n^2)$ ，具体算法描述如下：

1. 从第一个元素开始，该元素可以认为已经被排序
2. 取出下一个元素，在已经排序的元素序列中从后向前扫描
3. 如果该元素（已排序）大于新元素，将该元素移到下一位置
4. 重复步骤 3，直到找到已排序的元素小于或者等于新元素的位置
5. 将新元素插入到该位置后
6. 重复步骤 2 5

插入排序和人们打牌时所用的排序方式类似：抽第一张牌，此时手上的牌只有一张，所以是有序的。再抽一张牌，和手上的那张牌的大小进行比较，比它大就放在后面，否则放在前面。再抽一张牌，和手上的牌进行比较，插入在合适的位置，保持手上的牌有序。不断重复，直到牌抽完。从宏观来看，插入排序把数组分割成两部分，前段有序后段无序，随着插入排序的进行，后段无序的牌也越来越少，直到后段全部融入前段，排序也就结束了。

```

#include <iostream>
#include <algorithm>
using namespace std;

template <typename T>
void insert_sort(T arr[], int size)
{
    // 从第二个元素开始循环遍历未排序数组

```



```

    for (int i = 1; i < size; i++)
    {
        int temp = arr[i];
        int j = i - 1; //已排序数组最大索引
        // 从后向前扫描已排序数组
        while (j >= 0)
        {
            if (arr[j] > temp)
                arr[j + 1] = arr[j];
            else
                break;
            j--;
        }
        arr[j + 1] = temp;
    }
}

int main()
{
    int arr[] = {61, 17, 29, 22, 34, 60, 72, 21, 50, 1, 62};
    int size = (int)sizeof(arr) / sizeof(*arr);
    insert_sort(arr, size);
    for (int i = 0; i < size; i++)
        cout << arr[i] << ' ';
    cout << endl;

    float arrf[] = {17.5, 19.1, 0.6, 1.9, 10.5, 12.4, 3.8, 19.7, 1.5, 25.4, 28.6, 4.4,
        ↪ 23.8, 5.4};
    size = (int)sizeof(arrf) / sizeof(*arrf);
    insert_sort(arrf, size);
    for (int i = 0; i < size; i++)
        cout << arrf[i] << ' ';
    return 0;
}

```

4.1.3 归并排序

复杂度: $O(n \log n)$

归并排序的操作有两步, 分割和归并

- 1、分割: 将数组二等分, 并将得到的子数组继续二等分, 直到每个子数组只剩下一个元素为止。
- 2、归并: 不断将原本属于同一个数组的两个子数组归并成一个有序的数组, 方法为不断比较子数组的首元素, 并弹出较小的放入合并后组成的数组中。直到所有子数组合并为一个数组。

4.1.4 选择排序

复杂度: $O(n^2)$

在未排序序列中找到最小（大）元素，存放到排序序列的起始位置，然后，再从剩余未排序元素中继续寻找最小（大）元素，然后放到已排序序列的末尾。以此类推，直到所有元素均排序完毕。

```
#include <iostream>
#include <algorithm>
using namespace std;

template <typename T>
void select_sort(T arr[], int size)
{
    for (int i = 0; i < size - 1; i++)
    {
        // 假设最小元素为未排序部分的第一个
        int min = i;
        // 遍历未排序部分
        for (int j = i + 1; j < size; j++)
        {
            // 找到当前循环内的最小值
            if (arr[j] < arr[min])
            {
                min = j;
            }
        }
        swap(arr[i], arr[min]);
    }
}

int main()
{
    int arr[] = {61, 17, 29, 22, 34, 60, 72, 21, 50, 1, 62};
    int size = (int)sizeof(arr) / sizeof(*arr);
    select_sort(arr, size);
    for (int i = 0; i < size; i++)
        cout << arr[i] << ' ';
    cout << endl;

    float arrf[] = {17.5, 19.1, 0.6, 1.9, 10.5, 12.4, 3.8, 19.7, 1.5, 25.4, 28.6, 4.4,
        ↪ 23.8, 5.4};
    size = (int)sizeof(arrf) / sizeof(*arrf);
    select_sort(arrf, size);
    for (int i = 0; i < size; i++)
        cout << arrf[i] << ' ';
```

```
    return 0;  
}
```

4.1.5 快速排序

复杂度： $O(n^2)$ （最坏情况）； $O(n \log n)$ （平均情况），光听名字就很快。它其采用了一种分治的策略，通常称其为分治法 (Divide-and-ConquerMethod)。分治法的基本思想是：将原问题分解为若干个规模更小但结构与原问题相似的子问题。递归地解这些子问题，然后将这些子问题的解组合为原问题的解。

1. 从数列中取出一个数作为基准数（枢轴，pivot）；
2. 将数组进行划分 (partition)，将比基准数大的元素都移至枢轴右边，将小于等于基准数的元素都移至枢轴左边；

3. 递归地 (recursively) 把子数列排序。

C++ 实现的递归版本：

```
#include <iostream>  
#include <algorithm>  
using namespace std;  
  
template <typename T>  
int partition(T arr[], int low, int high)  
{  
    int pivot = arr[high]; // 挑最后一个元素做 pivot  
    int storeIndex = low - 1;  
  
    // 遍历数组  
    for (int j = low; j < high; j++)  
    {  
        if (arr[j] <= pivot)  
        {  
            storeIndex++; // 移动储存点  
            swap(arr[storeIndex], arr[j]);  
        }  
    }  
    // 将 pivot 放到两个子数组中间  
    swap(arr[storeIndex + 1], arr[high]);  
  
    return (storeIndex + 1); // 返回索引  
}  
  
template <typename T>  
void quick_sort(T arr[], int low, int high)  
{  
    // 数组尺寸大于 1
```

```

    if (low < high)
    {
        int p = partition(arr, low, high);
        quick_sort(arr, low, p - 1);
        quick_sort(arr, p + 1, high);
    }
}

int main()
{
    int arr[] = {61, 17, 29, 22, 34, 60, 72, 21, 50, 1, 62, 99, 90};
    int size = (int)sizeof(arr) / sizeof(*arr);
    quick_sort(arr, 0, size - 1);
    for (int i = 0; i < size; i++)
        cout << arr[i] << ' ';
    cout << endl;

    float arrf[] = {17.5, 19.1, 0.6, 1.9, 10.5, 12.4, 3.8, 19.7, 1.5, 25.4, 28.6, 4.4,
        ↪ 23.8, 5.4};
    size = (int)sizeof(arrf) / sizeof(*arrf);
    quick_sort(arrf, 0, size - 1);
    for (int i = 0; i < size; i++)
        cout << arrf[i] << ' ';
    return 0;
}

```

Python 实现的递归版本:

```

def quicksort(items, p, r):
    if p < r:
        q = partition(items, p, r)
        quicksort(items, p, q - 1)
        quicksort(items, q + 1, r)

def partition(items, p, r):
    x = items[r]
    i = p - 1
    for j in range(p, r):
        if items[j] <= x:
            i = i + 1
            items[i], items[j] = items[j], items[i]

```

```
    items[i + 1], items[r] = items[r], items[i + 1]
    return i + 1

items = [2, 5, 9, 3, 7, 0, -1]
quicksort(items, 0, len(items) - 1)
print(items)
```

问题：为什么临时储存点的索引是 low-1？

第五章 Git

5.1 如何给 Git 仓库添加一个空文件夹？

默认情况下，空文件夹不被记录，也不能被推送。特殊需求参见[How can I add an empty directory to a Git repository? - Stack Overflow](#)