

数字脚印与“社群智能”

郭 斌¹ 张大庆^{2,1} 於志文¹ 周兴社¹

¹西北工业大学

²法国国立电信学院

关键词：数字脚印 社群智能

什么是社群智能

互联网和社会网络服务（Social Network Service）正在快速增长。各种内嵌传感器的移动手机大量涌现，全球定位系统（GPS）接收器在日常交通工具中逐步普及，静态传感设施（如Wi-Fi、监控摄像头等）在城市大面积部署，人类日常行为的轨迹和物理世界的动态变化情况正以前所未有的规模、深度和广度被捕获成为数字世界。我们把收集来的各种数字轨迹形象地称为“数字脚印”（Digital Footprints）。通过对这些数字脚印进行分析和处理，一个新兴的研究领域——“社群智能”（Social and Community Intelligence）正在逐步形成。

社群智能的研究目的在于从大量的数字脚印中挖掘和理解个人和群体活动模式、大规模人类活动和城市动态规律，把这些信息用于各种创新性的服务，包括社会关系管理、人类健康改善、公共安全维护、城市资源管理和

环境资源保护等。图1通过一个“智慧校园”的例子展示了社群智能给我们的工作和生活带来的影响。

在大学校园里，学生A经常会遇到一些困扰：当他想去打球时，不知道谁有时间能陪他去玩；要去上自习时，不知道在哪个教学楼里可以找到空位。另外，作为人口密集场所，当严重流感如H1N1来袭时，如何寻求有效办法限制其传播？当确定B患上某疑似病例后，需要及时地把最近接触过B的人找到。在现有条件下，获取这些有关个人活动情境、空间动态、人际交互的信息还没有较好的



图1 社群智能在大学校园内的应用

技术解决方案,需依赖耗时且易出错的人工查询来完成。例如,A需要通过电话或网上通讯方式和多个朋友联系,来确定谁可以一起去打球。社群智能的出现将改变这一切。上面提到的问题都可以通过分析来自校园的静态传感设施和移动电话感知数据(蓝牙,加速度传感器等)以及发布在社会万维网(Web)上的人与人之间关系信息来解决。以流感防控问题为例,记录谁和B接触过、接触时的距离以及时间长短、社会关系(如亲戚、朋友或陌生人)等是非常重要的,这些信息可以通过分析移动电话感知数据得到。

社群智能是在社会计算(Social Computing)^[1]、城市计算(Urban Computing)^[2]和现实世界挖掘(Reality Mining)^[3]等相关领域发展基础上提出来的。从宏观角度讲,它隶属于社会感知计算(Socially-Aware Computing)范畴^[7]。社会感知计算是通过人类生活空间逐步大规模部署的多种类传感设备,实时感知识别社会个体行为,分析挖掘群体社会交互特征和规律,辅助个体社会行为,支持社群的互动、沟通和协作。社群智能主要侧重于智能信息挖掘,具体功能包括:(1)多数据源融合 即要实现多个多模态、异构数据源的融合。综合利用三类数据源:互联网与万维网应用、静态传感设施(Static Sensing Infrastructure)、移动及可携带(Wearable)感知设备,来挖掘“智能”信息;(2)分层次智能信息提取 利用数据挖掘和机器学习等技术从大规模感知数据中提取多层次的智能信息:在个体(Individual)级别识别个人情境(Context)信息,在群体(Group)级别提取群体活动及人际交互信息,在社会(Social or Community)级别挖掘人类行为模式、社会及城市动态变化规律等信息。

社群智能研究的演化过程

人类早期对自身日常行为模式、社会交互及城市动态变化规律的研究是通过调查问卷、个人观察等物理方式实现的^[3]。由于采样范围小、可获得的数据量少且工作量大,得到的结果往往是片面的、

不准确的,生成结果的时间也比较长。计算机、互联网、无线通信及传感技术的出现和发展,为实施大规模、近乎实时的人类行为模式、社会交互及城市动态变化规律的研究提供了可能。随着技术的演进,“社群智能”相关的研究形成了三个不同的方向:互联网与万维网挖掘,静态传感器感知、移动及可穿戴计算。

