

# 基于 SVM-2DPCA 的 X 光胸片异常筛查

王彦明<sup>1</sup>, 钱建忠<sup>2</sup>, 潘晨<sup>1</sup>

(1. 宁夏大学数学计算机学院, 银川 750021;

2. 银川市第二人民医院放射科, 银川 750001)

**摘要:** 基于统计学习理论的支持向量机分类算法, 提出一种 X 光胸片异常筛查系统, 能够自动判别胸片的正常和异常。为了提高 SVM 算法的效率, 利用小波变换等预处理手段去除对判读无用的图像冗余信息, 采用二维主成分分析进一步降低图像特征维数。实验结果表明, SVM 用于医学 X 光片异常筛查可行且有效、识别率高。

**关键词:** X 光片; 图像分类; 支持向量机; 二维主成分分析

## Abnormality Judgment of X-ray Chest File Based on SVM-2DPCA

WANG Yan-ming<sup>1</sup>, QIAN Jian-zhong<sup>2</sup>, PAN Chen<sup>1</sup>

(1. School of Mathematics and Computer Science, Ningxia University, Yinchuan 750021;

2. Radioactive Bureau, The Second People's Hospital of Ningxia, Yinchuan 750001)

**【Abstract】** Based on Support Vector Machine(SVM), the system for the abnormality judgment of X-ray chest file is presented, which can classify the X-ray picture normal and abnormal automatically. In order to improve the efficiency of the SVM, the wavelet transform is adopted in the system to eliminate the redundancy information in image. Two-Dimensional Principal Component Analysis(2DPCA) is used for feature extraction. Experimental results show that the SVM-based method is feasible in X-ray abnormality judgment, and has good classification ability.

**【Key words】** X-ray file; image classification; Support Vector Machine(SVM); Two-Dimensional Principal Component Analysis(2DPCA)

### 1 概述

医学是实践性和经验性很强的学科, 对于来源复杂、情况各异的病理数据, 每个医生的判读可能不尽相同。比如在边远地区医院, 或者某些大批量、重复性的疾病普查中, X 光片数据的判读非常容易受到诸如医生专业水平、人员疲劳等多方面因素的影响, 准确性和客观性都存在问题。在这样的场合, 对于拍摄到的图片, 如果能够发展一种自动化辅助手段识别出正常和异常情况, 对异常情况提示医生仔细观察和慎重诊断, 将大大提高疾病检查的效率。

图像识别是典型的模式识别问题。然而由于图像具有信息冗余, 图像场景的复杂性导致数据分布的任意性、样本数量有限、大规模数据和高维特征等具体问题, 利用传统的模式识别方法处理图像常常导致高昂的计算和存储开销。基于这种情况, 本文提出一种高效、准确的 X 光片自动识别方法。

### 2 系统简介

以支持向量机(Support Vector Machine, SVM)作为关键技术, 开发 X 光胸片计算机辅助识别系统。X 光胸片通常大致指人的锁骨以下、肚脐以上躯干区域的 X 光成像正面或侧面图片, 通常是灰度图。拍摄胸片时通常定位拍摄, 虽然图片大小有所不同, 但不同人的胸部位置大致相同。根据上述特点, X 光片的识别同灰度人脸图像的识别<sup>[1]</sup>具有许多相似之处。

本研究借鉴应用了人脸识别的一些相关技术, 如图像标准化方法、主成分分析 PCA 提取图像特征的思路等。图 1 为识别系统框架。

由图 1 知, 该系统由训练和识别 2 个阶段组成。其中, 训练阶段算法是本系统的核心, 主要利用有限数量的 X 光片

样本作为训练数据。首先通过图像预处理模块使得不同大小的图像规范为统一的尺寸(文中为 150×150 像素), 利用图像亮度变换和小波变换等手段去除图像冗余细节, 减少数据干扰; 随后使用二维主成分分析(Two-Dimensional Principal Component Analysis, 2DPCA)模块进一步降低图像维数, 用最有效特征描述图像; 最后通过 SVM 训练得到最优的分类模型, 并在训练过程中用“留一法交叉验证”得到最佳分类模型, 用该分类模型判读图像。

