

# 字典学习和稀疏表示

——数值计算方法大作业设计

组员 1：余鹏 学号：15352394

组员 2：袁子豪 学号：15352396

组员 3：詹巽霖 学号：15352399

组员 4：张海涛 学号：15352405

组员 5：张镓伟 学号：15352408

院系：数据科学与计算机

专业：软件工程(移动信息工程)

指导老师：纪庆革

组员合照：



从左到右依次为：张镓伟、余鹏、詹巽霖



张海涛

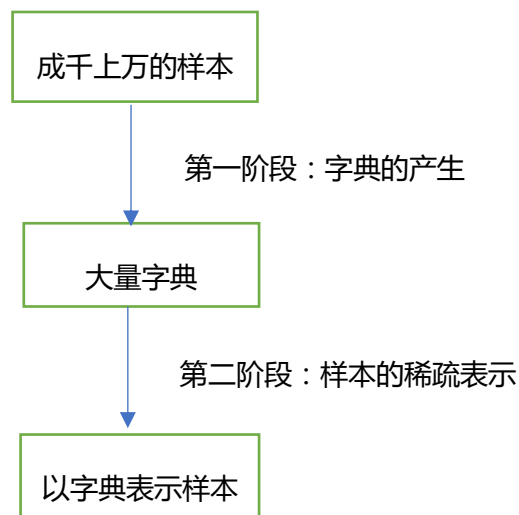


袁子豪

2017 年 5 月 28 日

## 一、引言

字典学习和稀疏表示一般被称为稀疏字典学习。该算法包括两个阶段：字典构建阶段和利用字典洗涤的表示样本阶段。如下图：



### 为什么要进行字典学习？

首先说一下字典的来源，我们知道人类社会的一切知识都必然要通过句子来表示出来，但无论有多少句子要被书写，对于每一个句子来说，它都有一个最本质上的特征，那就是构成这个句子字。可以这样说，无论人类的知识多么渊博，一本小小的字典就足以表达所有的知识，那些知识只不过是字典中的字的组合而已。所以我们看到字典学习的好处，它实质上是对于庞大数据集的一种降维表示。另外，正如同字是句子的最质朴的特征一样，字典学习总是尝试学习蕴含在样本背后最质朴的特征。所以我们要探讨字典学习。

### 为什么需要稀疏表示？

我想大家都有过这样的体会，在接触一个新的知识点时会有一种非常累的感觉，但是当我们把这些知识点完全掌握后，再次遇到相同的问题就会感觉到轻松，这是因为和初次接触知识点时相比，在完全掌握知识点后大脑可以调动尽可能少的脑区消耗尽可能少的能量进行同样有效的计算，此时大脑的计算速度也会变快，可以这样说，这是因为大脑学会了知识的

稀疏表示。因为稀疏表示的本质就是用尽可能少的资源表示尽可能多的知识，这种表示还能带来计算上的加快。

## 二、字典学习的方法

字典学习称之为稀疏编码。从矩阵分解角度看字典学习过程：给定样本数据集  $Y$ ， $Y$  的每一列表示一个样本；字典学习的目标是把  $Y$  矩阵分解成  $D$ 、 $X$  矩阵：

$$Y \approx D * X$$

同时满足约束条件： $X$  尽可能稀疏，同时  $D$  的每一列是一个归一化向量。

$D$  称之为字典， $D$  的每一列称之为原子； $X$  称之为编码矢量、特征、系数矩阵；字典学习可以有三种目标函数形式

(1)第一种形式：

$$D, X = \operatorname{argmin} \frac{1}{2} \|Y - D \cdot X\|^2 + \lambda \cdot \|X\|_0$$

这种形式因为  $L_0$  难以求解，所以很多时候用  $L_1$  正则项替代近似。

(2)第二种形式：

$$\begin{aligned} D, X &= \arg \min_{D, X} \{\|X\|_0\} \\ \text{st.} \quad &\|Y - DX\|^2 \leq \varepsilon \end{aligned}$$

