

Transfer - Defect - Learning

矩阵 X X^T X 的转置 $\text{tr}(X)$ X 的迹
 X^{-1} X 的逆矩阵 矩阵 $X \in \mathbb{R}^{n \times m}$ 具有 n 个(列) 向量
 x_i 表示第 i 个实例 $x^{(i)}$ 表示第 i 个特征 m 列(特征值)

令 源项目数据集为 $D_S = \{(x_{Si}, y_{Si})\}_{i=1}^{n_1}$
 且 $x_{Si} \in \mathbb{R}^{1 \times m}$ 表示输入(例如: 一个文件被表示为向量)
 $y_{Si} \in$ 表示缺陷信息(例如: clean 或 buggy) 度量的

令 目标项目数据集为 $D_T = \{(x_{Ti})\}_{i=1}^{n_2}$
 (注意, x_{Ti} 与 x_{Si} 的度量单位相同)
 软件

令 n_1, n_2 分别为源程序和目标程序的文件数量

令 $P(X_S)$ 和 $P(X_T)$ 是 $X_S = \{x_{Si}\}_{i=1}^{n_1}$ 和 $X_T = \{x_{Ti}\}_{i=1}^{n_2}$ 的边
 缘分布为 (X_S, X_T 分别为来自源程序和目标程序)

通常 $P(X_S)$ 与 $P(X_T)$ 可以不同

TCA 的目的是准确地学习到一种 迁移方法 将源
 程序和标程序的数据映射到一个 特征空间?

潜在的 特征空间 使得 $P(X_S)$ 和 $P(X_T)$ 的不同变小, 从而





Mo Tu We Th Fr Sa Su

Memo No.

Date

在 X_S $p(X_S)$ 和 Y_S 训练的模型可以对 $\varphi(X_T)$ 作出预测

Transfer Component Analysis (TCA)

TCA is to learn a transformation φ to map the original data of source and target domains to a latent space where the difference between domains $\text{Dist}(\varphi(X_S), \varphi(X_T))$ is small and the data variance after transformation $\text{Var}(\{\varphi(X_S), \varphi(X_T)\})$ is large

很多现实应用中, φ 可以看作是线性的, 此时 $\varphi(x) = \frac{x^T \Theta}{x \cdot \Theta}$
 $\Theta \in \mathbb{R}^{m \times d}$ 由此将 x 降为 d 维

因而 latent feature $\varphi(X_S) = X_S \Theta$ $\varphi(X_T) = X_T \Theta$

之后用 $X_S \Theta$ 与 Y_S train classifier f
 再用 f 预测 $X_T \Theta$

不同正则方法对 TCA 在不同情况下的性能有影响

来自 扫描全能王免费版

手机上的文档、证件扫描识别利器



扫描快速下载到智能设备

TCA+

就是选取合适的正则方法的TCA
通过一定规则

规则设定基于数据集的特征的相似性

Data set characteristic vector (DCV)
通过计算实例之间的距离 (欧几里得距离)
一对

每个项目测量项目里每对实例的欧氏距离

$$DIST = \{d_{ij} : U_{i,j}, 1 \leq i, j \leq n, i \neq j\}$$

此外还将 numInstances 加入 DCV

$$DCV = (dist_mean, dist_median, dist_min, dist_max, dist_std, numInstances)$$

Source Project DCV: C_S

Target Project DCV: C_T

为了测量 C_S 与 C_T 的距离, 定义了 similarity vector (SV)
并据此为距离分级

(s_{te}) 表示 s 中第 t 个元素的值 s_{te} C_{tree} 同理

给每对对应元素赋值 $(s_{te} \times a < C_{tree} < s_{te} \times b)$

$S \Rightarrow T$, s 与 t 中元素距离的等级 ($b > a > 0$)

来自 扫描全能王免费版

手机上的文档、证件扫描识别利器



扫描快速下载到智能设备

Benchmark Sets

① ReLink 论文中表IV 概括性信息
表VI 为在数据集中被使用的度量

② AFEEA 表V 概括性信息
表VII 使用的度量

Experiment Design (实验设计)

- ① within-project defect prediction (同软件缺陷预测)
- ② 无迁移学习跨软件缺陷预测
- ③ TCA CDD
- ④ TLAt CDD

① ~~一半~~ 数据集将随机分成两半，一部分为数据集train set
一部分为test set

另一种代替方案是使用 10 次交叉验证 (train set:
test set = 9:1)

为了避免片面，将执行100次整个过程，取平均预测结果

来自 扫描全能王免费版

手机上的文档、证件扫描识别利器



扫描快速下载到智能设备



② 将数据集中的项目两两结对，每对有两种情况($a \Rightarrow b$ 或 $b \Rightarrow a$)，而后例($a \Rightarrow b$) a 为训练集， b 为测试集

③ 与②类似，但是应该先用 TLA 训练(为 HL) source and target data set. 而后再训练 ~~训练~~ 预测模型，应用不同规则得到多个 ~~结果~~ ~~同实例~~ 结果 类别

④ 应用 TLA+ 规则和方式选则适宜正则方法再用 TLA 训练 source and target data set 最后训练预测模型

Machine Learning Classifier (机器学习分类器)

均使用 logistic regression (逻辑回归)

特别地使用 LIBLINEAR

线性执行时，我们选用参数 "-S 0" 和 "-B 1"

⑤ F-measure 作为评价指标