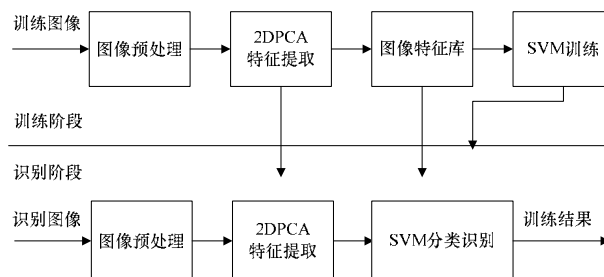


图 1 X 光胸片异常筛查系统框架

### 3 关键算法的基本原理及算法集成的思路

#### 3.1 SVM

SVM 是统计学习理论的一种具体实现。SVM 分类器应用了结构风险最小化原则, 它不仅要求最优分类面无错误地

**基金项目:** 国家自然科学基金资助项目(60663003); 宁夏自然科学基金资助项目(NZ0610)

**作者简介:** 王彦明(1984-), 男, 硕士研究生, 主研方向: 图形图像处理, 多媒体技术; 钱建忠, 硕士; 潘晨, 教授、博士

**收稿日期:** 2009-06-10 **E-mail:** wym-hero@126.com

分开各类(达到经验风险最小),而且要使类间间隔最大(达到分类器复杂度最小),从而保证使真实风险最小<sup>[2]</sup>。SVM的一般结构如图2所示。

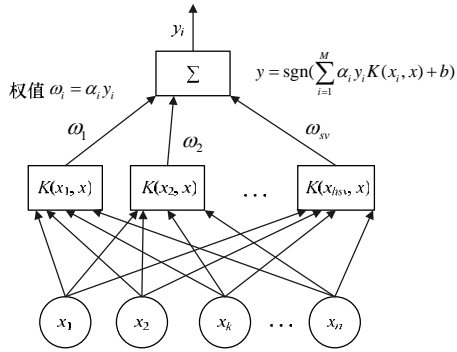


图2 SVM的一般结构

尽管通过非线性函数将样本数据映射到具有高维甚至无穷维的特征空间,并在特征空间中构造最优分类超平面,但在求解最优化问题和计算判别函数时并不需要显式计算该非线性函数,而只需计算核函数,从而避免特征空间维数灾难问题。

常见的核函数有

(1)线性函数:

$$K(x, x_i) = x \cdot x_i \quad (1)$$

(2)多项式核函数:

$$K(x, x_i) = [(x \cdot x_i) + 1]^d, \quad d=1, 2, \dots \quad (2)$$

(3)Sigmoid核函数:

$$K(x, x_i) = \tanh[v(x \cdot x_i) + c] \quad (3)$$

(4)高斯径向基核函数:

$$K(x, x_i) = \exp\{-q\|x - x_i\|^2\} \quad (4)$$

图像是典型的高维数据。尽管SVM的算法复杂度对高维数据不敏感,但是高维的数据运算和存储仍然需要占用大量的计算资源和空间。为了提高整体算法的效率,借鉴人脸识别技术的成功经验,提出利用二维主成分分析(2DPCA)手段进一步提取图像有效特征的方法。

### 3.2 2DPCA

主成分分析是一种经典的统计方法<sup>[3]</sup>,它对多元统计观测数据的协方差结构进行分析,以期求出能简约地表达这些数据依赖关系的主分量。设  $x_i \in R^{p \times q}$  是  $m$  个  $p \times q$  维的观测样本,  $X = (x_1, x_2, \dots, x_m)^T$  为观测样本矩阵,其总体样本的均值为

$$\mu = \frac{1}{m} \sum_{i=1}^m x_i \quad (i=1, 2, \dots, m) \quad (5)$$

用  $(x_i - \mu)$  来表示样本均值,这里的讨论是以样本的零均值情况进行的。样本集的协方差矩阵定义为

$$C = \sum_{i=1}^m [(x_i - \mu)(x_i - \mu)^T] \cong \frac{1}{m} (x_i - \mu)(x_i - \mu)^T \quad (6)$$

然而PCA提取图像特征时,大都是将二维图像按像素扫描顺序排列为  $p \times q$  维列向量,然后用一组基向量表示它们。这样使得二维图像变成了高维向量,丢失了图像结构信息,而且图像尺寸越大,该向量维数越高,容易造成维数灾难。

为了解决此类问题,文献[4]提出2DPCA。2DPCA的主要改进在于将某一个图像待测试样本  $x_{test}$  看作二维矩阵,这样不仅图像结构信息的完整性得到保留,而且可以有效地降低压缩特征。