$\varepsilon$  是重构误差所允许的最大值。

(3)第三种形式：

$$\begin{aligned} D, X &= \arg \min_{D, X} \|Y - DX\|^2 \\ \text{st.} \quad &\|X\|_0 \leq L \end{aligned}$$

$L$  是一个常数，稀疏度约束参数，上面三种形式相互等价。

## SRC 基本模型：

1、对  $x$  进行编码，对系数加  $l_1$  范数约束，取最小值。

$$\mathbf{a} = \underset{\mathbf{a}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{D}\mathbf{a}\|_2^2 + \lambda \|\mathbf{a}\|_1$$

2、对  $x$  进行分类判决。

$$c = \underset{i}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{X}_i \delta_i(\mathbf{a})\|_2^2$$

其中  $\delta_i(\mathbf{a})$  是提取与第  $i$  个类别相关的元素的向量指示函数。

## 传统字典学习框架：

$$\begin{aligned} \{\mathbf{A}, \mathbf{D}\} &= \underset{\substack{\mathbf{D} \in \mathbb{R}^{d \times K} \\ \mathbf{A} \in \mathbb{R}^{K \times N}}}{\operatorname{argmin}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{D}\mathbf{a}_i\|_2^2 + \lambda \|\mathbf{a}_i\|_1 \\ &= \underset{\substack{\mathbf{D} \in \mathbb{R}^{d \times K} \\ \mathbf{A} \in \mathbb{R}^{K \times N}}}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{D}\mathbf{A}\|_F^2 + \lambda \|\mathbf{A}\|_1 \\ &\text{s.t. } \|\mathbf{d}_i\|_2^2 \leq 1, \text{ for } \forall i = 1, \dots, N, \end{aligned}$$

其中，共  $N$  个信号。 $\mathbf{A}=[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N]$  是编码系数矩阵。矩阵  $\mathbf{A}$  的  $l_1$  范数等价于  $\mathbf{A}$  的各个列向量的  $l_1$  范数之和。

但 SRC 方法是有缺点的：

- 1、预定字典包含冗余和琐碎的不利于人脸识别的信息。
- 2、训练数据增加，稀疏编码的计算量增加。

为了解决以上的问题，可以使用以下两类已经存在的基于 DL 的方法：一类是直接学习具有识别力的字典，另一类是稀疏化稀疏，是字典具有识别力。

## 1、直接学习具有识别力的字典

### Meta-Face Learning 方法：

该方法是针对每一个类别学习得到一个自适应的字典。

$$\begin{aligned} \mathbf{D}_i = \underset{\mathbf{D}_i}{\operatorname{argmin}} \quad & \|\mathbf{X}_i - \mathbf{D}_i \mathbf{A}_i\|_F^2 + \lambda \|\mathbf{A}_i\|_1, \\ \text{s.t.} \quad & \|\mathbf{d}_j^i\|_2 \leq 1, \forall j = 1, \dots, K, \end{aligned}$$

其中矩阵  $\mathbf{X}_i$  包含第  $i$  个类别的所有样本， $\mathbf{D}_i$  是第  $i$  个类别对应的字典。

## Dictionary Learning with Structured Incoherence ( DLSI)方法：

不同类型的子字典具有连贯性，则重构查询图像时的原子是可以互相代替的。这导致无法用重构误差进行判决，为了解决这一问题，Ramirez 等人增加了一个不连贯项的约束，使不同类型的子字典之间尽可能互相独立。

不连贯项：

$$Q(\mathbf{D}_i, \mathbf{D}_j) = \|\mathbf{D}_i^T \mathbf{D}_j\|_F^2.$$

最终的字典学习算法为：

$$\min_{\{\mathbf{D}_i, \mathbf{A}_i\}_{i=1, \dots, C}} \sum_{i=1}^C \left\{ \|\mathbf{X}_i - \mathbf{D}_i \mathbf{A}_i\|_F^2 + \lambda \|\mathbf{A}_i\|_1 \right\} + \eta \sum_{i \neq j} \|\mathbf{D}_i^T \mathbf{D}_j\|_F^2$$

