

文章编号:1001-9081(2006)12-2877-03

基于仿射变换和线性回归的 3D 人脸姿态估计方法

邱丽梅,胡步发

(福州大学 机械工程及自动化学院,福建 福州 350002)

(qiulimei1981@21cn.com)

摘要:提出了一种由仿射变换关系到线性回归的 3D 人脸空间姿态估计方法。即跟踪到人脸特征点后,根据仿射变换关系得到人脸姿态的粗估计值,以这个粗估计值作为人脸姿态的初始值,再通过线性回归迭代求得人脸姿态的精确值。实验结果表明,该方法在较大的姿态变化范围内,具有良好的估计精确度和鲁棒性。

关键词:人脸姿态估计;仿射变换;线性回归

中图分类号: TP391.41 **文献标识码:** A

3D face pose estimation based on affine transformation and linear regression

QIU Li-mei, HU Bu-fa

(College of Mechanical Engineering and Automation, Fuzhou University, Fuzhou Fujian 350002, China)

Abstract: A method for estimating 3D space face pose was proposed based on affine transformation and linear regression. Namely, after tracking face feature points, a roughly estimated value was obtained based on affine transformation relationship, and then this rough value was taken as the starting value, precise value of face pose was obtained based on linear regression iteration. Finally, the experimental results show that the method achieves better estimation accuracy and robustness in a wide range of face poses.

Key words: estimation of face pose; affine transformation; linear regression

0 引言

人脸姿态估计是确定输入图像序列中的人脸在 3D 空间中姿态的过程。姿态与视角有直接的对应关系,因此,姿态估计问题也称为人脸视角估计。现有人脸姿态估计的方法大体上可以分为两类:

1) 基于人脸外观的学习方法。假设 3D 人脸姿态与人脸图像的某些特性(图像密度、颜色、图像梯度值等)存在唯一的对应关系,用大量已知 3D 人脸姿态的训练样本,通过统计方法来建立这种关系^[1,2]。这种方法在样本数量不充分时,结果往往不够准确。

2) 基于模型的方法。利用某种几何模型或结构来表示人脸的结构和形状,建立模型和图像之间的对应关系,然后通过几何或者其他方法实现人脸空间姿态估计^[3-6]。应用深度数据获取设备,可以获取精度较高的人脸 3D 模型,但是设备价格昂贵,给研究带来很大的障碍。这种方法的优点在于:操作简单、精度和效率高。但由于对噪声敏感,使得估计鲁棒性比较差。

文献[5]提出利用仿射变换关系得到了人脸平面法线方向,但不能唯一确定人脸 3D 空间姿态。基于对现有人脸姿态估计方法的认识和文献[5]中存在的不足,本文提出了一种简单、低成本的人脸姿态估计新方法,既弥补了文献[5]中存在的不足,又通过迭代优化提高了估计精度。该算法对人脸图像去除噪声及平滑处理后,进行特征点提取和跟踪,然后估算仿射变换参数,根据圆—椭圆之间的仿射变换关系得到人脸姿态的粗估计值,再以这个粗估计值作为初始值,通过线性回归迭代求得人脸姿态估计的精确值。实验结果表明,在

较大的姿态变化范围内,该算法具有良好的估计精确度和鲁棒性。

1 基于仿射变换的人脸姿态粗估计

人脸姿态有 6 个自由度的变化,即沿 X 、 Y 、 Z 轴的平移和绕 X 、 Y 、 Z 轴的旋转。对沿 X 、 Y 的平移,在图像上表现为人脸的位置变化,可以通过统一坐标系实现;对沿 Z 轴的平移,在图像上表现为比例的变化,可以通过比例归一化实现。所以本文只研究人脸绕 X 、 Y 、 Z 三轴的旋转问题,旋转角分别为 α 、 β 和 γ 。

本文首先采用仿射变换粗估人脸旋转角。仿射变换算法是一种近似估计算法,图像在图像平面中的旋转可以转化为坐标轴的旋转。这种方法有两个基本假设:1) 当人脸离摄像机一定距离以上时,可以假设人脸是一个平面;2) 假设人脸在自然表情状态时,是一个刚性体。

