

Crowdsourcing: Descriptive Study on Algorithms and frameworks for prediction

Dhinakaran K¹, R.Nedunchelian²

¹Research Scholar, Assistant Professor, Department of Computer science and Engineering,
Rajalakshmi Institute of Technology, Anna University.

²Professor, Department of Computer science and Engineering, Sri Venkateswara College of Engineering,
Anna University,
maildhina.k@gmail.com , nedun@svce.ac.in

Abstract

The three inter-related process data mining, data analytics and data processing are carried out on large volume of datasets. The value of the data may be text, numeric, ontology, alpha-numeric, images, video, and other multi-dimensional datasets. The dataset from the people is the one of the famous dataset from the above datasets. Crowdsourcing is used to solve the large size of data with people. The crowdsourcing input will be from a group of people by collecting a large number of people and analysis it is one the emerging technology, which initiate a new model for big data mining process. To define the data, data mining is the one of the traditional process for the exert in analytics domain. Data mining is an expensive process and it also take long time to complete the process. In industry and research area crowdsourcing has become a very active. The crowdsourcing uses the smart phone users as volunteers and share their annotation process for different type of contributions. This paper is used to review about the bigdata mining from crowdsourcing in recent years. Using crowdsourcing the opportunities and challenges of data analytics are reviewed, and summarize the data analytics framework. Then it is discussed several algorithms of including applications, cost control, quality control, latency control and big data mining framework which must be consider in the field of crowdsourcing. Finally, the conclusion of this project tell about the data mining limitation and give some suggestions for future research in crowdsource data analytics.

Keywords: Crowdsourcing, Bigdata analytics, Data mining.

Introduction

Crowdsourcing is used to connect the large number people through the internet. Using distributed knowledge it is used to solve many problems and produce large things by connecting people through internet. The information is collected from large number people to sov and complete business based tasks. Crowd Sourcing (CS) is a practical experiment to carry out a crowd of data for solving problems. It is used in novel and latest technologies, social media and web-2.0. CS is applied in different stages and several industries. CS is connected with optimizing tasks, consumer management, novel ideas and so on. It makes closeness among the organization, social media, new collaborations and stakeholders. The millions of peoples are connected to the internet suggestions, information regarding small or big projects and sharing their own ideas. CS can deliver local information and solve tricky problems . If the information and decision used in CS are right, then the using the wisdom of CS is easy. CS enter into interaction of e-business and all social networks . It also makes changes in research, create, human work and market. The innovations are applied on the healthcare problems, citizens empowered are democratized. The crowd data analytics fields are used since CS as more advantages. This paper provides a detailed survey about crowdsourcing and gives better understanding about the applications.

Howe, (2006), Faridani et al. (2009), proposed a model where different type of technology are which are available in the internet, in this model anyone to collect and persist extraordinary amount of data. There are many models used for analysing, concluding and learning on the data. Pattern mining in big data, the crowd sourcing is considered as a better model for analysing. To conclude the data and extract high relevant pattern the data miners makes use of better tools. To understand the various type of model we need to study about the CS. Howe, (2006) said a group of customers will perform some task based on outsourcing the CS is used. The assigning of task can be done in offline or in online. The requester who is going to work on the task allocation should not have any prior knowlegde about the assignment of task allocation. In recent days, CS is applied in image labelling, Faridani et al. (2009), VonAhn (2006), question answering, eliciting inspired works and solution providing for scientific problems, and various micro tasks like posting on CS markets such as AMT. MobileWorks is the most successful idea generation through the CS platforms.

Many data mining tasks cannot be performed efficiently by using existing machine-only algorithms, which includes image classification [1], sentiment analysis [2], and opinion mining [3]. For We have to cluster things based on some category like place, cluster them based on the country where they belong to. It is easy to identify the places with the human knowledge but it is hard for machines to understand the places based on the pictures. But by using the information from crowd source it is easy for machines to know about the places. Really thanks to the public CS platforms, e.g., AMT and CrowdFlower, the accessing the crowd became very easy. Fig 1 show the processing crowdsourcing application.

A crowd of data for solving problems is called as Crowd Sourcing (CS). It is used in novel and latest technologies, social media and web-2.0. CS is applied in different stages and several industries. CS is connected with novel ideas, optimizing tasks, consumer management and so on. It makes closeness among the social media, organization, stakeholders and new collaborations. Using crowd millions of customers are connected in the internet, sharing their own ideas, suggestions and information regarding small or big projects. CS can solve tricky problems and deliver local information. If the information and decision used in CS are right, then the using the wisdom of CS is easy. CS enter into all social networks and interaction of e-business. Also, it makes changes in human work, research, create and market. Citizens empowered, healthcare problems are democratized and apply innovations. Because of the advantages of CS is more, the crowd data analytics fields are using CS. This paper provides a detailed survey about crowdsourcing and gives better understanding about the applications.

Howe, (2006), Faridani et al. (2009), presented, various technologies available in the internet make anyone to collect and persist extraordinary amount of data. Also, it enables different models for learning, analysing and concluding on the information. Crowdsourcing is considered as a better model in this work, for analysing big-data regarding pattern mining. Bigdata miners can collaborate with efficient tools can make them to extract high relevant patterns and draw conclusion over the data. In order to understand various models and methods are studied for crowdsourcing. Howe, (2006) said that CS is a technique for outsourcing the task to a set of people. Task assigning can be done in online or in offline. Before task allocation in online or in offline, the requester who is going to work on the task are not known about the task. In recent days, CS is applied in image labelling, Von Ahn, (2006), Faridani et al. (2009), producing inspired works, question answering, and solution providing for scientific problems, and various micro tasks like posting on CS markets such as Amazon Mechanical Turk.

Many data mining tasks cannot be performed efficiently by using existing machine-only algorithms, which includes image classification [1], sentiment analysis [2], and opinion mining [3]. For instance, given a set of pictures of famous places of interest in the world, we want to cluster them according to the country they belong to. With human knowledge it is easy to categorize the pictures into countries like “India” or “America”, but it is quite hard for machines to do so. Favorably, by making use of huge numbers of ordinary workers which is the crowd, crowdsourcing has risen as an effective way to attempt such machine-hard tasks. Really thanks to the public crowdsourcing platforms, e.g., AMT and CrowdFlower, the access to the crowd becomes easier. Fig 1 show the processing crowdsourcing application.

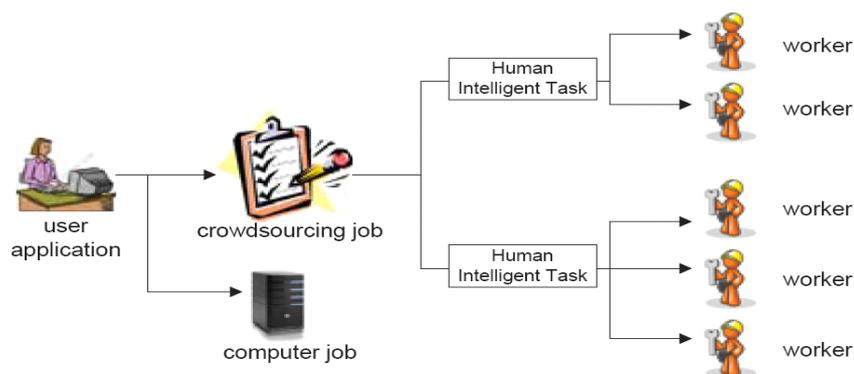


Fig 1. Crowdsourcing Application

In the field of research and industry the crowdsourcing has become a lively area. In a CS platform example AMT, a crowdsourcing platform is used for to the “requesters” publish tasks on , while “workers” perform some tasks and result will be return.. The images are classified into a hierarchical order for classification of the problem, the requester has to perform the "task design" and example the user interface will be designed of a task (e.g., providing workers with an image and a category, and asking them to check whether the image is a subset of the category), and some functionalities will be set up for the tasks example a task price, the total number workers who will answer for the task, time duration to perform the task. The platform is used to publish the requesters task. They should accept,answer and submit the task to that same platform.. The answer will be collected from the platform and the report that to asker. If the task has been completed by the worker, the requester can disapprove or approve the answer from the workers side, from the requester only the approved workers will get their payment. The common characteristics and models of crowdsourcing models are compared in Table 1.

Table 1. Comparison of Common Characteristics of Crowdsourcing Techniques.

Common Characteristics ----- Modes of Crowdsourcing	Cost	Anonymity	Scale of Crowd	Implementation Time	Task Quantity	Crowd Reliability
Virtual Labor-Markets	Flexible	High	High	Low	Simple	Medium
Tournament-Based Collaboration	Fixed	Medium	Medium	Medium	Complex	High
Open Collaboration	Free	Medium	Flexible	Flexible	Flexible	Flexible

The next thing in crowdsourcing is allocating task to the workers. The task allocation process is clearly defined in Fig 1. The tasks are classified into following methods based on task properties:

Task Specificity: Task specificity ensures that a task is well defined before publishing in a crowdsourced environment.

Task Complexity: This is to indicate the knows`ledge, experience and amount of skills that will help to solve CS task. The range of complexity is from simple to complex.

Task Contribution type: This type indicates that to whom that this task will specifically assigned to solve the crowdsourcing task. The task can be performed in an individual or collaborative manner.

Task granularity: Task granularity is in any crowdsourced task is assigned for individual or collaborative manner whether the task is divisible to micro-task or not.