20世纪90年代以来,各种互联网服务,如电子邮件、实时消息和万维网的普及极大地改变了人类交流和获取信息的方式。这一阶段的研究主要集中在信息检索和信息抽取方面,例如新闻推荐、个人或企业信息提取等。也有少量针对人类交互方式进行的研究,比如有人通过大量的电子邮件记录来分析人们之间的关系及交互模式^[18]。随着互联网进入Web 2.0时代,研究人员开始把注意力转移到各种新型的在线社区服务,如社会网络服务网站、维基(Wiki)、博客等。Web 2.0时代最大的特点是增强了用户参与度,这为研究社会行为和城市动态信息提供了由用户生成的海量数据。这一阶段的主要研究热点包括社会网络分析(Social Network Analysis)^[4]、人际关系预测^[5-6]、社会和灾难事件监测^[7-8]等。中国科学院自动化研究所的王飞跃研究员最早于2005年在IEEE智能系统杂志上把基于互联网和万维网的社会学研究称为“社会计算”^[1]。

社会计算主要是针对人们在虚拟世界的交互行为进行的。随着各种静态传感设施(如监控摄像头、室内定位系统、射频识别(RFID)等)的出现和普及,感知人们在现实世界中的行为和社会交互成为可能。在早期阶段,传感器主要用在一些重要的地点进行环境和异常事件监测,如森林火险报警等。技术的发展逐步实现了传感器的小型化和廉价化,人们开始把传感器部署到日常生活环境中创建各种“智能空间”(Smart Spaces)^[9],这为感知和分析室内环境中人们的活动和交互行为打开了大门。例如,AT&T剑桥实验室的研究人员使用超声波技术对室内物体和人进行定位^[10]。英特尔西雅图研究中心最早把射频识别标签贴附在各种室内物体上(如牙刷、椅子等),利用人和物体的交互来

识别人们的活动情况（如刷牙，吃饭等）^[11]。基于静态感知设施的人类行为和环境感知的最大不足在于其范围局限性，只能在为固定传感器所覆盖的环境下进行。

可穿戴传感器的出现弥补了静态传感设施的不足。可穿戴计算通过把各种小型传感器，如加速度传感器、心搏传感器、无线摄像头和微型麦克风等“穿戴”在人

身体上，把人转变成“可移动的传感器”（Mobile Sensor）以感知其个人行为、健康状况、活动情境（如在开会、在和朋友谈话等）和周边环境信息（噪音强度，亮度等）。虽然可穿戴传感器携带方便，但是还没有得到普及，不足以支持大规模数据的采集与分析。随着配备各种传感器（如GPS、蓝牙、WiFi，加速度传感器等）的移动电话的出现，情况进一步发生了改变。由于移动电话已经得到了广泛的应用，它们收集到的大量多模态数据为分析大规模人类行为模式、社会和群体动态信息开辟了一条新的现实可行的途径。从2006年开始的麻省理工学院“实时罗马”项目率先利用大规模移动电话数据来分析城市动态信息，如人的移动模式、城市热点区域（Hot Spots）随时间变化的规律等^[12]。来自麻省理工学院的另一个项目“现实世界挖掘”，同样利用移动电话感知数据分析人与人之间关系^[3]。美国达特茅斯（Dartmouth）学院安德鲁·坎贝尔（Andrew T. Campbell）研究组也在其后提出了“人本感知”（Human-Centric Sensing）的概念，通过以人为中心的移动电话感知来进行社会关系分析和周边环境监测^[13]。

其中，“社群智能”综合利用三种信息源（互联网和万维网应用、静态感知设施、移动和可穿戴设备）研究个人、群体及社会行为模式，人与人之间的关系，以及大规模人群和社会的动态变化规律。从表1可以看出，前面提到的几个相关研究领域与社群智能最大的不同点是，它们都只依赖单一

表1 社群智能及相关研究领域定义

研究领域	定义
社会计算 (Social Computing)	通过人类与互联网和Web的交互来进行社会学及人们交互行为的研究。
现实世界挖掘 (Reality Mining)	通过对移动电话感知数据对人们的社会行为（包括人与人之间关系，社会行为模式等）进行研究。
人本感知 (Human-Centric Sensing)	通过对移动电话中得到的数据进行挖掘来获取人的活动，交互及周边环境信息。
城市计算 (Urban Computing)	通过传感器网络技术来研究人与环境（包括城市，公园，森林等）的交互及环境的变化。
社群智能 (Social and Community Intelligence)	综合利用三种信息源：互联网和Web应用、静态感知设施，以及移动和可穿戴设备，来研究个人、群体及社会行为，人与人之间的关系，以及大规模人群和城市的动态信息。