1.1 特征点提取和跟踪

首先对输入图像进行维纳滤波,再自动提取正面平行参考人脸图像的两内外眼角和两嘴角^[7],因为这 6 个特征点的周围纹理信息较丰富,适合下一步的特征点跟踪。

在获取特征点的初始位置后,用 KLT 法对所需要的特征点进行跟踪。KLT 算法是一种以待跟踪窗口 w 在图像帧间的灰度差平方和(Sum of Squared intensity Differences, SSD)作为度量的跟踪算法^[8],目的是求出窗口中每个像素点的位移 d 。本文不是直接对每个特征点进行跟踪,而是以每一个特征点为中心,选取合适的矩形特征窗口 w ,通过计算出窗口内的像素点的水平位移 d_x 和垂直位移 d_y ,从而得到输入图像特征点的位置。

收稿日期:2006-06-23;修订日期:2006-09-04

作者简介:邱丽梅(1981-),女,福建三明人,硕士研究生,主要研究方向:图像处理、模式识别; 胡步发(1963-),男,福建宁德人,副教授,博士,主要研究方向:计算机视觉、模式识别

在特征窗口 w 内图像的 SSD 为:

$$e = \iint_w [J(X) - I(X - d)]^2 d_x \quad (1)$$

其中, J 和 I 分别是 $t + 1$ 和 t 时刻在矩形特征窗口 w 内像素的灰度值。但 KLT 算法有时会出现跟踪丢失和漂移,这在一定程度上是由于 d 太大造成的。对此本文根据对称性,提出用 $[J(x + \frac{d}{2}) - I(x - \frac{d}{2})]$ 灰度差算子代替 $[J(X) - I(X - d)]$, 这样既适应了人脸姿态较大旋转角度的变化,也确保了特征点跟踪的效果。所以式(1) 变成:

$$e = \iint_w [J(x + \frac{d}{2}) - I(x - \frac{d}{2})]^2 d_x \quad (2)$$

则:

$$\frac{\partial e}{\partial d} \approx \iint_w [J(X) - I(X) + \frac{1}{2}g^T(X)d]g(X)d_x \quad (3)$$

其中:

$$g(X) = \begin{bmatrix} \frac{\partial}{\partial x}(I + J) \\ \frac{\partial}{\partial y}(I + J) \end{bmatrix}$$

$$X = (x, y)^T \quad d = (d_x, d_y)^T$$

为了求 d , 令 $\frac{\partial e}{\partial d} = 0$, 则:

$$\iint_w [J(X) - I(X)]g(X)d_x = -\frac{1}{2}[\iint_w g(X)g^T(X)d_x]d \quad (4)$$

$$\text{有方程: } Z d = e \quad (5)$$

其中 Z 是一个 2×2 的矩阵, e 是 2×1 的矢量, 分别为:

$$Z = \iint_w g(X)g^T(X)d_x \quad (6)$$

$$e = 2 \iint_w [I(X) - J(X)]g(X)d_x$$

用牛顿-拉斐森(Newton-Raphson)方法求解方程(5), 解得 d , 该方法收敛速度快。

1.2 估算仿射参数

以 $(\mu_1, \mu_2, \dots, \mu_i, \dots, \mu_N)$ 表示正面平行参考人脸图像上的 N 个特征点, 用 $(\mu'_1, \mu'_2, \dots, \mu'_i, \dots, \mu'_N)$ 表示等待估计的人脸图像上相应的特征点。根据人体测量学, 一般人的两外眼角连线和两嘴角连线是相互平行的, 且把这两条平行线确定的平面称为人脸平面, 则任意的特征点对满足仿射变换有:

$$\mu'_i = A\mu_i + b, \quad i = 1, 2, 3, \dots, N \quad (7)$$

其中, A 是仿射参数的线性部分, b 是平移部分, A, b, μ_i 和 μ'_i 分别定义为:

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \quad b = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \quad \mu_i = \begin{pmatrix} \mu_{ix} \\ \mu_{iy} \end{pmatrix}$$

$$\mu'_i = \begin{pmatrix} \mu'_{ix} \\ \mu'_{iy} \end{pmatrix}, \quad i = 1, 2, 3, \dots, N$$

由这两组 N 个特征点对定义的仿射变换关系可表示为:

$$Kp = U \quad (8)$$

其中, $p = [a_{11} \ a_{12} \ b_1 \ a_{21} \ a_{22} \ b_2]$, 为仿射参数向量, K 为参考人脸图像上特征点的坐标值构成的矩阵, U 为待估计人脸图像上的特征点构成的向量。

$$K = \begin{bmatrix} \mu_{1x} & \mu_{1y} & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu_{1x} & \mu_{1y} & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mu_{Nx} & \mu_{Ny} & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu_{Nx} & \mu_{Ny} & 1 \end{bmatrix} \quad U = \begin{bmatrix} \mu'_{1x} \\ \mu'_{1y} \\ \vdots \\ \mu'_{Nx} \\ \mu'_{Ny} \end{bmatrix}$$

可以用最小二乘法求解这个最优化问题, 得到的仿射变换参数向量:

$$\hat{p} = (K^T K)^{-1} K^T U \quad (9)$$

1.3 粗略估算人脸 3D 姿态

利用圆-椭圆之间的仿射变换关系, 通过三次旋转得到粗略的人脸姿态估计。假设在正面平行参考人脸平面上存在一个假想圆, 关键是圆的方向, 而圆的大小对研究没有影响。当人脸姿态发生改变后, 它在图像平面上的投影就变成了一个椭圆, 所以确定人脸的空间姿态问题就转化为对摄像机坐标系的旋转使该圆的投影椭圆再次变为圆。基于文献[5]的仿射变换思想, 通过构造一个椭圆锥, 使它与图像平面的交线就是仿射变换之后的椭圆。通过对摄像机坐标系做两次旋转并得到旋转矩阵, 使该椭圆锥与图像平面的交线又变成圆, 得到人脸平面的法线方向, 再通过一次平面旋转, 得到唯一确定的人脸姿态。

第一次旋转摄像机坐标系, 使投影椭圆的主轴沿着 X 轴方向。根据文献[5]有矩阵 M :

$$M = \begin{bmatrix} A & \frac{B}{2} & \frac{D}{2f} \\ \frac{B}{2} & c & \frac{E}{2f} \\ \frac{D}{2f} & \frac{E}{2f} & \frac{F}{f^2} \end{bmatrix} \quad (10)$$

其中, $A = a_{11}^2 + a_{21}^2, B = 2(a_{11}a_{12} + a_{21}a_{22}), C = a_{12}^2 + a_{22}^2, D = 2(a_{11}b_1 + a_{22}b_2), E = 2(a_{12}b_1 + a_{22}b_2), F = b_1^2 + b_2^2 - 1, f$ 为摄像机焦距。矩阵 M 有特征值 $\lambda_1, \lambda_2, \lambda_3$ 且 $\lambda_1 > \lambda_2 > \lambda_3$, 对应的特征向量分别为 m_1, m_2, m_3 , 得第一次旋转矩阵为:

$$R_1 = [m_2 \ m_1 \ m_3] \quad (11)$$

假设绕 X, Y, Z 三轴分别旋转了 $\alpha_1, \beta_1, \gamma_1, R_1(\alpha_1), R_1(\beta_1), R_1(\gamma_1)$ 分别为相对应的旋转矩阵。其中:

$$R_1(\alpha_1) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos\alpha_1 & \sin\alpha_1 \\ 0 & -\sin\alpha_1 & \cos\alpha_1 \end{bmatrix}$$

$$R_1(\beta_1) = \begin{bmatrix} \cos\beta_1 & 0 & -\sin\beta_1 \\ 0 & 1 & 0 \\ \sin\beta_1 & 0 & \cos\beta_1 \end{bmatrix}$$

$$R_1(\gamma_1) = \begin{bmatrix} \cos\gamma_1 & -\sin\gamma_1 & 0 \\ -\sin\gamma_1 & \cos\gamma_1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R_1 = R_1(\alpha_1)R_1(\beta_1)R_1(\gamma_1) \quad (12)$$