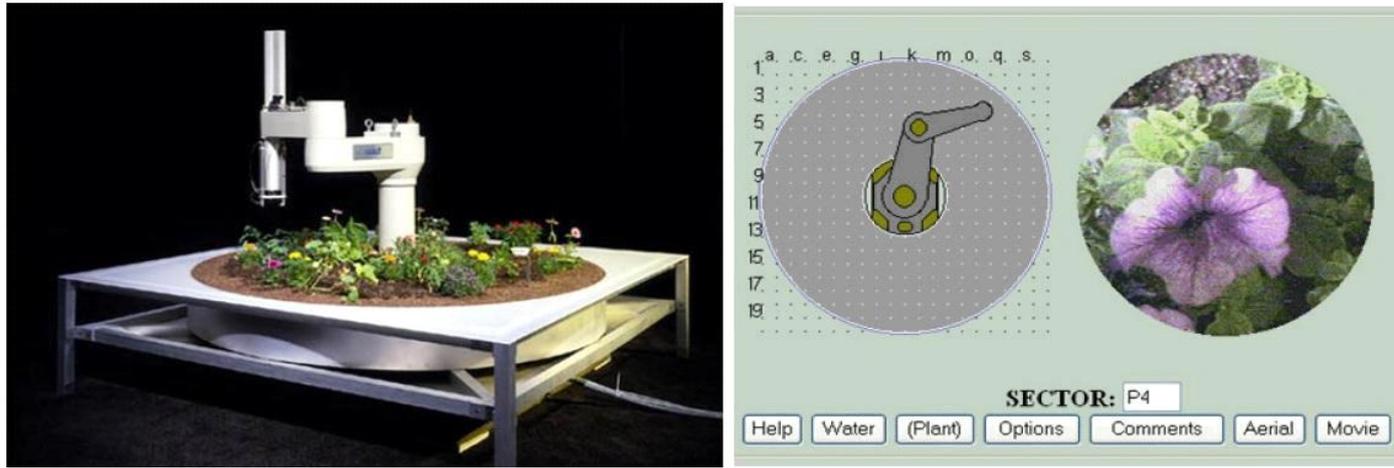
Task requirement: Task requirement is based on any human or computing applications that is required to complete the crowdsourced job.

Task incentives: Task incentives are purely based on the size of the crowd and the workers will be paid either money or coupons.

Task Problem Type: The type of problem in crowdsourcing is based on data. for example, data processing, Location based tasks or annotation task.

Real Time Applications

Some of the crowdsourcing applications used in the real-time environment are (Dahl, (2007a), CONE, (Rappole and Faridani. (2011)), Gamification, (Deterding et al. (2011b), BWL system – (Bird Watching Learning) (Chen et al. (2003)), frame selection algorithm, Deterding et al. (2011a)), Faridani et al. (2009) and Song et al. (2008). These environments are used in real time industry effectively.



(a). Telegarden Robot

(b). Web-based Interface of the Telegarden

Figure-2. Crowdsourcing Applications Based on CONE

Telerobotics

After the creation of WWW with a short interval of time, the internet based telerobotics is a project is created. Online softwares is Telegarden, Goldberg et al. (1995) where online users can do watering and planting over internet using web-based-interface. The entire amount of contestants is 9000 in the year of 2004. Marin et al. (2005) is the biggest telerobotic project. The frontview of the Telegarden robot and the web-based interface is shown in Figure-2. Goldberg and Siegwart, (2002) presented a detailed survey about various telerobotics and web-based robotic projects. Some of the telerobotics used in various fields are telerobotic-surgery (Arora et al. (2011), Tele-actors (Goldberg et al. (2003)), Haidegger et al. (2011)) and explosive handling and multi-user-robotic cameras for online-video conferencing (Kimber et al. (2002), Liu et al. (2002)). Some of the latest examples are (Kim et al. (2002), Song et al. (2008b), Schiff et al. (2009). For more discussion of collaboration procedures on shared robots like (Xu, (2012), Corredor and Sofrony, (2011)).

The frame selection problem

The CONE a frame selection algorithm is used in order to choose a telerobotic camera for multiple participants. Song et al. (2006), this algorithm shares camera among multiple users participating in the telerobotics. Multiple participants can share a single robotic camera. Web-based-interface, is used for N number of participants who submit their frame request to camera. The frame-selection algorithm satisfies the participants by optimal frame allocation (Song, (2009a)). Only 65% of participants can satisfy by this algorithm. Therefore, new algorithm were proposed by Song et al. for new frame distribution algorithm for frame selection and allocation.

Robotic Avian Observation

Chen and colleagues implemented BWL system to capture and classify the images using PDAs connected in WLAN (Chen et al. (2003)). The author stated that children can improve their learning ability in BWL. The various kinds of streaming is not possible in the remote locations, the bandwidth of the internet connection is not sufficient. The solution for this problem is proposed by Dahl, (2007b) as relay server. The relay-server allows only the user who got permission in the network by converting the frames into JPEG format and sends it to the users. Hence the bandwidth is managed in the network potentially.

Gamification

One of the important reasons for using a gaming element in a non-gaming context is defined by Deterding et al. (2011a, 2011b). The CONE-Welder used a crowdsourced avian classification model and it is inspired by Google's image labelling by Weber et al. (2008) and the Von Ahn et al. (2006). One image shared for two participants online as well as offline simultaneously. Therefore, the sharing and time are improved. From the above application's development and deployment, it is understood that there are some merits and

demerits faced by the earlier researchers. In order to understand the problems, a general challenge faced in crowdsourcing is given here.

Challenges in Crowdsourcing

The crowd differs from the machines by some characteristics. (1) Not Free. Some payment will be done for the workers who will answer for the task, and it is important for cost monitoring. (2) Susceptible to errors. Workers may provide some noisy outcomes and authorization will be done for that noisy results and quality will be improved. Moreover, workers have collective contextual data, leading to diverse accuracies to answer a profusion of tasks. We need to capture workers' characteristics to achieve high quality. (3) Dynamic. All the workers won't be available online to answer tasks and the latency should be controlled by us. Thus, three fundamental methods must be considered in crowdsourcing: "quality control", "cost control", and "latency control".

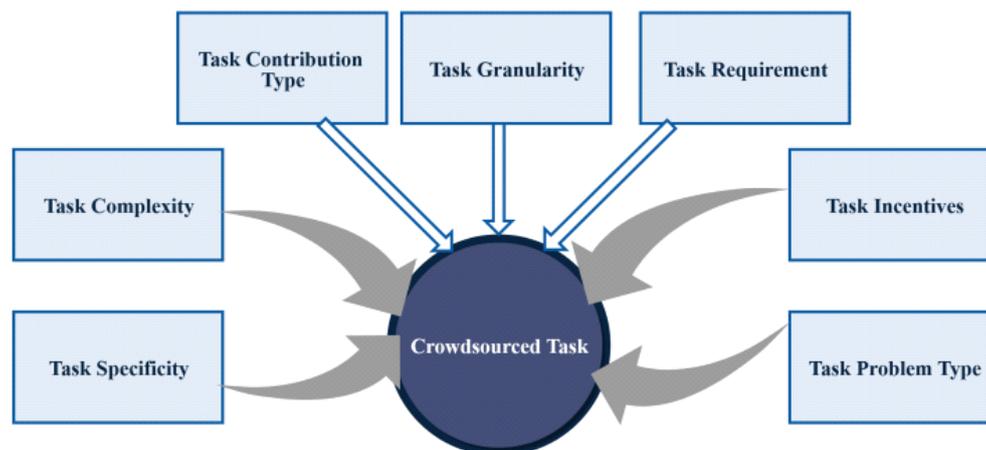


Fig 3 Crowdsourced task features

It mainly only focuses on the quality control on generating superior answers from workers (possibly noisy answers), by differentiating a aggregating workers' answers and worker's quality [5]. Cost control aims to maintain the good result quality by reducing human costs[6]. Latency control achieves on reducing the latency by modelling and measuring the worker's arrival rate and their latency [7]. Note that there are compromises among quality, cost, and latency, and current training motivation is on how to stabilize them, e.g., optimizing minimizing the cost under latency and quality constraints, lessen the latency given a static cost, the quality given a static cost etc.

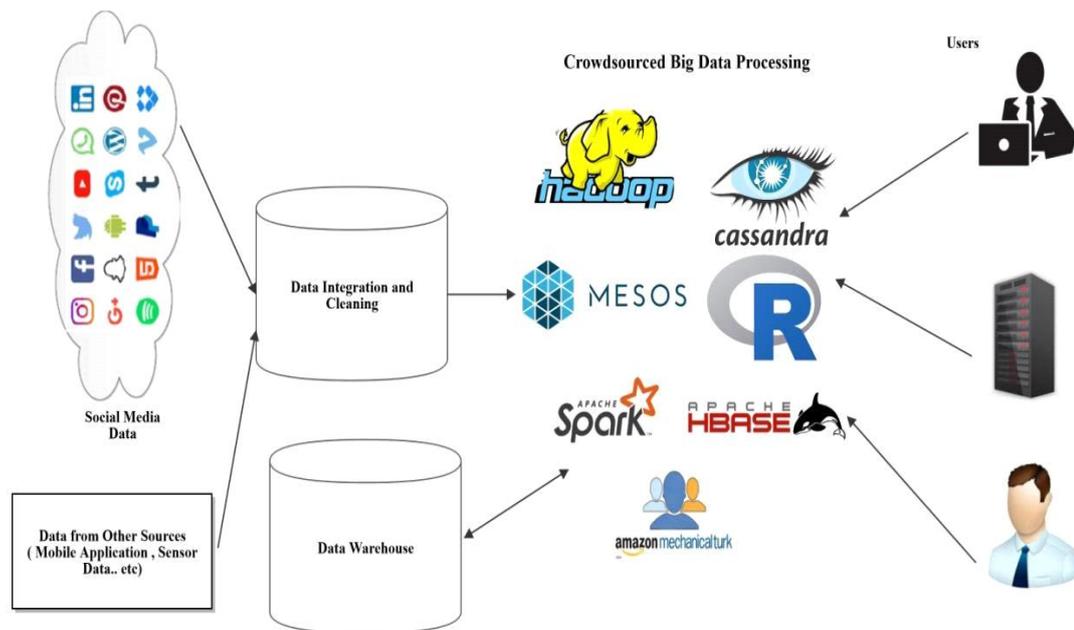
Literature Survey

This section presents the various aspects of crowdsourcing based on the applications and functional definitions. It presents various processes and methods related to crowdsourcing with issues and challenges. The initial stage of crowdsourcing research needs to understand the issues and challenges faced in crowdsourcing platform to obtain research problem. Hence, this section provides a detailed survey about crowdsourcing methods.