不连贯项的含义是：在重构误差是，忽略与公共原子相关的系数的绝对值，一次提高系统的判决能力。

## 2、使系数具有识别力的方法

该类方法使稀疏系数具有识别力，间接使字典具有识别力，该类方法只需要学习一个整体的字典，不需要每个类别都学习一个相应的字典。

## 监督字典学习 ( SDL ) 方法：

Mairal 等人提出讲逻辑回归与传统字典学习框架结合。优化公式为：

$$\begin{aligned}
(\mathbf{A}, \mathbf{D}) = \underset{\substack{\boldsymbol{\theta} \\ \mathbf{D} \in \mathbb{R}^{d \times K} \\ \mathbf{A} \in \mathbb{R}^{K \times N}}}{\operatorname{argmin}} \sum_{i=1}^N (\mathcal{C}(y_i f(\mathbf{x}_i, \mathbf{a}_i, \boldsymbol{\theta})) + \lambda_0 \|\mathbf{x}_i - \mathbf{D} \mathbf{a}_i\|_2^2 + \lambda_1 \|\mathbf{a}_i\|_1) + \lambda_2 \|\boldsymbol{\theta}\|_2^2 \\
\text{s.t. } \|\mathbf{d}_i\|_2^2 \leq 1, \text{ for } \forall i = 1, \dots, N,
\end{aligned}$$

where  $\mathcal{C}$  is the logistic loss function ( $\mathcal{C}(x) = \log(1 + e^{-x})$ ),

$\lambda_2$  是防止过拟合的正则项参数， $f$  是与系数  $\mathbf{a}$  呈线性关系的函数：

$$f(\mathbf{x}, \mathbf{a}, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{a} + b, \boldsymbol{\theta} \in \mathbb{R}^L$$

或者与  $\mathbf{a}$  成双线性关系的函数：

$$f(\mathbf{x}, \mathbf{a}, \boldsymbol{\theta}) = \mathbf{x}^T \mathbf{W} \mathbf{a} + b, \boldsymbol{\theta} = \{\mathbf{W} \in \mathbb{R}^{d \times L}, b \in \mathbb{R}\}$$

## 用于字典学习的具有识别力的 K-SVD 方法：

D-KSVD 在传统 DL 框架上加了一个简单的线性回归作为惩罚项。

$$\begin{aligned}
(\mathbf{D}, \mathbf{W}, \mathbf{A}) = \underset{\mathbf{D}, \mathbf{W}, \mathbf{A}}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{D} \mathbf{A}\|_F^2 + \lambda_1 \|\mathbf{Q} - \mathbf{G} \mathbf{A}\|_F^2 + \lambda_2 \|\mathbf{H} - \mathbf{W} \mathbf{A}\|_F^2 + \lambda_3 \|\mathbf{A}\|_1 \\
\text{s.t. } \|\mathbf{d}_i\|_2^2 \leq 1, \text{ for } \forall i = 1, \dots, N,
\end{aligned}$$

$\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N]$  是训练图像的标签。 $\mathbf{h}_n = [0, \dots, 0, 1, 0, \dots, 0]$ , 非零元素的位置即为所属

的类别。 $\mathbf{W}$  是分类器的参数。由上式可见，前两项可以混合为一项，最后一项根据 KSVD 可以去掉。最终优化得到  $\mathbf{D}$  和  $\mathbf{W}$  后，即可迅速对查询图像进行分类。

## 接下来详细介绍 K-SVD 方法

K-SVD 是字典学习的一种经典算法，其求解方法是固定其中一个，然后更新另外一个变量，交替迭代更新。字典  $\mathbf{D}$  的每一列就相当于一个聚类中心，稀疏编码  $\mathbf{X}$  的每一列允许有几个非零元素，给定训练数据  $\mathbf{Y}$ ,  $\mathbf{Y}$  的每一列表示一个样本，我们的目标是求解字典  $\mathbf{D}$  的每一列。