根据式(11) 和(12), 得到 $\alpha_1, \beta_1, \gamma_1$ 。

第二次把 Y, Z 轴绕 X 轴旋转 θ , 使椭圆变成圆, 得:

$$\sin^2\theta = \frac{\lambda_2 - \lambda_1}{\lambda_3 - \lambda_1} \quad (13)$$

这样得到四个可能的姿态, 根据人脸的法线方向必须指向照相机, 可以排除两个姿态, 再通过图像实际的大致姿态, 排除第三个姿态, 得到 θ 。

两次旋转后, 得到人脸平面的法线方向, 但不能唯一确定人脸的 3D 空间姿态。因此需要第三次旋转, 即在人脸平面内绕 Z 轴的旋转, 利用两外眼角来估计旋转角 γ' 。

设两外眼角坐标分别为 $(\mu'_{1x}, \mu'_{1y}), (\mu'_{2x}, \mu'_{2y})$, 则:

$$\gamma' = \arctan \frac{\mu'_{1y} - \mu'_{2y}}{\mu'_{1x} - \mu'_{2x}} \quad (14)$$

综上所述, 由式(11) ~ (14) 可得到人脸姿态的粗略估计值 α, β, γ 。

2 基于线性回归的姿态估计值求精

仅基于仿射变换的人脸姿态估计值不够精确,因此用线性回归算法对粗略值进行迭代求精^[6]。假设一特征点的 3D 空间坐标为 $x = (x_1 \ x_2 \ x_3)^T$ 和相应的 2D 投影坐标为 $y = (y_1 \ y_2)^T$,则在不考虑平移时,它们之间的关系为:

$$y = sP\Phi x \tag{15}$$

s 为比例因子, P 为投影矩阵, $\Phi = R_x R_y R_z$ 为三维坐标系下的变形矩阵,其中, R_x, R_y, R_z 分别为绕 X, Y, Z 三轴旋转的旋转矩阵。假设人脸上有 n 个特征点 $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ 及相应的 2D 投影坐标为 $\{y^{(1)}, y^{(2)}, \dots, y^{(n)}\}$ 。根据式(15),则 3D 人脸姿态估计问题可用线性回归描述为:

$$\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} = s [I_n \otimes (P\Phi)] \begin{bmatrix} \check{x}^{(1)} \\ \check{x}^{(2)} \\ \vdots \\ \check{x}^{(n)} \end{bmatrix} \tag{16}$$

I_n 为单位矩阵, \otimes 代表矩阵张量积。把仿射变换得到的姿态估计粗略值作为线性回归的初始值。如已知 α, β 粗略值,则:

$$\check{x} = R_\alpha R_\beta x \tag{17}$$

因此由式(15)得到计算旋转角 γ 的线性回归模型如下:

$$Y = A\varphi \tag{18}$$

其中, $A = \check{x}_1 \otimes I_2 + \check{x}_2 \otimes \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$, $Y = y_1 \otimes \begin{pmatrix} 1 \\ 0 \end{pmatrix} + y_2 \otimes \begin{pmatrix} 0 \\ 1 \end{pmatrix}$, $\varphi = \begin{pmatrix} s \cdot \cos\gamma \\ s \cdot \sin\gamma \end{pmatrix}$ 。对于 s ,如前所述,可以通过归一化处理,则可得:

$$\hat{\varphi} = (A^T A)^{-1} A^T Y = (\|\check{x}_1\|^2 + \|\check{x}_2\|^2)^{-1} \begin{pmatrix} \check{x}_1^T y_1 - \check{x}_2^T y_2 \\ \check{x}_1^T y_2 + \check{x}_2^T y_1 \end{pmatrix} \tag{19}$$