类别数据源进行某个方面信息的获取。例如，社会计算强调从万维网数据中发掘人们在虚拟世界的交互信息，但不关心物理世界的人类交互。与社会计算一样，现实世界挖掘和人本感知主要也是分析人们之间的社会关系，但它仅仅依赖于来自移动电话的数据。城市计算则主要依靠静态感知设施来研究人与环境的交互及环境动态信息。与这些领域不同，社群智能综合利用人类与信息物理空间内（Cyber-Physical Spaces）多种信息源交互留下的数字脚印，以挖掘更为广泛的情境信息。从小的角度讲包括个人情境、小范围群体行为、周边环境信息，从大的角度讲包括大规模人群、城市及社会的动态变化情况和规律（如交通阻塞、突发事件、热点地区监测等）。

三种不同信息源数据融合的作用

社群智能的目标是综合利用互联网和Web应用、静态感知设施、移动和可穿戴设备的信息源来挖掘个人、群体及社会动态信息。三类信息源具有各自的特点：互联网和Web应用中包含更多静态的或不常改变的信息，如个人或企业基本信息，网络社区中人与人之间的关系等；静态感知设施能在部署传感器的智能空间中感知个人或群体的行为以及空间状态信息；移动和可穿戴计算以人为中心，能够在没有静态感知设施的环境中感知人的行为及位置、人



图2 表现多信息源融合作用的三个例子

们之间关系、社会情境 (Social Context) 和周边环境信息。

由于这三种信息源具有各自的特点, 对它们进行集成和融合后可以起到单一类别信息源所无法达到的效果。图2给出三个例子来说明多信息源融合的效用。

通过万维网获取的知识来辅助基于传感器的行为识别 从万维网获取的社会关系信息, 可以用于物理世界中的群体行为识别。例如, 通过传感器可以感知到某空间内有很多人, 这些人正在从事什么活动却很难由传感器感知到。但如果知道这些人是朋友关系, 则可以推测这里正进行社交聚会; 如果是管理者和员工的关系, 则很有可能是在开会。

通过传感器感知到的人类交互信息来增强在线社会关系网社区建设 当前社会关系研究主要依靠用户输入 (朋友列表设定、相互间留言及评论等) 来推测人们之间的关系。由于用户输入的信息往往是不完全的、主观的, 挖掘出的社会关系往往也是不完全的、粗粒度的 (往往只能推测两人间是否存在关系)。如果用传感器对物理世界中人们的交互轨迹进行记录, 则可以很客观且较为全面地揭示人们之间的关系, 更精确地推测关系的类型、强弱等。例如, 两人仅仅在工作场合见面, 则可能仅仅是同事关系; 如果在工作场合之外见面, 则很

可能是亲密的朋友关系。麻省理工学院科学家亚历克斯·彭特兰 (A. Pentland) 研究组利用移动电话蓝牙设备来感知用户之间的邻近 (Proximity) 情况, 通过分析验证和描述了在线社会关系网中人与人之间的关系^[14]。

融合移动感知和万维网数据来提供公共服务 来自不同信息源的数据常常只能反映自然和社会事件的某个方面, 对这些数据进行融合后可以更好地刻画某个事件的全貌。例如, 通过对从Twitter用户所发布的帖子中挖掘主题信息, 综合利用从移动电话中获取的GPS地址信息, 日本东京大学的研究人员已经成功地开发了接近实时的地震报告服务^[8]。

社群智能的体系架构

社群智能的基本体系架构共分为五层 (见图3)。其中, 感知层 (Sensing Layer) 负责从三种信息源来获取原始数据。由于这些原始数据极有可能暴露用户的行踪和隐私, 数据处理之前需要通过数据匿名保护层 (Data Anonymization Layer) 进行匿名化的工作; 混合学习层 (Hybrid Learning Layer) 采用各种机器学习和数据挖掘算法将原始感知数据转换为高级特征或情境 (Context) 信息; 语义推理层 (Semantic Inference Layer) 与混合学习层相辅相成, 通过基

于专家知识明确设定好的逻辑规则对不同的特征或情境信息做进一步的集成,并最终得到社群智能信息;应用层包含大量的社群智能应用程序,它们可以从社群智能库中查询自己需要的信息来提供各种创新服务。

创新应用领域

社群智能为开发一系列创新性的应用提供了可能。从用户角度来看,它可以开发各种社会关系网络服务来促进人与人之间的交流。从社会和城市管理角度来看,它可以实时感知现实世界的变化情况来为城市管理、公共卫生、环境监测等多个领域提供智能决策支持。