### 1. 随机初始化字典 $\mathbf{D}$ ：

从样本集  $\mathbf{Y}$  中随机挑选  $k$  个样本作为  $\mathbf{D}$  的原子，并且初始化编码 矩阵  $\mathbf{X}$  为 0 矩阵。

## 2.固定字典，求取每个样本的稀疏编码：

编码过程可以采用如下公式：

$$D, X = \arg \min_{D, X} \{\|X\|_0\}$$
$$st. \quad \|Y - DX\|^2 \leq \varepsilon$$

其中  $\varepsilon$  是重构误差所允许的最大值。

假设我们的单个样本是向量  $y$ ，为了简单起见我们就假设原子只有这 4 个，也就是字典  $D=[\alpha_1, \alpha_2, \alpha_3, \alpha_4]$ ，且  $D$  是已经知道的；我们的目标是计算  $y$  的编码  $x$ ，使得  $x$  尽量稀疏。

(1)首先从  $\alpha_1, \alpha_2, \alpha_3, \alpha_4$  中找出与向量  $y$  最近的那个向量，也就是分别计算点乘：

$$\alpha_1*y, \alpha_2*y, \alpha_3*y, \alpha_4*y$$

然后求取最大值对应的原子  $\alpha$ 。

(2)假设  $\alpha_2*y$  最大，那么我们就用  $\alpha_2$ ，作为我们的第一个原子，然后我们的初次编码向量就为：

$$x_1 = (0, b, 0, 0)$$

$b$  是一个未知参数。

$$(3) \text{求解系数 } b : y - b*\alpha_2 = 0$$

方程只有一个未知参数  $b$ ，是一个高度超静定方程，求解最小二乘问题。

$$(4) \text{然后我们用 } x_1 \text{ 与 } \alpha_2 \text{ 相乘重构出数据，然后计算残差向量：} y' = y - b*\alpha_2$$

如果残差向量  $y'$  满足重构误差阈值范围  $\varepsilon$ ，那么就结束，否则进入下一步。

(5)计算剩余的字典  $\alpha_1$ 、 $\alpha_3$ 、 $\alpha_4$  与残差向量  $y'$  的最近的向量，也就是计算

$$\alpha_1 * y' \text{ 、 } \alpha_3 * y' \text{ 、 } \alpha_4 * y'$$

然后求取最大值对应的向量  $\alpha$ ，假设  $\alpha_3 * y'$  为最大值，那么就令新的编码向量为：

$$x_2 = (0, b, c, 0)$$

b、c 是未知参数。

(6)求解系数 b、c,于是我们可以列出方程：

$$y - b * \alpha_2 - c * \alpha_3 = 0$$

方程中有两个未知参数 b、c，我们可以进行求解最小二乘方程，求得 b、c。

(7)更新残差向量  $y'$ ：

$$y' = y - b * \alpha_2 - c * \alpha_3$$

如果  $y'$  的模长满足阈值范围，那么就结束，否则就继续循环，就这样一直循环下去。

### 3、逐列更新字典，并更新对应的非零编码

通过上一步我们已经知道样本的编码，接着我们的目标是更新字典，同时还要更新编码，K-SVD 采用的方法是更新字典，就是更新第 k 列原子的时候其他的原子固定不变，假设我们当前更新第 k 个原子  $a_k$ ，矩阵 X 对应的第 k 行为  $x_k$ ，则目标函数是：

$$\|Y - DX\|^2 = \left\| Y - \sum_{j=1}^K \alpha_j \cdot x_j \right\|^2 = \left\| (Y - \sum_{j \neq k} \alpha_j \cdot x_j) - \alpha_k \cdot x_k \right\|^2 = \|E_k - \alpha_k \cdot x_k\|^2$$