由于2DPCA将  $x_{test}$  看作二维矩阵而非一维向量,因此观测样本的协方差矩阵可以有2种表示方式:

$$C_{row} = \frac{1}{m} (x_i - \mu)(x_i - \mu)^T \quad (7)$$

$$C_{col} = \frac{1}{m} (x_i - \mu)^T (x_i - \mu) \quad (8)$$

分别对式(7)、式(8)表示的协方差矩阵求得前  $d$  个最大特征向量组成的基矩阵  $feature_{row}, feature_{col}$ ,从而对于任何一个人X光片测试样本  $x_{test}$ ,可以按照下式来提取特征:

$$Y_{test} = (feature_{row}) x_{test} feature_{col} \quad (9)$$

对所有的样本按照式(9)进行投影变换,这样得到的图像特征就是  $d \times d$  维的。由于  $d$  值通常远小于图像的宽和高,因此利用2DPCA对二维图像数据进行特征抽取,不仅保留了一定图像结构信息,而且能为后续分类任务带来便利。

然而,2DPCA也存在一定局限性:(1)同PCA算法等一样,都是一种从信号整体特征考虑的特征提取算法。这样的算法对影响信号整体特征的干扰或变化敏感,如光照变化。(2)2DPCA在特征选择上是从特征最大的角度出发进行的,而不是从特征最优的角度出发,因而提取的特征对识别率的贡献也不一定是最优的,即并非所有提取的特征对识别都有利,过多的不相关特征有可能会降低识别性能。人和物理仪器对于图像都有一定的分辨率,对低于一定尺度的信号细节是无法认识的,因此,对所有尺度信号进行研究是没有必要的。为了解决这个问题,引入小波变换等预处理手段尽量减少对分类无用的噪声信息的干扰。

### 3.3 小波变换

小波变换是一种有效的多分辨率方法,通过对图像的高低通滤波可以将图像分解为对应不同尺度的近似分量(低频分量)和细节分量(高频分量),这样的方法有助于针对有意义的某种尺度信号展开研究。如文献[5]研究了人脸外观变化与小波频谱变化之间的关系,指出人脸的光照、少许遮罩、旋转扭曲和面部表情只影响图像中的高频部分,人脸图像的低频部分仍然保持稳定。

在计算机中可利用金字塔算法实现小波变换。在离散小波变换中,可将图像信号分解为许多低分辨率的成分,其中基本包含了原有图像的所有信息。同时在这些分辨率下,随机噪声和冗余信息被大大抑制。本文的具体做法是通过对图像样本预先进行小波分解,提取其低频分量子图作为原始图像的近似。图3是一个小波分解示意图。

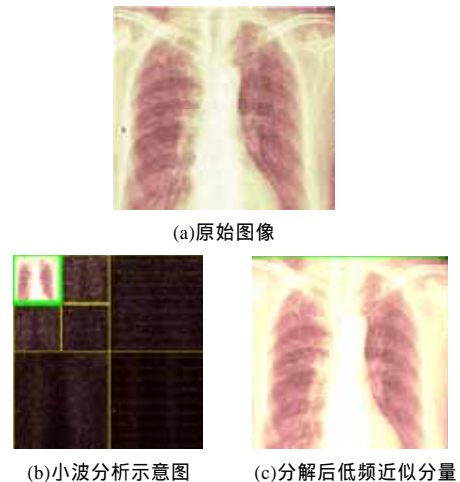


图3 小波分解示意图

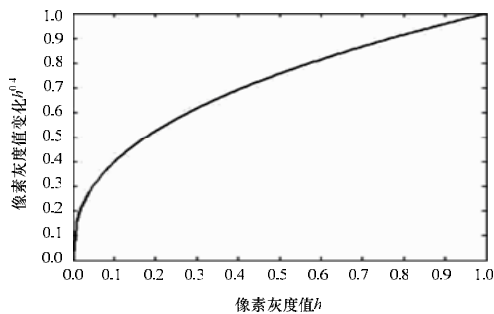
### 3.4 光照的改善

图像亮度变化同人的视觉感受之间具有指数特性。也就

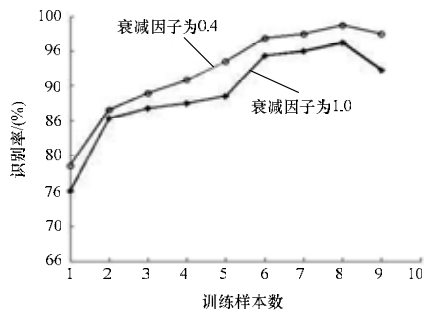
是说,当整体光照强度合适时,人对图像灰度变化较敏感,而光照变得太暗或太强,人对图像灰度变化也变得不敏感。PCA 是一种提取全局特征的方法,此类算法无论光照强弱,都对光照变化敏感。为了抑制光照影响,对样本图像做了简单的幂次变换处理,也就是对训练图像的像素亮度作一种指数衰减,使得机器视觉模拟人的感受:

$$h = h^a \quad 0 < a < 1 \quad (10)$$

这样的处理类似于图像显示时的伽马校正。实际效果是通过适当增加暗区的对比度,稍微提高整体图像亮度,抑制亮度较高时光照变化的幅度,从而提高识别精度。图 4 显示了采用上述策略能够提高识别正确率。在实验中,当指数取  $a=0.4$  时,识别的效果最佳。在后续的实验中,将这种方式作为预处理手段处理 X 光片图像。



(a)指数衰减因子



(b)衰减因子对识别率的影响

图 4 指数衰减因子对识别率的影响

## 4 实验数据分析

### 4.1 实验数据数量分布

系统由正常胸片和异常胸片各 25 个作为训练样本,是典型的小样本训练集。测试数据有 100 张正侧 X 光胸片,样本分布如表 1 所示。

表 1 实验测试样本分布

数据名称	数量/张
正常训练样本	25
异常训练样本	25
测试样本总数	100
测试样本正常图片	30
测试样本异常图片	60
不确定的图片	10

### 4.2 实验结果

在软件测试中,对医生已经明确标定类别的 90 张图片正确识别 84 张,正确率达到 93.25%。对于 10 张医生不确定类别的图片,软件测试结果为正常 1 张,异常 9 张;其中,异常的结果占多数,这样的结果比较符合胸片异常筛查的目的。

对于分错的测试样本,有可能是对训练样本分类较粗糙造成的。例如,大夫认为部分胸片判别的标准和年龄有关,不同年龄段衡量标准不同。对于这样的情况,通常将其归类为异常样本。实验测试样本结果如表 2 所示。

表 2 实验测试样本结果

数据名称	识别率/(%)
正常样本测试	93.25
异常样本测试	96.14

## 5 结束语

本文将不同算法有机地集成为一体,相互补充,实现机器视觉同人类视觉相匹配,能够得到好的训练和测试结果。从实验结果看,异常图片筛选方法精度较高、计算简单,能够帮助医生提高判读的准确率,具有很好的实用价值。

目前系统只实现了正常和异常 2 类图像的区分,功能上还存在一定局限性。进一步的改进应该引入更多的病人相关信息,如男、女、大人、小孩等;能够区分更多类别的肺部疾病,如肺癌、肺结核、肺炎等。相应的,通过对训练样本集按照上述要求分层分级进行处理,能够不断增加筛查系统的功能。虽然进一步改进还需要进行更多更细致的图像样本采集和多分类系统的设计工作,但是本文的研究明确地证明了基于 SVM 算法的计算机辅助 X 光片自动识别系统的有效性和可行性。利用这样的筛查模型作为一个插件,将来能够同商品化的 X 光机接口匹配,无疑为创造新的医学智能仪器打下了良好的基础。

## 参考文献

- [1] Belhumeur P N, Hespanha J P, Kriegman D J. Eigenfaces vs Fisherfaces: Recognition Using Class Specific Linear Projection[J]. IEEE Trans. on Pattern Analysis and Machine Intelligence, 1997, 19(7): 711-720.
- [2] Vapnik V. Statistical Learning Theory[M]. New York, USA: John Wiley & Sons, 1998.
- [3] Duda R O, Hart P E, Stork D G. Pattern Classification[M]. 2nd ed. [S. l.]: John Wiley & Sons, Inc., 2003.
- [4] Yang Jian, Zhang D, Frangi A F, et al. Two-dimensional PCA: A New Approach to Appearance-based Face Representation and Recognition[J]. IEEE Trans. on Pattern Analysis and Machine Intelligence, 2004, 26(1): 131-137.
- [5] Nastar C, Ayache N. Frequency-based Non-rigid Motion Analysis[J]. IEEE Trans. on Pattern Analysis and Machine Intelligence, 1996, 18(11): 1067-1079.

编辑 顾逸斐