

Prédiction de valeurs manquantes dans les bases de données — Une première approche fondée sur la notion de proportion analogique

William Correa Beltran, H  l  ne Jaudoin, Olivier Pivert

Universit   de Rennes 1 – Irisa, Lannion, France
{William.Correa_Beltran@irisa.fr, Helene.Jaudoin@irisa.fr, Olivier.Pivert@irisa.fr}

R  sum  . Cet article pr  sente une m  thode originale de pr  diction de valeurs manquantes dans les bases de donn  es relationnelles, fond  e sur la notion de proportion analogique. Nous montrons en particulier comment un algorithme propos   dans le cadre de la classification automatique peut   tre adapt      cette fin. Deux cas sont consid  r  s : celui d’une base de donn  es transactionnelle (attributs bool  ens), et celui o   les valeurs manquantes peuvent   tre de type num  rique.

1 Introduction

Dans cet article, nous proposons une solution originale    un probl  me classique de bases de donn  es qui consiste    pr  dire/estimer les valeurs manquantes dans une base de donn  es relationnelle incompl  te. De nombreuses approches ont   t   propos  es pour traiter cette question,    la fois dans la communaut   des bases de donn  es et dans celle de l’apprentissage automatique, fond  es sur des d  pendances fonctionnelles (Atzeni et Morfuni (1986)), des r  gles d’association (Ragel (1998)), des r  gles de classification (Liu et al. (1997)), des techniques de clustering (Fujikawa et Ho (2002)), etc. Nous explorons quant    nous une nouvelle id  e, issue de l’intelligence artificielle, qui consiste    exploiter les *proportions analogiques* (Prade et Richard (2012)) pouvant exister dans les donn  es.

La suite de l’article est organis  e comme suit. Dans la section 2, nous rappelons les notions de base concernant les proportions analogiques. La section 3 pr  sente le principe g  n  ral de l’approche que nous proposons pour estimer les valeurs manquantes, inspir  e par la technique de classification propos  e dans (Bayouddh et al. (2007); Miclet et al. (2008)). La section 4 est consacr  e    une exp  rimentation visant      valuer les performances de la m  thode et    comparer cette derni  re avec une technique classique d’estimation (kNN). Finalement, la section 5 rappelle les contributions principales de l’article et trace quelques perspectives de recherche.

2 Rappels sur les proportions analogiques

La pr  sentation qui suit est tir  e principalement de Miclet et Prade (2009). Une proportion analogique est une proposition de la forme « A est    B ce que C est    D », ce qui sera not   : ($A : B :: C : D$). Dans la suite, les objets A , B , C , et D seront suppos  s   tre des n -uplets