我们需要注意的是  $x_k$  不是把  $X$  一整行都拿出来更新（因为  $X_k$  中既有非零元素又有非零元素，如果全部取出来的话那么计算的时候  $x_k$  就不再保持以前的稀疏性），所以我们只能抽取非零的向量。然后  $E_k$  只保留  $x_k$  对应的非零元素项。

对于上面的方程，我们可以通过求解最小二乘的方法来求解  $\alpha_k$ ，不过这样存在一个问题，我们求解的  $\alpha_k$  不是一个单位向量，因此我们需要采用 svd 分解，才能得到单位向量  $\alpha_k$ 。经过 svd 分解后，我们以最大奇异值所对应的正交单位向量，作为新的  $\alpha_k$ ，同时我们还需要更新系数编码  $x_k$  中的非零元素（根据 svd 分解）。然后算法就在 1 和 2 之间一直迭代更新，直到收敛。

### 三、字典学习和稀疏表示的演示和实现

下面使用基于 python3.5 实现的代码演示，测试平台为 Anaconda 的 spyder IDE。代码及意义如下：

第一步先导入各种工具包和测试样例，此测试样例为 scipy 库自带的图片

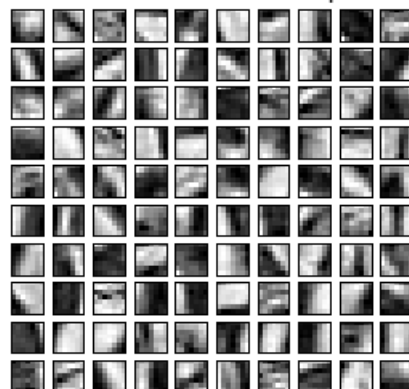
```
1 # -*- coding: utf-8 -*-
2 """
3 step1: 首先是各种工具包的导入和测试样例的导入
4 """
5 from time import time # 导入time模块, 用于测算一些步骤的时间消耗
6 import matplotlib.pyplot as plt # 导入画图模块
7 import numpy as np # 导入矩阵计算模块
8 import scipy as sp # 导入科学计算模块
9 from sklearn.decomposition import MiniBatchDictionaryLearning # 导入一种字典学习方法
10 from sklearn.feature_extraction.image import extract_patches_2d
11 # 导入碎片提取函数, 碎片提取函数将图片分割成一个个小块(patch)
12 from sklearn.feature_extraction.image import reconstruct_from_patches_2d
13 # 导入图片复原函数, 可通过patch复原一整张图片
14 from sklearn.utils.testing import SkipTest
15 # 导入测试工具nose下的异常抛出函数SkipTest
16 from sklearn.utils.fixes import sp_version
17 # 导入SciPy版本检测函数sp_version用于检测版本高低, 版本低于0.12的SciPy没有我们需要的样本测试
18 if sp_version < (0, 12):
19     raise SkipTest("Skipping because SciPy version earlier than 0.12.0 and "
20                  "thus does not include the scipy.misc.face() image.")
21 # 下面是测试样例导入, 我们使用的是scipy库自带的测试图片
22 try:
23     from scipy import misc
24     face = misc.face(gray=True)
25 except AttributeError:
26     # Old versions of scipy have face in the top level package
27     face = sp.face(gray=True)
```