Bigdata Mining From crowdsourcing

Due to many advances in technology, large volume of valuable data (e.g streaming of financial, banking and shoppers market basket data) are generated at from wide varieties of structure, unstructured and semi-structure data at high velocity in a various real – life business environment, scientific application, Engineering and healthcare application in society and organizations. Due to their high volumes, the accuracy and quality of this data based on their veracity. This leads us into a new era of Bigdata. Many data mining tasks and analytics tasks can be achieved by making use of crowdsourcing. e.g. clustering, classification, MapReduce algorithm, association rule mining, Machine learning techniques and some of the algorithms have

some difficulty in handling this problem, for lack of knowledge extraction. In this situation the people volunteers that is crowds can accomplish efficiently, submissively and accurately than the prevailing algorithms. we need to solve such a kind of problems and discuss how crowdsourcing will be used to solve a problem.



Entity Resolution Model for Crowdsourcing

Due to the increasing data model complexity heuristic, meta-heuristic, machine learning and deep learning approaches used for analyzing, clustering and classifying the crowdsources. Each approach has their own style and ability in mining. Thinking about crowdsources, entity resolution (ER) is one of the important methods used for identifying a record among all the records in database. ER attains all the similar underlying records, and are therefore each other duplicates. Because of the inherent poor quality of data ER and ambiguity of data representation becomes a challenging method for automatic processes. Hence, human-powered ER (HPER) through crowdsourcing becomes a popular. The time and cost of the crowdsourcing is reduced as much as possible, by answering queries using crowd. It may provide fault answers sometime, crowd-based ER (CBER) methods are used to reduce the human interactions without affecting the quality and using a computerized similarity value. Several practical methods are performed well but theoretical analysis for crowdsources very less. Fundamental task of crowdsources is ER used for searching and identifying the records in a large size database refer to the similar underlying real-life entity, (Verykios 2007; Getoor and Machanavajjhala 2012; Christen 2012). Identifying and representing real-world entities is a complicated task. For example, user profile management in e-business, information about e-products, services and websites and data collection and management in social media websites are huge in volume to be resolute. These kinds of data have more error, missing elements, mis-matched attributes and other conflicts like redundancy. ER is the main and chief task for pre-processing the data which improves the quality of the data. Several earlier approaches have been focused on ER techniques using machine learning approaches like SVM, unsupervised learning, decision tree, ensemble classifier and conditional random fields and so on (Getoor and Machanavajjhala 2012). Still ER method is used for automated process to improve accuracy. Most of the approaches considers ER as a clustering problem. The ER method is represented mathematically as given as: number of elements (n) is clustered in disjoint parts Cluster is said to be true cluster if n are unknown to the data owners. e is the entity of the data. The data element e , which has a number of attributes A . In order to obtain the similar attributes, a similarity function is used for calculating same entities. Similarity function is applied among any two-attribute set e_1 and e_2 . The similarity function returns "0" if the entities are same, else it returns "1". Though, in experimental practices obtaining the similarity is not possible sometimes due to error in attributes. Any of the automatic processes which use similarity function including error prone task, In order to eliminate this criticality, a human knowledge based crowdsourcing method is proposed for increasing the accuracy through ER (Davidson et al. 2014; Firmani et al. 2016; Verroios and Garcia-Molina 2015; Gruenheid et al. 2015; Wang et al. 2012; Wang et al. 2013; Vesdapunt et al. 2014; Yi et al. 2012; Whang et al. 2013). The human knowledge-based ER can compare and find out the

matched and un-matched entities successfully. But it got failed in automatic strategies. Wang et al. (2013), Wang et al. (2012); Demartini et al. (2012), Whang et al. (2013) proposed Hybrid Human Interaction (HHI) approach for ER. HHI used transitive-relationship based entity selection to reduce query processing time and number of queries. Hence, it becomes a staple in each future works following or derived from HHI in Firmani et al. (2016), Verroios and Garcia-Molina (2015), Gruenheid et al. (2015), Vesdapunt et al. (2014). Wang et al. used transitive-relation for classifying the matched and non-matched entities in crowdsourcing using some trained matched and non-matched entries. Wang et al. used crowdsourcing method, were function well on real dataset. Arya Mazumdar and BarnaSaha (2016) presented a theoretical analysis about the complexities faced during query process referred from (Wang et al. 2013; Vesdapunt et al. 2014). The analysis result is measured using similarity values of the algorithms based on various constraints. It helps to understand the superiority of the algorithms and it derives the query complexity in accordance to different condition, by comparing the results among the methods by obtain the near-optimality or sub-optimality of the heuristic methods, (Mazumdar and Saha 2016). The two heuristic methods compared are edge ordering algorithm, Wang et al. (2013), and node ordering algorithm, Vesdapunt et al. (2014). All these kinds of comparisons and analysis makes more complexity regarding time and cost. Hence ER based pre-processing improves the accuracy of mining with any algorithms.

Task Allocation of Crowdsourcing

Crowdsourcing can provide a better solution for the applications like on-demand transportation, online shopping, and on-demand local delivery [10]. Crowdsourcing is attracted highly by various academic and industry people, due to increasing usage the data tremendously. For example, the information about shipment, warehouses, customers and delivery information in an unpredictable, hence Amazon uses crowdsourcing. There are three stakeholders are considered in crowdsourced delivery such as workers, customers and matching platform. The spatial information from the customers are assigned to the platform. Then platform compares and match the tasks with the workers, by analysing the availability of the workers. If the worker is free then the task will be allocated. But task allocation is possible, if and only if the spatio-temporal requirements are matched with one another.

Big Crowd Data Classification

Victor S. Sheng et al. briefly described about the frequent attainment of labels for data items when the labeling is flawed. Using the method of recurrent labeling, they identified the lack in the quality of data and concentrate particularly on the development of training labels for supervised induction. Through minimizing the cost of labeling, separating the part of data which is unlabeled could convert into significantly further limited than labeling. They provide repeated-labeling approaches of accumulative complication, hence provide major results. The authors concluded that when there exist flaws in labeling, the data miners must make use of selective attainment of multiple labels scheme in their collection. Their focus in this paper is to improve the supervised learning's data quality; however, the outcomes have inferences for data mining (Victor S. Sheng et al., 2008).

HesanSalehian et al. discussed about the problem of relating properly arranged menu item of a restaurant to a huge database consisting of less structured items of food through crowd-sourcing (HesanSalehian et al., 2017). They established an original, real-world, and scalable machine learning solution architecture, involving two main steps. Query generation approach was used which was built on a Markov Decision Processing algorithm in order to lessen the difficulty of time while looking for identical candidates. The deep learning techniques are then used to track them by a re-ranking step. The grouping of strengthening via MDP for query generation at initial, building synthetic training data using SVM, and CNN architectures for relevance learning, all originate organized to generate an influential tool for petite text matching in the absenteeism of context and/or user feedback. Across the three main variations of test sets – random partition withheld from the training set, hand-labeled data by humans, and observed responses from users – the model outperforms the basic SVM technique once the scale of training data is on the scale of 100K or more examples, and becomes really evident once the data set size reaches 1M.

Peter Welinder et al. presented a method for accessing the underlying value of individual image from comments given by numerous annotators. Their technique was created on a classical form of the annotation process and

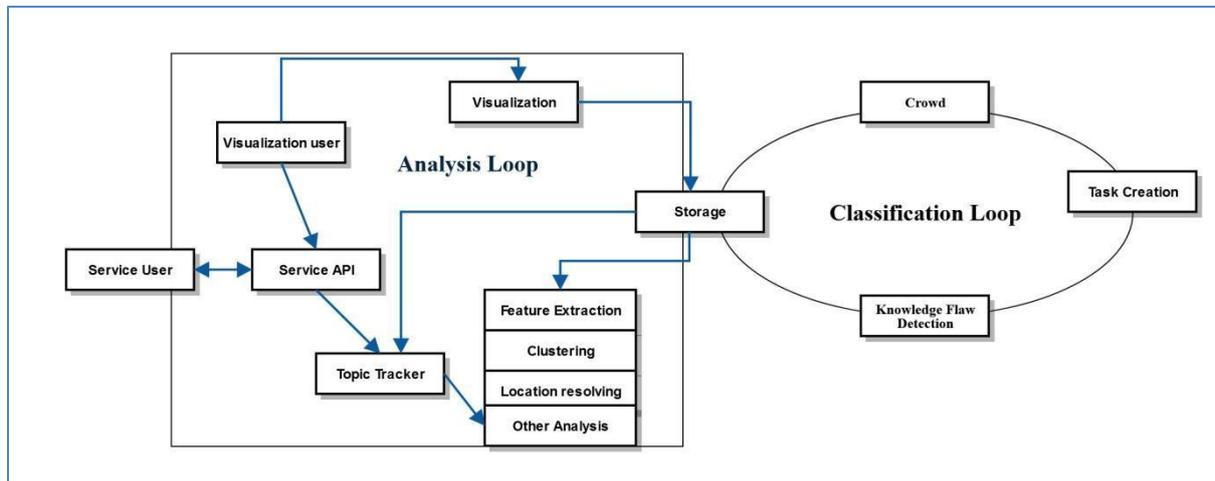
image formation (Peter Welinder et al., 2012). Each image represented in an abstract Euclidean space has different characteristics while each annotator is modeled to signify clusters of annotators that have mixed groups of services and information as a multidimensional object with variables indicating capability, knowledge and favoritism. They found out that the ground fact labels on both synthetic and real data that is guessed by the model is more precise than state of the art methods. They demonstrated that each model that start with a set of binary labels, may turn up with substantial content, such as different “schools of thought” amongst the annotators, and can group the images that are associated with separate classes. The authors gave a result about their model as it provides values for defining loss functions and for training classifiers, can accurately estimate the ground truth labels by integrating the labels provided by many annotators with totally different skills, and it will therefore higher than the present state of the art ways.