γ 被确定后, α, β 以类似的方法来确定。这样 α, β, γ 就不断地被更新,直到误差小于事先调整好的阈值或循环次数达到了最大,循环结束。

3 实验

在检测范围为: $\alpha \in [-30^\circ, +30^\circ], \beta \in [-50^\circ, +50^\circ], \gamma \in [-90^\circ, +90^\circ]$ 内,获取 5 组人脸图像序列,对本文的算法进行了测试。

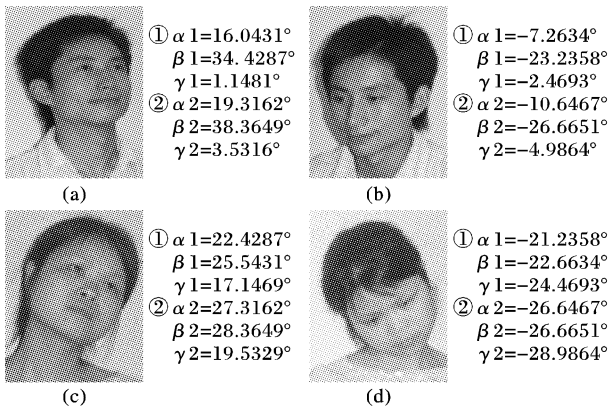


图 1 部分图像三个角度估计结果

由于人脸三个姿态角的真值很难同时测量,当人脸发生三姿态角变化时,测得粗略估计值和经过迭代 50 次后的估计值。其中两组人脸图像序列中的四幅图像估计结果如图 1

(a)~图 1(d)所示:①为粗略估计值;②为最终估计值。

实验表明,基于仿射变换的粗略估计值与经过线性回归迭代的估计值,两者差值的绝对平均值分别为: $\Delta\alpha = 3.4926^\circ, \Delta\beta = 2.9753^\circ, \Delta\gamma = 2.0128^\circ$ 。这些数据表明了迭代优化的必要性及本文算法的有效性。

在一般应用系统中对绕 Y 轴旋转的斜视图应用较多,所以通过测试 β 来检测算法的精确程度,最大迭代次数为 100 次。其中一组人脸图像序列中的三幅图像估计结果如图 2(e)~图 2(g)所示:①为粗略估计值;②为最终估计值;③为真实值。

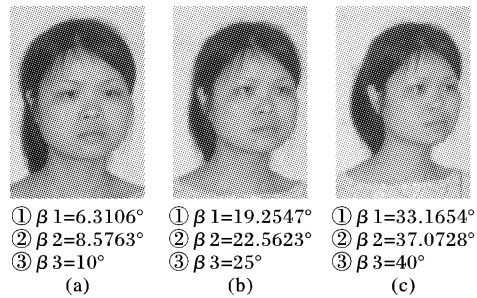


图 2 部分图像 β 角估计结果

图 3 是本组图像中当 $\beta = 10^\circ$ 时对应的迭代曲线。图 4 是本组图像在 $\beta \in [0^\circ, +50^\circ]$ 内,每间隔 5° 的误差对比曲线。