社会关系网服务

通过对人们在物理或数字世界中进行的各个方面的交互(如见面对象、时间及地点、聊天内容等)进行记录,并对用户行为模式进行挖掘(如用户兴趣),社群智能为大规模开发社会关系网服务提供了基础支持,如朋友推荐^[6]、基于行为感知的增强型在线交互^[13]等。

城市计算

当无线传感器设备遍布于整个城市空间及城市居民身上时,就可以利用收集来的大规模数据解决一些日益严重的城市问题:如城市动态信息监测、交通规划、公共设施管理等。例如,美国麻省理工学院的“实时罗马”项目,利用从罗马市内收集来的移动电话、公交车及出租车信息来分析城市热点地区动态变化规律^[12]。美国加州大学洛杉矶分校的Biketastic项目(<http://biketastic.com>)利用城市大量自行车的运行轨迹、运行速度等为其他自行车用户进行路径推荐。

环境监测

在社群智能感知过程中,配有移动或可穿戴设备的人成为一个特殊的信息获取源,称之为感知(Citizen Sensing)。人本身具有的移动性和参与性,为在静态传感设施没有覆盖的区域

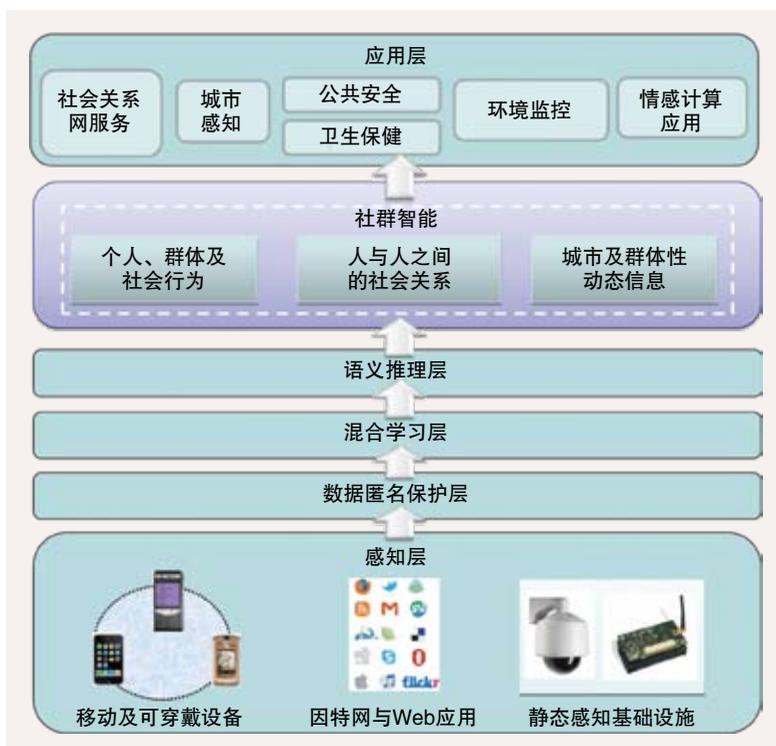


图3 社群智能的体系架构

进行环境感知创造了机会。美国麻省理工学院的“猫头鹰”项目(<http://web.mit.edu/newsoffice/2008/tracking-0822.html>),利用安装在智能移动设备上的GPS和指向传声器等,在人移动的过程中估计森林中猫头鹰的分布情况。美国达特茅斯学院在这方面也做了很多研究,BikeNet项目利用安装在自行车的二氧化碳检测器来测量自行车行走路线上的空气污染情况^[13]。

卫生保健

社群智能还可以用于公共卫生和个人健康辅助。例如,谷歌研究人员对用户提交的健康搜索关键词进行分析,推测全球各地区流行病的实时传播情况,该研究成果发表在2006年的《自然》(Nature)杂志上^[15]。美国休斯顿大学的Neat-o-Games系统通过可穿戴的加速度传感器感知人的运动情况(如走路、跑步等),利用感测到的运动量来控制虚拟社区中参加赛跑游戏的用户“替身”(Avatar),以此促进用户进行更多的运动^[16]。