Thomas Bonald and Richard Combes deliberated the problem of precisely assessing the dependability of labors based on noisy characteristics they provide, which is a topical query in crowdsourcing. They proposed a novel lower bound on minimizing the guess estimate error which is applicable to any measured technique and an named Triangular Estimation (TE) for estimating the dependability of workers (Thomas Bonald et al., 2017). TE has low complexity, and it proved to have a minimal optimization that matches the lower bound, since it does not depend on an iterative procedure establishing in a flowing situation when labels are given by workers in actual time. Therefore, they concluded by assigning the performance of Triangular Estimation and other advanced algorithms on both artificial and real-life data.

Ofer Dekel and Ohad Shamir explained that with the repetition of search engines and crowdsourcing websites, machine learning practitioners face datasets that are labeled by a large varied set of teachers. These datasets test the limits of our current learning theory, which largely assumes that data is sampled from a fixed distribution (Ofer Dekel et al., 2009). In many cases, the number of teachers actually equals with the number of examples, with each teacher providing just a handful of labels, impeding any statistically reliable assessment of an individual teacher’s quality. In order to increase the label quality of our training set they have further elaborated the problem of clipping low-quality teachers in a crowd. Despite the obstacles mentioned above, they found that this is in fact achievable with a simple and efficient algorithm. Thus, they have provided a theoretical analysis of the algorithm and back their findings with empirical evidence.

Author Shi Zh et al. proposed an experiment in which crowd sourced data clumped task, there happens struggles in the responses given by huge collection of bases on redundant questions. The primary goal for this job is to guess cause dependency and pick replies that are given by good quality sources. Current method resolves this problem by concurrently measuring sources’ consistency and assuming queries’ factual responses. However, these procedures infer that a source has the similar dependency degree on all the queries, ignoring the detail that sources’ dependability may differ amongst various topics. To combine numerous knowledge points on diverse subjects, they proposed FaitCrowd, a fine-grained truth discovery model for the mission of gathering varying information composed from numerous bases. This method generates a fusion of query content and causes’ providing responses in a probabilistic model to guess both topical expertise and accurate responses concurrently, thus leading to a extra exact assessment of source reliability. Thus, results on 2 real-life datasets demonstrates that FaitCrowd will significantly decrease the flaw rate of accumulation associated with the progressive multi-source aggregation and demonstrated improved capability to acquire true responses for the queries associated with prevailing methodologies.

Jakob Rogstadius et al. describe the novelty approach for integrating crowdsourcing architecture into analyse the web and social media contents posted the small bloggers and service users. This given figure clearly defines the real-time working process (Jakob Rogstadius et al., 2011).



Clustering

Max/top-k and clustering problems have been studied by Susan Davidson and others for erroneous answers on comparison operations which can moreover be type or value: provided two data essentials, the response to a type assessment is “yes” if the fundamentals have the similar type and consequently fit to the same collection. To achieve accurate results with high probability, they gave well-organized algorithms that are guaranteed to analyze the cost in relation of the whole amount of judgments (i.e., using a fixed-cost model), and demonstrated that they are fundamentally the finest probable. They also came up with a prediction that less judgments are desired when values and types are associated, or when the error model is one in which the error reduces as the distance amongst the two elements in the sorted order upsurges. Since these difficulties are inspired by top-k and group-by database queries in the crowdsourced environment, where the standards used for combining and collecting are hard to assess by machineries but much easier by the crowd.

Ryan Gomes et al. proposed a feasible solution for cloud categorization. Amongst the challenges he proposed three other trials: (a) individual worker has only a limited view of the data, (b) dissimilar workers may have diverse grouping criteria and may provide dissimilar records of categories and (c) the fundamental category construction may be ordered. They used a Bayesian model to calculate how workers may tactic clustering and show how one may assume clusters / categories, as well as worker parameters, using that model (Ryan Gomes et al., 2011). Their experiments, carried out on bulky collections of images, suggest that Bayesian crowd clustering works efficiently and may be superior to single-expert annotations. Therefore, exhibiting both data entry properties and the employees’ annotation process and limits seems to provide performance that is higher to prevailing clustering aggregation methods.

Arya Mazumdar and Barna Saha initiated an exhaustive conceptual study of clustering with noisy queries. They said, even if the clusters are unknown, theoretic lower bound on the number of queries for clustering with noisy oracle in both situations and designed novel algorithms that closely match that query complexity lower bound (Arya Mazumdar and Barna Saha; 2017). Moreover, they designed computationally efficient algorithms that can be applied for both adaptive and non-adaptive settings. The problem simplifies multiple application scenarios. The crowd constitutes the noisy oracle, and the number of queries directly correlates to the weight of crowdsourcing. Furthermore, clustering with noisy oracle is closely bound with the correlation clustering, leading to improvement therein. This proposed model establishes a new course of study in the popular stochastic block model where one has a partial stochastic block model matrix to recover the clusters.

The task of clustering items using answers from non-expert crowd workers was considered by Ramya Korlakai Vinayak and Babak Hassibi in 2016. In such scenarios, the workers are often unable to name the items directly, however, it is acceptable to assume that they can differentiate items and come to a conclusion whether they are similar or not. They worked with fractional observations subject to a fixed query budget restriction since it is far too expensive to query all possible edges and/or triangles. The cost of a query by its entropy is measured first, when a generative model for the data is available; when such models do not exist the

average response time per query of the workers as a substitute for the cost is applied. In addition to conceptual reasoning, through multiple simulations and experiments on two real data sets on Amazon Mechanical Turk, they practically demonstrated that, for a fixed budget, triangle queries uniformly outclass edge queries. Even though, in contrast to edge queries, triangle queries reveal dependent edges, they provide more reliable edges and, for a fixed budget, many more of them. Compared to edge queries, it is proved that triangle queries reveal dependent ones. However, due to their error correcting abilities, triangle queries result in more consistent edges. In particular, experiment on two real datasets suggests that clustering items from random triangle queries significantly outperforms random edge queries when the total query budget is fixed.

AnttiUkkonen mentions that crowdsourcing usually depend on relative distance comparisons, as these are at ease to stimulate from human workers than absolute distance information. He overcame the obstacle in existing work by making use of correlation clustering, which is a well-known non-parametric approach to clustering. He first defined a novel variant of correlation clustering that is based on relative distance comparisons which is very much suitable for human computation. He goes on to show that his new problem is closely related to simple correlation clustering, and use this property to project an approximation algorithm for our problem and empirically compared against existing methods from literature by proposing a more practical algorithm. While conducting experiments with synthetic data, the result suggested that their approach can outperform more complex methods and also that their method efficiently finds good and intuitive clustering from real relative distance comparison data.

Fabian L. Wauthier et al. proposed an active learning algorithm for spectral clustering to remove uncertainty in a transitional clustering solution that successfully measures similarities that are mainly expected. They enlarge their algorithm to preserve running estimates of the accurate similarities, as well as estimates of their accuracy. Using this information, the algorithm updates only those estimates which are moderately inaccurate and whose update would most likely eliminate clustering uncertainty. Comparing the methods on several datasets, including a practical example where similarities are expensive and noisy, the results showed a significant enhancement in performance compared to the alternatives. They proposed an extension of the algorithm by taking the accuracies of result into account during query selection which can potentially avoid unnecessary repeat measurements and speed up the learning process in noisy settings.

Jinfeng Yi et al. combined the low-level features of objects with the manual annotations of a subset of the objects obtained via crowdsourcing by deriving a new approach for clustering which is called as semi-crowdsourced clustering. Their main idea was to learn an appropriate similarity measure, based on the low-level features of objects and from the manual explanations of only a small portion of the data to be clustered. One complexity in learning the pair wise resemblance measure is that there is a significant amount of noise and inter-worker variations in the manual annotations obtained via crowdsourcing (Jinfeng Yi et al., 2012). They addressed this difficulty by developing a metric learning algorithm based on the matrix completion method. Their empirical study with two real-world image data sets shows that the proposed algorithm outperforms state-of-the-art distance metric learning algorithms in both clustering accuracy and computational efficiency.

He Jiang et al. describes that the Fuzzy Clustering Test Reports (FULTER) problem which makes the test results and their root causes complex to diagnose. In order to resolve FULTER, sequences of barriers need to be conquered. Thus, they proposed a new framework named Test Report Fuzzy Clustering Framework (TERFUR) by aggregating excess and multi-bug test reports into clusters to decrease the number of tested reports of test. The efficiency of TERFUR is validated in prioritizing test reports for manual inspection. The experimental results show that TERFUR can cluster redundant test reports with high accuracy and significantly outperform comparative methods. In addition, trial results also reveal that TERFUR can greatly reduce the cost of test report inspection in prioritizing test reports Jiang et al., 2018).

Patterns Mining

Based on the workers response the pattern mining observes the significant patterns. In pattern mining process discovering the significant pattern is the challenging task. For example, a health researcher tries to discover the rules of association by analyzing its performance of medicine in traditional manner and she discovers that "Garlic can be used to treat flu". But in this case, she cannot use the database since it understands only treatment and symptoms for a specific disease and the list of diagnosed cases from the healers. On the

survey we cannot get all transaction list, only summary can be provided. For example, they may have a perception that “Once I have flu, most of the time I will take Garlic because it indeed is useful to me”. Given the summaries of individuals or the personal rules answered by diverse people, they can be summed up together to find a complete important rule (or the general trends). So, the crowd pattern mining intends to collect the personal rules from crowd workers, cumulate them and find the complete set of rules (i.e., general trends). Crowd pattern mining typically involves generation of a large amount of pattern that occur frequently without actually providing information which is needed for interpreting the pattern. It also provide semantic annotation for the frequently occurring pattern, it will help to understand patterns in a better manner.