## 第二步计算样例图片的字典 V

```
29 """
30 step2: 通过测试样例计算字典V
31 Convert from uint8 representation with values between 0 and 255 to
32 a floating point representation with values between 0 and 1.
33 """
34 face = face / 255.0
35 #读入的face大小在0~255之间，所以通过除以255将face的大小映射到0~1上去
36 face = face[:, :, 2] + face[:, :, 2] + face[:, :, 2] + face[:, :, 2] + face[:, :, 2]
37 face = face / 4.0
38 #以上两行对图形进行采样，把图片的长和宽各缩小一倍。
39 #记住array矩阵的访问方式 array[起始点: 终结点 (不包括): 步长]
40 height, width = face.shape
41 print('Distorting image...')
42 distorted = face.copy()
43 #将face的内容复制给distorted，这里不用等号因为等号在python中其实是地址的引用
44 distorted[:, width // 2:] += 0.075 * np.random.randn(height, width // 2)
45 #对照片的右半部分加上噪声，之所以左半部分不加是因为想要产生一个对比的效果
46 print('Extracting reference patches...')
47 t0 = time() #开始计时，并保存在t0中
48 patch_size = (7, 7) #tuple格式的patch大小
49 data = extract_patches_2d(distorted[:, :width // 2], patch_size) #对图片的左半部分（未加噪声的部分）提取patch
50 data = data.reshape(data.shape[0], -1)
51 #用reshape函数对data(94500, 7, 7)进行整形，reshape中如果某一位是-1，则这一维会根据
52 #（元素个数/已指明的维度）来计算这里经过整形后data变成（94500, 49）
53 data -= np.mean(data, axis=0)
54 data /= np.std(data, axis=0)
55 #以上两行作用在于每一行的data减去均值除以方差，这是zscore标准化的方法
56 print('done in %.2fs.' % (time() - t0))
57 print('Learning the dictionary...')
58 t0 = time() #开始计时，并保存在t0中
59 dico = MiniBatchDictionaryLearning(n_components=100, alpha=1, n_iter=500)
60 #上面这行初始化MiniBatchDictionaryLearning类，并按照初始参数初始化类的属性
61 V = dico.fit(data).components_
62 #调用fit方法对传入的样本集data进行字典提取，components_返回该类fit方法的运算结果，也就是我们想要的字典V
63 dt = time() - t0
64 print('done in %.2fs.' % dt)
65 plt.figure(figsize=(4.2, 4)) #figsize方法指明图片的大小，4.2英寸宽，4英寸高。其中一英寸的定义是80个像素点
66 for i, comp in enumerate(V[:100]): #: 循环画出100个字典V中的字
67     plt.subplot(10, 10, i + 1)
68     plt.imshow(comp.reshape(patch_size), cmap=plt.cm.gray_r,
69               interpolation='nearest')
70     plt.xticks(())
71     plt.yticks(())
72 plt.suptitle('Dictionary learned from face patches\n' +
73             'Train time %.1fs on %d patches' % (dt, len(data)),
74             fontsize=16)
75 plt.subplots_adjust(0.08, 0.02, 0.92, 0.85, 0.08, 0.23) #left, right, bottom, top, wspace, hspace
76
```

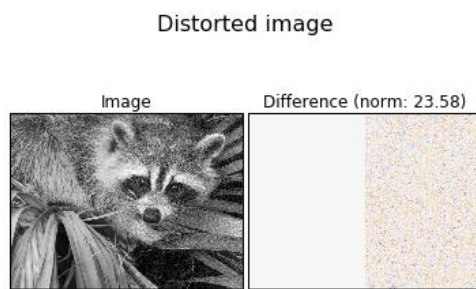
本步骤之后可以看到字典 V 如下：

Dictionary learned from face patches  
Train time 5.6s on 94500 patches



第三步我们画出标准图像和真正的噪声，方便同之后字典学习学到的噪声相比较

```
77 '''
78 step3: 画出标准图像和真正的噪声，方便同之后字典学习学到的噪声相比较
79 '''
80 def show_with_diff(image, reference, title):
81     """Helper function to display denoising"""
82     plt.figure(figsize=(5, 3.3))
83     plt.subplot(1, 2, 1)
84     plt.title('Image')
85     plt.imshow(image, vmin=0, vmax=1, cmap=plt.cm.gray,
86               interpolation='nearest')
87     plt.xticks(())
88     plt.yticks(())
89     plt.subplot(1, 2, 2)
90     difference = image - reference
91     plt.title('Difference (norm: %.2f)' % np.sqrt(np.sum(difference ** 2)))
92     plt.imshow(difference, vmin=-0.5, vmax=0.5, cmap=plt.cm.PuOr,
93               interpolation='nearest')
94     plt.xticks(())
95     plt.yticks(())
96     plt.suptitle(title, size=16)
97     plt.subplots_adjust(0.02, 0.02, 0.98, 0.79, 0.02, 0.2)
98 show_with_diff(distorted, face, 'Distorted image')
99
```