公共安全

利用多种信息源,可以对影响公民安全的事件,如恐怖袭击、自然灾害等进行预警或

及时响应。例如,安装在建筑物内或街道上的视频监控设备为检测各种异常或突发事件提供了支持。美国波士顿的警察通过分析互联网社区中用户发布的数据来进行犯罪预防^[7]。日本的研究人员利用Twitter用户发布的信息进行地震实时报告^[8]。

社群智能研究面临的挑战

作为一个新的研究领域,社群智能在感知、数据管理和智能信息抽取等多个方面都面临着新的问题和挑战。

参与感知还是机会感知 在社群智能中,人类将成为一个重要的信息载体,人应该承担什么样的角色呢?比如,当移动电话成为感知源时,是否需要打断使用者的工作,让他控制(接受或中止)某个感知任务呢?目前有两种解决方案^[13]:一是**参与感知**,让用户参与到感知决策中来,判定是否要共享数据,是否参与某任务的控制,以及如何对不同数据设置安全级别等;二是**机会感知**,它不需要用户主动参与到感知决策中来,而是在设备所处状态(如所处位置,资源可用性等)满足感知任务需求的情况下自动开启感知任务。这两种方案各有优缺点。参与感知虽然具备较为可靠的安全隐私控制,但需要频繁向用户发出请求从而会降低用户的使用意愿;机会感知虽然减轻了用户负担,但却增加了安全风险及设备的资源开销(需要消费大量资源来做决策)。就社群智能感知来说,还需要做很多工作,在涉及的各个因素方面寻求一个平衡解决方案。

数据安全、质量和可信度 传统计算机安全领域主要对虚拟个人账户等进行保护,社群智能可能暴露更多物理层面的个人信息,包括当前活动、所处位置、兴趣爱好等。如果没有很好的隐私保护制度,用户将不愿意共享自己的数据,也将使得社群智能的提取无法进行。目前有两种隐私保护方法:一是采用数据匿名技术,如美国达特茅斯学院的MetroSense系统采用k-匿名方法来保护用户提交的位置信息^[13];二是增强用户控制功能,使其能够对数据进行授权管理。数据质量和可信度是社群智

能数据管理的另一个重要问题。例如,万维网上收集的用户提交的数据很多是不准确的,甚至是虚假的。从移动设备或静态感知设施中得到的数据在质量上也有很大的差别。例如,移动电话放置在口袋里还是在用户手中,所获取的数据质量会不一样,对用户活动的识别准确率也会造成很大影响。因此,需要开发可信计算和异常数据监测方法来保证收集数据的质量和可信度。

智能信息抽取 这一过程是要通过多种数据源收集的数字轨迹挖掘高级智能信息,如个人情境、社会事件、人与人之间关系、带语义的位置(如在火车上、在闹市区)、城市动态信息等等。这个把原始数据转换到高级智能信息的过程称为智能信息抽取。这里涉及到两个问题:

1. 智能信息抽取框架的结构设计

该框架需要对来自大量感知节点的大规模实时数据进行分析处理。有些信息可以从单个节点得到(如个人活动),但有些则需要对某空间内多个或更大范围分布的节点得到的数据进行融合(如城市动态信息等)。有两种常用结构设计方式可以采用:

分布式结构 在每个节点进行个人语义信息处理,处理的结果发送到服务器端进一步挖掘群体性语义信息,其缺点是对感知节点的处理能力要求很高,会占用大量本地资源;

云计算结构 把所有数据传送到服务器端进行统一处理,其缺点是通信要求比较高。

这两种设计方案都有各自的优点,但都不能很好地满足社群智能系统的要求。未来需要综合考虑社群智能系统中网络通信和感知节点处理的代价,各种智能信息的特点,寻求一个较为平衡的且能综合多种结构优势的设计方案。

2. 逻辑推理方法和统计机器学习方法如何有机结合

逻辑推理和统计机器学习是智能信息抽取的两种主要方法。逻辑推理由用户定义的规则库和推理器两部分组成。规则一般通过专家或用户经验设定,实现比较简单。缺点是由于基本假设过于理想,会产生冲突和不确定性问题。统计机器学习一

般需要对样本进行训练以推测新的样本类型。在社群计算背景下,由于训练数据获取的困难、人类个体行为的差异、设备所处环境的变化等多方面的因素,对统计分类系统的性能会造成很大的影响,需要对其进行进一步研究。