Yukino Baba and Hisashi Kashima from the University of Tokyo came up with a solution to overcome the challenge of quality control in crowdsourcing. The prevalent and existing models that overcome this issue are by introducing redundancy, where a number of workers vote their responses. But this solution is hypothetical in case of unstructured response formats. Yukino Baba and Hisashi Kashima have proposed an unsupervised statistical approach for unstructured responses. This approach involves two stages namely the creation stage which involves the unstructured responses of the crowd workers and the review stage involves voting via the multiple-choice questions. It is proved to deliver high quality crowdsourcing with low cost.

AmnaBasharat, I. BudakArpinar and KhaledRasheed of the University of Georgia gave a unique illustration on how to leverage crowdsourcing to create workflows by thematic annotation of the special case of the widely read manuscript Qur’an. The Qur’an is rich in morphology and semantics. Hence it involves knowledge intensive and specialized domains. This model involves several stages such as ontology design, Task generation and design. The task is entered in the Amazon Mechanical Turk. This proposal involves sub-verse level annotation along with explicit and implicit assertions. The result of the crowdsourcing is promising with 96% approved for explicit assertions and 81% approved for implicit assertions. This framework can be generalized to other knowledge and domains.

Outlier Detection

The major goal of the Crowd-powered Outlier Detection is to sense outliers from crowd answers. For the Quality of Experience (QoE) assessment in done by outlier detection, which deals with the user’s subjective expectation, perception, satisfaction and feeling with respect to the content in multimedia. Workers are requested to specify an evaluation score ranging from “Bad”, “Medium” to “Excellent” in order to grade the quality of a multimedia. But the crowd may generate noise during such enquiry. Research has been done to assess the Quality of Experience (QoE) for Outlier detection.

QianqianXu et al. (2017) have proposed a simple iterative algorithm using non-convex optimization principle to evaluate the QoE for outlier detection. They have come up with two approaches (1) for a known outlier sparsity size, they proposed iHT and iLTS methods which provides the same performance as LASSO and ninety times faster computational speed than LASSO (2) for an unknown outlier sparsity size, they propose aLTS which is an adaptive method that can used to estimate the number of outliers without any prior dataset and is proved to be nearlythree-eight times faster than LASSO. They have shown the proposed method effectiveness with the help of data which is simulated on known ground-truth outliers, followed by a real-world crowd sourcing dataset without ground-truth outliers. Thus, they have proposed an approach for the people in multimedia community to exploit crowd source paired comparison data for robust ranking.

HongleiZhuang et al. (2015) have proposed a specific type of explanation bias in crowdsourcing where the data submitted for the workers can be judged simultaneously only through grouping them into batches. They have come up with a model to characterize the annotating behavior on data batches and also train the worker model based on the annotated datasets. They have formalized a method for de-biasing crowdsourced batches to eliminate the effect of annotation bias from unfavorably affecting the accuracy of labels. Their experimental results reflect on both synthetic data and real-world data which demonstrate the effectiveness of their proposed method.

Crowd Sourced Machine Learning

Chong Sun et al. described about their solution to classify millions of products into thousands of product types at Walmart Labs using Chimera, which employs a combination of learning, rules and crowdsourcing to achieve accurate, continuously improving, and cost-effective classification. The authors claim that bulky products in crowdsourcing is critical but must be used in combination with learning, rules, and in-house analysts. They also argue that usage of protocols is essential, and that more attention in research should be paid to helping analysts create and manage large quantity of rules more effectively. They also state that it is significant to explore further hybrid human-machine systems such as Chimera, which have proven successful in solving certain classes of real-world Big Data problems. Their key message is that all of the components such as learning, rules, crowdsourcing, analysts, and developers are major for large-scale classification.

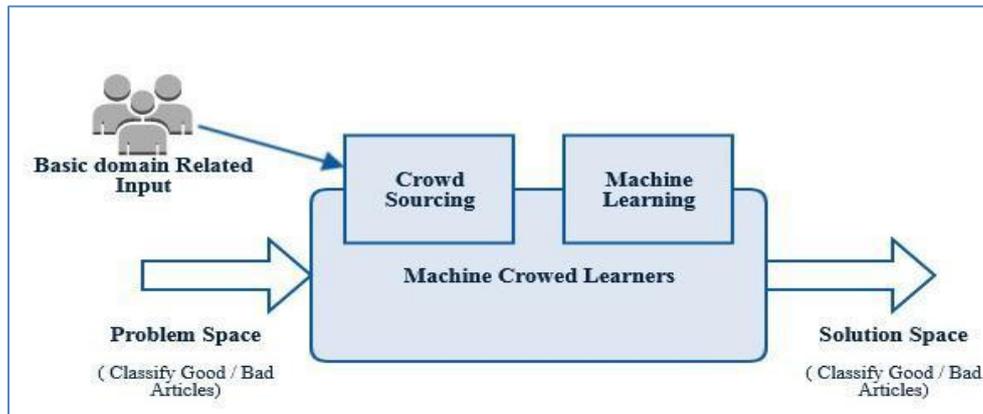
Matthew Lease mentions about the two particular aspects of crowdsourcing- data quality control (QC) and ML. He states that the advent of crowdsourcing has created its own opportunities for improving over the traditional methods of data collection and annotation, which paved the path for the arrival of data-driven machine learning. Crowd-based human computation has supplemented the automated machine learning (ML). The author compares and analyses the advent of automation over the crowd-based work. He questions about the sustainability of the crowd labor over the growing demand of applications which requires human expertise or with privacy, security, or intellectual property. The author is concerned about leveraging the human interactive computation over the crowd manual work. He concludes that the Computational wisdom of crowds (WoC) and collaborative thinking may help to understand better about how to mine and aggregate human wisdom while learning active theories which might provide deep insights and focused learning.

Steven Burrows et al. proposed the Web is Crowd Paraphrase Corpus 2011 (Webis-CPC-11) for paraphrasing and plagiarism detection by focusing on two aspects of paraphrase acquisition via crowdsourcing and passage-level samples. Since crowdsourcing paradigm is not effective without the objective of quality assurance, the creation of text corpus is unacceptably expensive. The second facet states the discrepancy that the majority of the previous add generating and evaluating paraphrases has been conducted exploitation sentence-level paraphrases or shorter. They show that the financial outcome is cost neutral. Machine learning experiments that discovers if passage-level paraphrases majorly contributes to identify a two-class classification problem using paraphrase similarity features. They concluded that the difference between paraphrased and no paraphrased samples can be correctly distinguish proposed method using k-nearest-neighbor classifier.

Justin Cheng and Michael S. Bernstein from Stanford University introduced Flock, an end-user machine learning platform that uses paid crowdsourcing to speed up the prototyping iteration and extend the performance of machine learning systems. The system allows the users to enhance as hybrid crowd-machine learners by performing three methods such as structuring a nomination process, grouping the suggested features and collecting labels on these new features. It loops and gathers more crowd features to improve performance level on subsets of the space where the model is misclassifying many examples. If a decision tree is considered that uses machine-readable features, Flock can dynamically grow sub-trees from nodes that have high classification error, or even replace whole branches. Moreover, these constraints can help focus the crowd to generate more informative features. The authors demonstrated Flock's effectiveness through an evaluation of six broadly various prediction tasks, including discerning videos of people telling the truth or lying and differentiating between paintings by impressionist artists. They found that aggregating crowd characteristics is more accurate or measureable than asking for direct opinions from the crowd. They conclude that hybrid crowd-machine learning systems may offer a route to rapid prototyping and consideration of the feature space, even in traditionally quite complex domains. Flock predicts that in future end user could create a machine learning system just by explaining the prediction goal into a free-form textbox.

EceKamar et al. demonstrated how machine learning and inference can be coupled to uplift the harmonizing human strengths and autonomous agents in a group to solve crowdsourcing tasks by constructing a set of Bayesian predictive structures from data. Moreover, an overall schematic crowdsourcing model that combines the efforts of people and machine vision on the task of segregating celestial bodies defined within a specific citizens' science project named Galaxy Zoo is explained. They found how to learn probabilistic models that can be used to combine human and machine contributions and predict the behavior of workers. A system is created that combines machine vision, machine learning, and decision-theoretic planning to make effective decisions

about when to hire workers and how to perform classifications. They labored a collection of inferences collectively to guide choices on hiring and routing staff to tasks. A prototype system was constructed and assessed the learning and decision-making techniques on real-world data collected during the operation of the Galaxy Zoo system. The experiments demonstrate that the methodology can solve consensus tasks accurately and achieve significant savings in worker resources by intelligently allocating resources.



Various Tools for Crowdsourcing

Type	Tool	Kind of Problem	Working
Distributed human intelligence tasking	Amazon Mechanical Turk	Large-scale data analysis where intelligence of human is more efficient.	Large amount of information is analyzed by a huge crowd through organizing the tasks.
Knowledge discovery and Management	Ushahidi	Crowdsourcing Mapping and Crowd feeding Tool	SMS, Web Submission, E-mail, Facebook, Twitter and Voice Mail.
Artificial Intelligence	CROWDFLOWER	Sentiment analysis and AI based problems	Collecting live internet data.
Knowledge discovery and Management	SeeClickFix	Ideal for information gathering. Creation of collective resources	Finding, gathering or collecting of crowd into a mutual location and format.
Broadcast Search	INNOCENTIVE	Design or aesthetic problems.	Empirical problems are solved by the crowd.
Peer vetted creative production	Threadless	Scientific problem solving	Selective and creative ideas of the crowd.