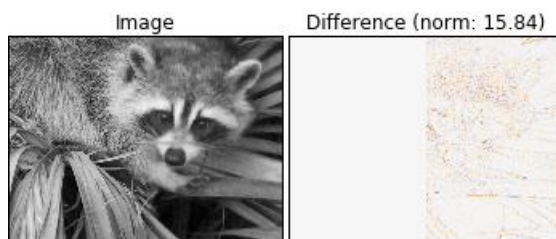


第四步我们测试不同字典学习方法对字典学习的影响

```
100 '''
101 step4: 测试不同的字典学习方法和参数对字典学习的影响
102 '''
103 print('Extracting noisy patches... ')
104 t0 = time()
105 data = extract_patches_2d(distorted[:, width // 2:], patch_size)#提取照片中被污染过的右半部进行字典学习。
106 data = data.reshape(data.shape[0], -1)
107 intercept = np.mean(data, axis=0)
108 data -= intercept
109 print('done in %.2fs.' % (time() - t0))
110 transform_algorithms = [ #这里是四中不同字典的表示策略
111     ('Orthogonal Matching Pursuit\n1 atom', 'omp',
112      {'transform_n_nonzero_coefs': 1}),
113     ('Orthogonal Matching Pursuit\n2 atoms', 'omp',
114      {'transform_n_nonzero_coefs': 2}),
115     ('Least-angle regression\n5 atoms', 'lars',
116      {'transform_n_nonzero_coefs': 5}),
117     ('Thresholding\n alpha=0.1', 'threshold', {'transform_alpha': .1})]
118 reconstructions = {}
119 for title, transform_algorithm, kwargs in transform_algorithms:
120     print(title + '...')
121     reconstructions[title] = face.copy()
122     t0 = time()
123     dico.set_params(transform_algorithm=transform_algorithm, **kwargs)#通过set_params对第二阶段的参数进行设置
124     code = dico.transform(data)
125     #transform根据set_params对设定参数的模型进行字典表示，表示结果放在code中。
126     #code总共有100列，每一列对应着V中的一个字典元素，所谓稀疏性就是code中每一行的大部分元素都是0。
127     #这样就可以用尽可能少的字典元素表示回去。
128     patches = np.dot(code, V)#code矩阵乘V得到复原后的矩阵patches
129     patches += intercept
130     patches = patches.reshape(len(data), *patch_size)#将patches从 (94500, 49) 变回 (94500, 7, 7)
131     if transform_algorithm == 'threshold':#
132         patches -= patches.min()
133         patches /= patches.max()
134     reconstructions[title][:, width // 2:] = reconstruct_from_patches_2d(
135         patches, (height, width // 2))#通过reconstruct_from_patches_2d函数将patches重新拼接回图片
136     dt = time() - t0
137     print('done in %.2fs.' % dt)
138     show_with_diff(reconstructions[title], face,
139                   title + ' (time: %.1fs)' % dt)
140 plt.show()
```

下面是 4 种不同的字典学习方法的结果：

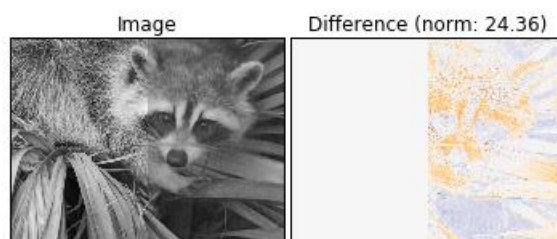
Orthogonal Matching Pursuit  
1 atom (time: 6.8s)



Orthogonal Matching Pursuit  
2 atoms (time: 14.9s)



Least-angle regression  
5 atoms (time: 70.9s)



Thresholding  
 $\alpha=0.1$  (time: 0.8s)



## 提交附件

test.py 该算法的实现代码