同时,社群智能在多模态数据管理与建模、大规模实时数据处理算法设计(采样优化、问题分解等方面)、情境感知不一致性处理等方面还面临很多问题需要解决。

结语

社群智能正在开辟一个新的多学科交叉研究领域。随着越来越多的“数字脚印”被收集到,社群智能的研究和应用范围在未来一段时间内会进一步扩展和延伸。作为一个新兴领域,社群智能的普及和发展还面临着感知过程控制、多模态数据管理、复杂智能信息抽取、安全隐私保护等多个方面的挑战,也为理论基础和创新性研究提供了新的机遇。虽然当前在社群智能方面的实践还主要面向单数据来源,人们期望在不久的将来会看到更多的三种数据来源融合的研究成果及创新应用。■

致谢

本文得到欧盟第七框架计划(EU FP7)、国家自然科学基金(60903125)、国家863高技术研究发展计划(2009AA011903)和教育部“新世纪优秀人才支持计划”(NCET-09-0079)资助。



郭斌

西北工业大学副教授。主要研究方向为普适计算、人机交互和社会计算。
guobin@nwpu.edu.cn



张大庆

法国国立电信学院教授,西北工业大学兼职教授。主要研究方向为普适计算、情境感知系统和服务计算。
Daqing.Zhang@it-sudparis.eu



於志文

CCF高级会员、本刊编委、2006CCF优秀博士学位论文奖获得者。西北工业大学教授。主要研究方向为普适计算、情境感知系统和智能信息技术。
zhiwenyu@nwpu.edu.cn



周兴社

CCF常务理事。西北工业大学教授。主要研究方向为嵌入式计算、普适计算和网格计算。
zhouxs@nwpu.edu.cn

参考文献

- [1] F. Wang et al., Social Computing: From Social Informatics to Social Intelligence, IEEE Intelligent Systems, vol. 22, no. 2, 2005, pp.79~83
- [2] K. Tim; C. Matthew, P. Eric, “Urban Computing”, IEEE Pervasive Computing, vol.6, no. 3, 2007, pp.18~20
- [3] N. Eagle, et al., Inferring Social Network Structure using Mobile Phone Data, Proceedings of the National Academy of Sciences (PNAS), vol. 106, no. 36, 2007, pp. 15274~15278
- [4] S. Staab et al., “Social networks applied,” IEEE Intelligent systems, vol. 20, no. 1, 2005, pp. 80~93
- [5] R. Xiang et al., Modeling relationship strength in online social networks, Proc. of WWW 2010 Conf., 2010
- [6] D. Quercia et al., Nurturing Social Networks Using Mobile Phones, IEEE Pervasive Computing, 2010
- [7] A. Sheth et al., Citizen Sensing, Social Signals, and Enriching Human Experience, IEEE Internet Computing, vol. 13, no. 4, 2009, pp. 87~92
- [8] Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo. Earthquake Shakes Twitter Users: Real-time Event Detection by Social Sensors, Proc. WWW 2010 Conf., ACM, 2010
- [9] X. Wang, et al., Semantic Space: An Infrastructure for Smart Spaces, IEEE Pervasive Computing, vol. 3, no. 3, 2004, pp. 32~39
- [10] H. Andy et al., The Anatomy of a Context-Aware Application, Wireless Networks, vol. 8, no. 2-3, 2002, pp.187~197
- [11] M. Philipose et al., Inferring ADLs from interactions with objects, IEEE Pervasive Computing, vol. 3, no. 4, 2004

- [12] J. Reades, F. Calabrese, A. Sevtsuk, C. Ratti, Cellular Census: Explorations in Urban Data Collection, *IEEE Pervasive Computing*, vol. 6, no. 3, 2007
- [13] A. T. Campbell et al., The Rise of People-Centric Sensing, *IEEE Internet Computing*, vol. 12, no. 4, 2008, pp. 12-21
- [14] A. Pentland, Socially Aware Computation and Communication, *IEEE Computer*, vol. 38 no.3, 2005, pp. 33-40
- [15] N.M. Ferguson et al., Strategies for mitigating an influenza pandemic, *Nature*, vol. 442 no, 7101, 2006, pp.448-452
- [16] Y. Fujiki, et al., NEAT-o-Games: Blending Physical Activity and Fun in the Daily Routine, *ACM Computers in Entertainment*, vol.6 no. 2, 2008
- [17] 於志文, 周兴社, 社会感知计算, *中国计算机学会通讯*, Vol. 6, No. 9, 2010
- [18] Kossinets, G., and Watts, D.J. Empirical analysis of an evolving social network, *Science*, 311(5757), January 2006, 88-90