Heuristic Methods based Crowdsourcing

Several algorithms have been used for computing processes over social aware data. This section says the name of the algorithms used for obtaining the region of mobile crowdsourcing. Most of the algorithms are executed iteratively based on cyclic models like query processing, Q and A, and feedback outputs using heuristic algorithms. These kinds of methods always adjust the query based on the customer's repeated input, where it takes more computational time and comparison. There are three different heuristics are considered here are incorporated with crowdsourcing, minimizing the computational costs and improving the efficiency in crowdsourcing applications. This process involves while using crowdsourcing is,

- (i) On expectation of failures in each round, additional questions are asked.
- (ii) Neighborhood associations are used in the circumstance where clusters are created based on regions of interest.
- (iii) Using spatial point processes to model the region of interest.

But heuristics can improve the performance using a stylish model like clustering operations in the data. It is well known that crowdsourcing is one of the influential systems for problem answering. It is also used for gathering critical information, binding the power of crowd and it combines human & machine computations, discussed in Brabham DC (2013), Law and Ahn (2011), Michelucci P (2013). Crowdsourcing is used in certain situations where a task which cannot be done only by a machine or human in better manner, for example CrowdDB, Franklin et al. (2011). Crowdsourcing to mobile users is named as "Mobile CrowdSourcing (MCS)". MCS offer more opportunities, issues and challenges for human calculations including tasks with spatial and time-based properties are discussed in Alt.F et al. (2010), Georgios et al. (2012), Gupta et al. (2012), Charoy et al. (2013), Kazemi and Shahabi (2012), Della et al. (2013), Yan et al. (2009). Algorithms help human to process, view and store the data processing, which can be explored for finding the maximum data, Guo et al. (2012), filtering a data is discussed in Parameswaran et al. (2012), searching a subset of data from the whole unstructured dataset is given in Das Sarma et al. (2014), and finally optimizing the cost, time and computational efficiency.

Classification Algorithms

The very first data classification method which is using workers participated in the medical context and all the patients are labelled with the clinicians, works and workers. Dawid and Skene, (1979) introduced an algorithm known as Expectation – Maximization (EM) for improving the accuracy of estimating each calculation in a problem. Different updates of the algorithms are proposed and experimented with different value settings, is discussed by, Hui and Walter (1980), Smyth et al., (1995), Albert and Dodd, (2004), Raykar et al., (2010), Liu et al., (2012).

Bayesian techniques are projected in a wide range and has been implemented by [Raykar et al., 2010, Welinder and Perona, 2010, Karger et al., 2011, Liu et al., 2012, Karger et al., 2014, 2013] and references therein. The belief-propagation (BP) is on a particular interest which was proposed by [Karger et al., 2011]. The algorithm is order-optimal in relation to the total number of workers needed to perform a task given for a particular target error rate, considering the boundary of an infinite number of tasks and workers. The other algorithm family deals with the matrix's spectral analysis which represents the associations between workers and tasks. [Ghosh et al., 2011] proposed the task-task matrix where the entries denote the quantity of workers who labeled two tasks in the similar way, while [Dalvi et al., 2013] proposed the worker-worker matrix where the entries denote the amount of tasks labeled in similar way by two workers. The two works obtain guarantee in their performance through perturbation analysis of the topmost eigenvector of the equivalent anticipated matrix. The BP algorithm of Karger, Oh and Shah is thoroughly associated to these spectral algorithms: message-passing scheme is comparable to the power-iteration method functional to the task-worker matrix, as observed in [Karger et al., 2011].

Two recent contributions that are notable are [Chao and Dengyong, 2015] and [Zhang et al., 2014]. The former delivers performance assurances for two forms of EM, and originates lower bounds on the

possible estimation error (the probability of assessing labels erroneously). The latter offers lower limits on the approximation error of the employees' dependability as well as performance assurances for a better type of EM trusting on spectral approaches in the initialization stage. Our lower bound cannot be associated to that of [Chao and Dengyong, 2015] because it applies to the workers' dependability and not the forecast error; and our lower bound is snuggier than that of [Zhang et al., 2014]. Our estimator shares some structures of the algorithm proposed by [Zhang et al., 2014] to reset EM, which proposes that the EM phase itself is not vital to attain minimax optimality. All these algorithms necessitate the storing of all labels in recollection and, to the best of our knowledge, the only known streaming algorithm is the recursive EM algorithm of [Wang et al., 2013], for which no performance guarantees are available.

The following table summarises the set of all crowdsourcing methods, merits and the limitations of the methods faced in the earlier research works. It helps to understand the overall issues and challenges and make us to decide about the research problem, which can be solved efficiently.

Author, year	Method proposed	Merits	Limitations
Victor S. Sheng et al., 2008	Examining the improvement in data quality through repeated Labeling and training labels.	Improves the labeled data quality and model learned from data in a wide range of conditions.	Different qualities of labelers were not assumed for analysis.
Peter Welinder et al., 2010	Bayesian generative probabilistic model used to model each annotator as multidimensional entity for estimating the underlying value of each image from multiple annotations.	Identifies different sets of skills and knowledge of annotators and parameters of the image.	Expertise annotators required.
Ryan Gomes et al., 2011	Proposed a solution on how workers may approach and infer clusters using Bayesian model.	May be superior to single-expert annotations.	Implementation is not carried out.
Matthew Lease 2011	Two features of crowdsourcing and relationship between them are considered. ML and Data Quality Control (QC)	Provides wide exposure since Machine Learning is applied in crowdsourcing.	Only the basics are covered. Detailed knowledge about future enhancement is not explained.
Fabian L. Wauthier et al., 2012	Proposed a dynamic algorithm for spectral clustering which totally measures the resemblances which are probable to eliminate the indecision in an midway clustering solution.	The method was used on several datasets. The similarities were noisy and expensive in a realistic example. The performance showed a huge improvement compared to the alternatives.	The pairwise similarities are expected to be identified originally. But supplementary labels or restrictions can be enquired.

Jinfeng Yi et al., 2012	A new approach called semi-crowd sourced clustering was proposed that combines the structures of objects with the labor-intensiveness of a sub set of the objects attained via crowdsourcing.	Both clustering accuracy and computational efficiency were outperformed using state-of-art distance learning algorithms.	Restricted to a linear similarity function.
Severin Hacker et al., 2012	A set of Bayesian predictive models were constructed and described their operation within crowdsourcing architecture.	The efficiency of large-scale crowdsourcing procedures are exploited built on predictable utility.	Challenges of this proposed system is not considered.
Steven Burrows et al., 2012	Subsidizes to paraphrase acquisition and emphasizes on two features that are not lectured by present study: (1) attainment via crowdsourcing, and (2) attainment of passage-level examples.	K-nearest neighbour classifier can appropriately differentiate between rephrased and non-rephrased samples	Quality of the paraphrases matters.
Yukino Baba et al., 2013	Unsupervised numerical quality assessment method for overall crowdsourcing responsibilities with unstructured reply arrangements which inhabit the widely held crowdsourcing marketplaces.	Proposed model achieved significantly advanced performance than further methods and could bring high-quality outcomes with subordinate prices in crowdsourcing.	In terms of correlation measure, the language translation task performed the best.
Susan Davidson et al., 2014	Showed that comparisons are needed when values are associated, or when the error model is one in which the error decreases as the distance between the two elements in the sorted order increases.	Provides a formal basis for max/top-k and clustering queries in crowdsourcing applications, that is, when the oracle is implemented using the crowd.	Consequences from only the source for learning additional compound models appropriate for actual crowd obtained submissions by providing lower bounds as well as algorithmic ideas.
Fenglong Ma et al., 2015	Conflicting data are collected from multisource and aggregated using Fain Crowd, a fine grained	The error rate of collection is reduced associated with the state-of-the-art multi-source combination owing to its	Considers both questions topic and fine-grained user-Expertise at once. This is not suitable for

	truth discovery model.	capability of learning skill by modelling questions and responses.	other Dataset.
Vladimir Stantchev et al., 2015	Artificial Immune System and abstract global knowledge representation model based on ontology is used to provide a novel way of taking advantage of information from user's social network and to recommend users.	Provides a method for optimal matching supply-demand which is the current augment collaborative learning environments	Vary the parameters of the algorithm to provide better results. Perform a deep analysis of cutting-edge ontology technologies to determine user characteristics and knowledge.
Justin Cheng et al., 2015	Flock's crowd feature generation was evaluated to recognize the usefulness of human-generated features and structures generated by hybrid human-machine systems.	The process of combining crowd-nominated structures outpaces estimates from the crowd and engineered structures.	Machine-only classifier's performance suffers when class involvement is possibly unpredictable.
Maryam M Najafabadi et al., 2015	Proposes a system for utilizing some important problems in Big Data Analytics and Deep Learning challenges introduced by Big Data Analytics.	Provides a solution to learning and analysis problems in huge volumes of input data	Works are required for adapting deep learning algorithms with big data.
Honglei Zhuang et al., 2015	A novel worker model to illustrate and train the explaining performance on data groups. Debiasing technique is applied to eliminate the result of annotation bias from harmfully disturbing the correctness of labels obtained.	The debiasing strategy can activate a wide variety of claims.	Incorporate the different behavior of each individual worker and adjust the debiasing strategy accordingly.
Shayan Doroudi et al., 2016	Training novice workers to achieve fine on answering errands in circumstances where the space of approaches is huge and laborers want to be positive through checking possibility and exercise workers in the nonappearance of domain skill.	Workers in the filtered medium-long set have a much-advanced average per task accuracy.	Did not find a universal mechanism to find the structures of high-quality authentication tasks in any area