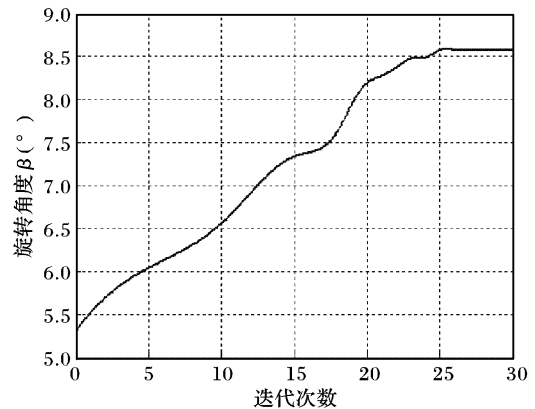


图 3 $\beta = 10^\circ$ 时对应的迭代曲线

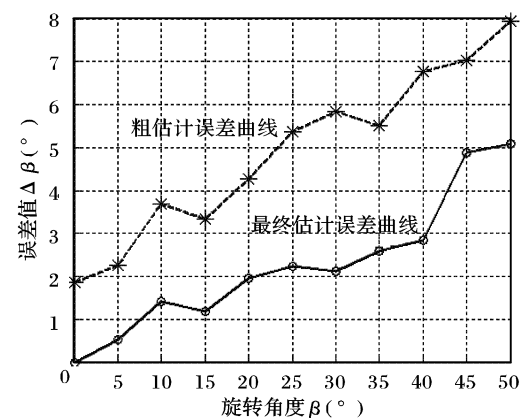


图 4 β 误差对比曲线

表 1 给出了用本文方法与文献[5,6]方法测得的 β 最终估计结果比较。

表 1 结果比较

方法	最终估计结果
文献[5]	只能得到人脸平面法线方向
文献[6]	β 角的估计精度为 6.54°
本文	β 角的估计精度为 2.7453°

从实验结果分析,本文算法与文献[6]有相似的视觉质量和相近的 PSNR 值,略高的 PESNR 值,但 CPU 处理插值的平均时间较小,介于 Bilinear 和 Bicubic 之间,小于文献[6]。因此,本文算法不但有较好的视觉质量,而且有较小的算法复杂度 and 相对更低的计算复杂性及低平均 CPU 处理时间。

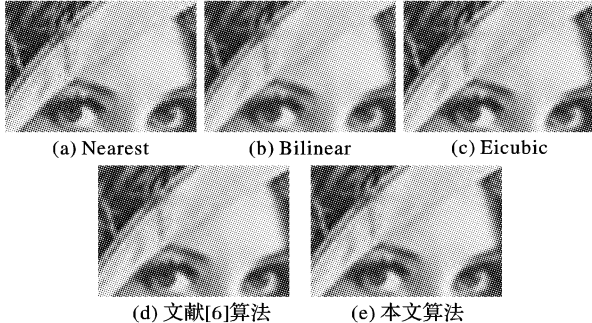


图 5 Lena 局部 5 种算法插值放大 2 倍视觉效果比较

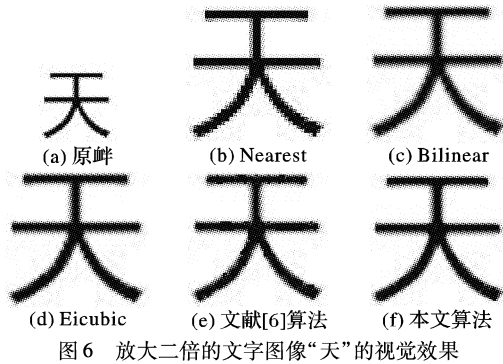


图 6 放大二倍的文字图像“天”的视觉效果

表 2 插值放大 2 倍的彩色图像 PSNR 各分量的性能比较

PSNR 标准	256 × 256 ~ 512 × 512 彩色图像					
	lena		pepper		football	
	文献[6]	本文	文献[6]	本文	文献[6]	本文
R	34.5832	34.629	34.0111	33.9488	34.3179	34.2877
G	34.2338	34.4307	34.4621	34.7271	34.1314	34.2877
B	33.3759	33.8034	34.1473	34.2266	33.7151	33.7803

(上接第 2879 页)

可见,本文方法不仅弥补了文献[5]只能得到人脸平面法线方向的不足,也具有有良好的估计精度。

4 结语

基于仿射变换和线性回归的 3D 人脸姿态估计方法是一种由粗到精的姿态估计方法。实验结果表明:利用该算法得到的 β 粗略估计值绝对平均误差约为 4.8976° , 最终估计值绝对平均误差约为 2.7453° , 且当人脸旋转角度在 40° 范围内时,估计误差都不超过 3° , 但超出 40° 时,误差有较明显增大;对于相同旋转角度下的不同人,姿态估计误差仅在大约 0.5° 的范围内波动。实验分析表明,该算法在较大的姿态变化范围内,具有良好的姿态估计精确度和鲁棒性。在该算法的基础上,下一步将通过 2D 人脸图像建立人脸 3D 模型,来估计人脸 3D 姿态,使算法适应“自遮挡”和多表情变化情况。