Amna Basharat et al., 2016	Influence crowdsourcing to make workflows for knowledge engineering in particular and knowledge rigorous domains.	96% errands were accepted for explicit assertions on submission, 81% were approved on implicit ones on without validation.	A substantial amount of tasks remained incomplete for the implicit annotations indicated lower crowd engagement, due to that fact that the job declarations were obtainable only in the Arabic language.
Mohammad Abu Alsheikh et al., 2016	Using deep learning in MBD, a context-aware action recognition application is proposed which analyzes and deliberates scalable learning framework over Apache Spark	Deep models avoid overfitting which is superior to the shallow context learning models which is the existing model. Shows speedup efficiency.	Collecting and labelling the MBD data is a challenging task.
Ramya Korlakai Vinayak and Babak Hassibi 2016	Two types of queries: random edge queries, where item pairs are revealed, and random triangles, where a triple is compared for partial observations.	Provides a sufficient state for the recovery of the adjacency matrix need for minimum cluster size based on the number of observations, edge densities outside and inside clusters.	Structure of the triangle query have not yet exploited.
Arya Mazumdar and Barna Saha 2017	Initiated a rigorous theoretical study of clustering with noisy queries (or a faulty oracle) to recover the true clustering by asking minimum number of pairwise queries to an oracle.	A new direction of study is introduced in the popular stochastic block model, where the clusters are recovered through an incomplete stochastic block model matrix.	In some cases, adaptive querying is difficult to apply.
Antti Ukkonen 2017	Proposed a work on correlation clustering which is a non-parametric approach to clustering.	More complex methods can be outperformed by making use of this approach.	Only the initial limitations have been taken.
Thomas Bonald et al., 2017	A novel lower bound and Triangular Estimation on the minimax estimation error which applies for estimating the reliability of workers	One can obtain both minimax optimality and better numerical performance by forgoing the EM phase altogether in the case of binary labels.	A large spectral gap of worker-worker matrix.
Hesam Salehian et al., 2017	Markov Decision Process algorithm, using deep learning techniques	A powerful tool for text matching is created by the combination of	The food data limitation can make it difficult for straightforward

	was proposed to decrease the complexity of time by searching for similar candidates trailed by a re-ranking step.	strengthening via MDP for query generation, SVM for training data building, and CNN architectures for learning relevance, in the absence of context or user feedback.	techniques.
QianqianXu et al.,2017	Some modest and fast algorithms were proposed for outlier detection and evaluation of robust QoEbased on the nonconvex optimization principle.	Algorithms produced with approximately 8- or 90-times speed-up, without or with a prior knowledge on the sparsity size of outliers which is similar in performance to the robust ranking using Huber-LASSO approach.	The dependability of one rater is built on comparing with the other raters.
He jiang et al., 2018	Proposed a novel framework named Test Report Fuzzy Clustering Framework (TERFUR) which reduces the quantity of inspected test reports by collecting jobless and multi-bug test reports in the form of clusters.	TERFUR provides the test report of clustering by up to 78.1% and outperforms other comparative methods.	Framework will be deployed in the future.

Various Applications in Crowd Sourcing

Min Chen et al. (2018) proposed “Urban Healthcare Big Data System based on Crowdsourced and Cloud-based Air Quality Indicators” which integrates the air quality data from numerous bases, in command to prepare data for the artificial intelligence based smart urban services. The increasing process of globalization urges half of the population to live in the cities which bring to bear a major influence in the air quality which eventually affects the health status of people. In this paper, an UH-BigDataSys is proposed which is an urban healthcare big data system that is used along with the data composed over meteorological sites, IoT sensing with user’s body signals and mobile crowdsourcing. The data for artificial intelligence based smart urban services is prepared by making use of the technique of mixing numerous source air quality data. A testbed of UH-BigDataSys is set up with the execution of healthcare applications which focus on air quality-awareness. The guidance to health by considering the quality of air is given to the people for better living.

Shixia Liu et al. (2018) introduced an “Interactive Method to Improve Crowdsourced Annotations” - a well organised method for authenticating and refining crowdsourced annotations. The outcomes of the supervised and semi-supervised learning show the training data quality to be a censorious factor. Due to the expensiveness of labelling large datasets, researchers have found it inefficient in terms of quality. An interactive method is hence proposed in order to assist specialists in validating indeterminate instance labels and untrustworthy workers.

Xiangjie Kong et al. (2018) suggested “Shared Subway Shuttle Bus Route Planning Based on Transport Data Analytics”. A two-stage approach is used to provide with an easy transportation mode for overcrowding urban traffic, which basically involves the dynamic route planning and travel requirement prediction, based on various

bus data shared by the crowd to produce instant ways for communal buses in the “last mile” act. The characteristics of the travelers and the buses are analysed and an algorithm (prediction) for dynamic routes are proposed. The results of the prediction is combined with the station properties and the user is updated with the optimal routes.

Mark Birkin (2018) initiated “Spatial data analytics of mobility with consumer data” which deals with heightening of bias selection which harms the eminence of many customer data bases. The information which is emerged through the interaction between the service provider and consumer is becoming omnipresent. These frequently collected and quickly released sets of data are more often used for the research. These data cover a wide variety of characteristics such as lifestyle, attitude and other behavioral features which are often dynamically restocked. Hence these data are analysed to provide the connection between consumer data and spatial data. New patterns such as spatial variation in channel preferences for are revealed for customer obtaining by investigating the flexibility designs and procedures in the marketing and leisure sectors.

Mohammad Masudur Rahman et al. (2018) proposed “Effective Reformulation of Query for Code Search using Crowdsourced Knowledge and Extra-Large Data Analytics” which addresses query reformulation targeting code search. Software developers often use the natural language for searching code snippets in the search engines. The results are not efficient because the quality of the query. This paper focuses on a method that uniquely recognizes specific and relevant API classes. The border count is used to rank and collect the top term weighting algorithms which makes use of pseudo-relevance feedback of candidate API classes from Stack overflow. Thus the relevance is identified and the results are presented to the developer.

Matthew Brehmer et al. (2018) brought forward “Visualizing Ranges Time on Mobile Phones: A Task-Based Crowdsourced Evaluation” uses a Linear or Radial design of range marks to compare the performance of the participants. It also identifies the set of drawbacks in terms of what range could viably be showed on a minor screen.

Yoonjung Kim et al. (2019) came up with “Quantifying nature-based tourism in protected areas in developing countries by using social big data”. The system identifies ‘where and when the people visit,’ to estimate the structures of Nature based tourism. The goal of this system is to classify the potentiality of the tourism area through geographical spatial data. Geographical data from flickr, i.e, from www.flickr.com, which is used as the key source for social big data. Hence the protected area management is evaluated based on the people interest to visit that particular place and in the same time improving ‘eco-tourism’.

Conclusion

The major objective of this work is to carry out a detailed study on various algorithms, methods and techniques used for crowdsourcing. Crowdsourcing is used huge amount of dynamic data where human participated, continuously changing and increasing regarding volume, variety and value. Crowdsourcing is used, applied and deployed in various computing industries with people involvement like academic research, hotel, medical, healthcare and environmental industries for managing, clustering, classifying and predicting a data required by a user dynamically. From the above discussion, it is clearly identified that crowdsourcing is mainly used in large size dataset changing its nature dynamically. From individual distance mapping into machine learning algorithms were highly used for crowdsourcing processes. Whereas, still the efficiency needs to be improved in terms of classification and prediction. Hence, this survey has given a suggestion that, deep learning based approaches can improve the accuracy in crowdsourcing.

References

- [1]. Mazumdar and B. Saha. A Theoretical Analysis of First Heuristics of Crowdsourced Entity Resolution. The Thirty-First AAAI Conference on Artificial Intelligence (AAAI-17), 2017.
- [2]. Elmagarmid, A. K.; Ipeirotis, P. G.; and Verykios, V. S. 2007. Duplicate record detection: A survey. IEEE Trans. Knowl. Data Eng. 19(1):1–16.