参考文献:

[1] LI SZ, FU QD, GU L, *et al.* Kernel Machine Based Learning for Multi-View Face Detection and Pose Estimation[A]. Proceedings of 8th IEEE International Conference on Computer Vision[C]. Vancouver, Canada, 2001, 674 – 679.
 [2] LI SZ, ZHU L, ZHANG Z, *et al.* Statistical Learning of Multi-View Face Detection[A]. Proceeding of The 7th European Conference on

表 3 各种插值算法平均 CPU 处理时间(s/图像)

彩色图像	Nearest	Bilinear	Bicubic	文献[6]	本文
Lena	0.010673	0.203708	0.24701	0.23585	0.226913
Pepper	0.010692	0.190914	0.252559	0.268894	0.235039

4 结语

基于边缘自适应最大相关性快速图像插值方法利用了图像边缘多个方向的最大相关性,在降低计算复杂度的同时,插值生成的图像边缘更清晰,消除了插值图像产生的锯齿,且具有自适应性和较好的视觉质量,保持了图像边缘和纹理特性。本文算法更大的优点是插值速度快,且算法复杂度小,计算量不大。

参考文献:

[1] 王耀南,李树涛,毛建旭. 计算机图像处理与识别技术[M]. 北京: 高等教育出版社, 2001.
 [2] KEYS RG. Cubic convolution interpolation for digital image processing[J]. IEEE Transactions on Acoustics, Speech, Signal Processing, 1981, 29(6): 1153 – 1160.
 [3] BATTIATO S, GALLO G, STANCO F. A locally-adaptive zooming algorithm for digital images[J]. Elsevier Image Vision and Computing Journal, 2002, 20(11): 805 – 812.
 [4] Li X, ORCHARD MT. New edge-directed interpolation[J]. IEEE Transactions on Image Processing, 2001, 10(10): 1521 – 1527.
 [5] 陈建辉,王博亮,徐中佑,等. 一种自适应最大相关性数字图像插值算法[J]. 厦门大学学报(自然科学版), 2005, 44(3): 355 – 358.
 [6] CHEN MJ, HUANG CH, LEE WL. A fast edge-oriented algorithm for image interpolation[J]. Image and Vision 2005, 23(9): 791 – 798.
 [7] 刘晓松,杨新,汪进. 基于统计特征的彩色图像快速插值方法[J]. 电子学报, 2004, 32(1): 29 – 33.
 [8] MUKHERJEE J, PARTHASARATHI R, GOYAL S. Markov random field processing for color demosaicing[J]. Pattern Recognition Letters, 2001, (22): 339 – 351.

Computer Vision (4)[C]. 2002. 67 – 81.

[3] CHOI KN, CARCASSONI M, HANCOCK ER. Recovering Facial Pose with the EM Algorithm[J]. Pattern Recognition, 2002, 35(10): 2073 – 2093.
 [4] YOSHINOBU EBISAWAL. Face Pose Estimation Based on 3D Detection of Pupils and Nostrils[A]. VECIMS 2005 IEEE International Conference on Virtual Environments, Human-Computer Interfaces, and Measurement Systems Giardini Naxos[C]. Itay, 2005.
 [5] YAO P, EVANS G, CALWAY A. Using Affine Correspondence to Estimate 3D Facial Pose[A]. Proceedings of the IEEE International Conference on Image Proceeding[C]. Thessaloniki, 2001, Vol 3: 919 – 922.
 [6] HU YX, CHEN LB, ZHOU Y, *et al.* Estimating Face Pose by Facial Asymmetry and Geometry[A]. Proceedings of the Sixth IEEE International Conference on Automatic Face and Gesture Recognition (FGR'04)[C]. IEEE, 2004.
 [7] DEND J, LAI F. Region-based Template Deformation and Masking for Eye-feature Extraction and Description[J]. Pattern Recognition, 1997, 30(30): 403 – 419.
 [8] SHI J, TOMASI C. Good Features to Track[A]. Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition[C]. Seattle, WA, 1994. 593 – 600.