- [3]. Getoor, L., and Machanavajjhala, A. 2012. Entity resolution: theory, practice & open challenges. *PVLDB* 5(12):2018–2019.
- [4]. Christen, P. 2012. *Data matching: concepts and techniques for record linkage, entity resolution, and duplicate detection*. Springer Science & Business Media.
- [5]. Getoor, L., and Machanavajjhala, A. 2012. Entity resolution: theory, practice & open challenges. *PVLDB* 5(12):2018–2019.
- [6]. Davidson, S. B.; Khanna, S.; Milo, T.; and Roy, S. 2014. Top-k and clustering with noisy comparisons. *ACM Trans. Database Syst.* 39(4): 35 :1–35:39.
- [7]. Firmani, D.; Saha, B.; and Srivastava, D. 2016. Online entity resolution using an oracle. *PVLDB* 9(5):384–395.
- [8]. Verroios, V., and Garcia-Molina, H. 2015. Entity resolution with crowd errors. In 31st IEEE International Conference on Data Engineering, ICDE 2015, Seoul, South Korea, April 13-17, 2015, 219–230.
- [9]. Gruenheid, A.; Nushi, B.; f, T.; Gatterbauer, W.; and Kossmann, D. 2015. Fault-tolerant entity resolution with the crowd. *CoRR* abs/1512.00537.
- [10]. Wang, J.; Kraska, T.; Franklin, M. J.; and Feng, J. 2012. Crowder: Crowdsourcing entity resolution. *PVLDB* 5(11):1483–1494.
- [11]. Wang, J.; Li, G.; Kraska, T.; Franklin, M. J.; and Feng, J. 2013. Leveraging transitive relations for crowdsourced joins. In *SIGMOD Conference*, 229–240.
- [12]. Vesdapunt, N.; Bellare, K.; and Dalvi, N. 2014. Crowdsourcing algorithms for entity resolution. *PVLDB* 7(12):1071–1082.
- [13]. Yi, J.; Jin, R.; Jain, A. K.; Jain, S.; and Yang, T. 2012. Semi-crowdsourced clustering: Generalizing crowd labeling by robust distance metric learning. In *NIPS 2012.*, 1781–1789.
- [14]. Whang, S. E.; Lofgren, P.; and Garcia-Molina, H. 2013. Question selection for crowd entity resolution. *PVLDB* 6(6):349–360.
- [15]. Demartini, G.; Difallah, D. E.; and Cudré-Mauroux, P. 2012. Zencrowd: leveraging probabilistic reasoning and crowdsourcing techniques for large-scale entity linking. In *WWW*, 469–478.
- [16]. Howe, J. (2006). The rise of crowdsourcing. *Wired magazine*, Vol. 14, No. 6, pp.1-4.
- [17]. Faridani, S., Lee, B., Glasscock, S., Rappole, J., Song, D., and Goldberg, K. (2009). A networked telerobotic observatory for collaborative remote observation of avian activity and range change, *International Federation of Automatic Control*. Elsevier.
- [18]. Von Ahn, L. (2006). Games with a purpose. *Computer*, Vol. 39, No. 6, pp.92-94.
- [19]. Brabham DC (2013) *Crowdsourcing*. The MIT Press, Cambridge.
- [20]. Law E, Ahn LV (2011) *Human computation*. Synthesis lectures on artificial intelligence and machine learning. Morgan & Claypool Publishers, San Rafael.
- [21]. Michelucci P (2013) *Handbook of Human Computation*. Springer Publishing Company, Incorporated, New York.
- [22]. Franklin MJ, Kossmann D, Kraska T, Ramesh S, Xin R (2011) CrowdDB: answering queries with crowdsourcing. In: *Proc. of the 2011 ACM SIGMOD International Conference on Management of Data, SIGMOD '11*, ACM, New York, pp 61–72.
- [23]. Alt F, Shirazi AS, Schmidt A, Kramer U, Nawaz Z (2010) Location-based crowdsourcing: extending crowdsourcing to the real world. In *Proceedings of the 6th Nordic Conference on Human-Computer Interaction: extending boundaries, NordiCHI '10*, New York, ACM, pp 13–22.
- [24]. Georgios G, Konstantinidis A, Christos L, Zeinalipour-Yazti D (2012) Crowdsourcing with smartphones. *IEEE Internet Computing* 16(5):36–44.
- [25]. Gupta A, Thies W, Cutrell E, Balakrishnan R (2012) mClerk: enabling mobile crowdsourcing in developing regions. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '12*, ACM, New York, pp 1843–1852.
- [26]. Charoy F, Benouare K, Valliyur-Ramalingam R (2013) Answering complex location-based queries with crowdsourcing. In: *Proc. of the 9th IEEE Int. Conf. on Collaborative Computing: Netw., App. and Worksharing*, IEEE Computer Society.
- [27]. Kazemi L, Shahabi C (2012) Geocrowd: enabling query answering with spatial crowdsourcing. In: *Proceedings of the 20th International Conference on Advances in Geographic Information Systems, SIGSPATIAL '12*, ACM, New York, pp 189–198.

- [28]. Della Mea V, Maddalena E, Mizzaro S (2013) Crowdsourcing to mobile users: a study of the role of platforms and tasks. In DBCrowd, pp 14–19.
- [29]. Tamilin A, Carreras L, Ssebagala E, Opira A, Conci N (2012) Context-aware mobile crowdsourcing. In: Proceedings of the 2012 ACM Conference on Ubiquitous Computing, UbiComp '12, ACM, New York, pp 717–720.
- [30]. Yan T, Marzilli M, Holmes R, Ganesan D, Corner M (2009) mCrowd: a Platform for Mobile Crowdsourcing. In Proc. Of the 7th ACM Conference on Embedded Networked Sensor Systems, SenSys '09, ACM, New York, pp 347–348.
- [31]. Guo S, Parameswaran A, Garcia-Molina H (2012) So who won?: dynamic max discovery with the crowd. In: SIGMOD Conference, pp 385–396.
- [32]. Parameswaran AG, Garcia-Molina H, Park H, Polyzotis N, Ramesh A, Widom J (2012) Crowdscreen: algorithms for filtering data with humans. In SIGMOD Conference, pp 361–372.
- [33]. Das Sarma A, Parameswaran A, Garcia-Molina H, Halevy A (2014) Crowd-powered find algorithms. In: Data engineering (ICDE), 2014 IEEE 30th International Conference on, pp 964–975.
- [34]. Paul S Albert and Lori E Dodd. A cautionary note on the robustness of latent class models for estimating diagnostic error without a gold standard. *Biometrics*, 60(2):427–435, 2004.
- [35]. Gao Chao and Zhou Dengyong. Minimax optimal convergence rates for estimating ground truth from crowdsourced labels. Tech Report <http://arxiv.org/abs/1310.5764>, 2015.
- [36]. Nilesh Dalvi, Anirban Dasgupta, Ravi Kumar, and Vibhor Rastogi. Aggregating crowdsourced binary ratings. In Proc. of WWW, 2013.
- [37]. A. P. Dawid and A. M. Skene. Maximum likelihood estimation of observer error-rates using the EM algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):pp. 20–28, 1979.
- [38]. Arpita Ghosh, Satyen Kale, and R. Preston McAfee. Who moderates the moderators?: crowdsourcing abuse detection in user-generated content. In Proc. of ACM EC, 2011.
- [39]. Sui L Hui and Steven D Walter. Estimating the error rates of diagnostic tests. *Biometrics*, pages 167–171, 1980.
- [40]. David R. Karger, Sewoong Oh, and Devavrat Shah. Iterative learning for reliable crowdsourcing systems. In Proc. of NIPS, 2011.
- [41]. David R Karger, Sewoong Oh, and Devavrat Shah. Efficient crowdsourcing for multi-class labelling. *ACM SIGMETRICS Performance Evaluation Review*, 41(1):81–92, 2013.
- [42]. David R Karger, Sewoong Oh, and Devavrat Shah. Budget-optimal task allocation for reliable crowdsourcing systems. *Operations Research*, 62(1):1–24, 2014.
- [43]. Qiang Liu, JianPeng, and Alex T Ihler. Variational inference for crowdsourcing. In Proc. of NIPS, 2012.
- [44]. Shmuel Nitzan and Jacob Paroush. Optimal decision rules in uncertain dichotomous choice situations. *International Economic Review*, pages 289–297, 1982.
- [45]. Vikas C Raykar, Shipeng Yu, Linda H Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca Bogoni, and Linda Moy. Learning from crowds. *The Journal of Machine Learning Research*, 11: 1297–1322, 2010.
- [46]. Lloyd Shapley and Bernard Grofman. Optimizing group judgmental accuracy in the presence of interdependencies. *Public Choice*, 43(3):329–343, 1984.
- [47]. Padhraic Smyth, Usama Fayyad, Michael Burl, Pietro Perona, and Pierre Baldi. Inferring ground truth from subjective labelling of venus images. *Advances in Neural Information Processing Systems*, pages 1085–1092, 1995.
- [48]. Alexandre B. Tsybakov. *Introduction to non-parametric estimation*. Springer, 2008.
- [49]. Dong Wang, Tarek Abdelzaher, Lance Kaplan, and Charu C Aggarwal. Recursive fact-finding: A streaming approach to truth estimation in crowdsourcing applications. In Proc. of IEEE ICDCS, 2013.
- [50]. Peter Welinder and Pietro Perona. Online crowdsourcing: rating annotators and obtaining cost-effective labels. In Proc. of IEEE CVPR (Workshops), 2010.
- [51]. Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier R Movellan, and Paul L Ruvolo. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In Proc. Of NIPS, 2009.
- [52]. Yuchen Zhang, Xi Chen, Dengyong Zhou, and Michael I Jordan. Spectral methods meet EM: A provably optimal algorithm for crowdsourcing. In Proc. of NIPS, <http://arxiv.org/abs/1406.3824>, 2014.
- [53]. Dengyong Zhou, Qiang Liu, John C Platt, Christopher Meek, and Nihar B Shah. Regularized minimax conditional entropy for crowdsourcing. arXiv preprint arXiv:1503.07240, 2015.

- [54]. Victor S. Sheng et al., for Get Another Label? Improving Data Quality and Data Mining Using Multiple, Noisy Labelers, KDD'08, August 24–27, 2008, Las Vegas, Nevada, USA : pp 614-622.
- [55]. Ofer Dekel and Ohad Shamir on VoxPopuli: Collecting High-Quality Labels from a Crowd in Twenty-Second Annual Conference on Learning Theory, 2009.
- [56]. Peter Welinder, Steve Branson, Serge Belongie, Pietro Perona for The Multidimensional Wisdom of Crowds in neural information processing systems, 2010.
- [57]. Fenglong Ma, Yaliang Li, Qi Li, Minghui Qiu, Jing Gao, Shi Zhi for FaitCrowd: Fine Grained Truth Discovery for Crowdsourced Data Aggregation in KDD'15, August 10-13, 2015, Sydney, NSW, Australia: 745-754
- [58]. Hesam Salehian, Patrick Howell and Chul Lee on Matching Restaurant Menus to Crowdsourced Food Data: A Scalable Machine Learning Approach, KDD 2017 Applied Data Science Paper.
- [59]. Thomas Bonald, Richard Combes. A Minimax Optimal Algorithm for Crowdsourcing. NIPS, 2017, Los Angeles, United States.
- [60]. Ryan Gomes, Peter Welinder, Andreas Krause, Pietro Perona for “Cloudclustering” in 2011.
- [61]. Fabian L. Wauthier, Michael L Jordan on Active Spectral Clustering via Iterative Uncertainty Reduction in August 12–16, 2012, Beijing, China.
- [62]. Jinfeng Yi, Rong Jin, Anil K. Jain, Shaili Jain, Tianbao Yang for Semi-Crowdsourced Clustering: Generalizing Crowd Labeling by Robust Distance Metric Learning, 2012.
- [63]. Susan Davidson and Sanjeev Khanna on Top-k and Clustering with Noisy Comparisons, ACM Trans. Datab. Syst. 39, 4, Article 35 (December 2014).
- [64]. Ramya Korlakoti Vinayak and Babak Hassibi for Crowdsourced Clustering: Querying Edges vs Triangles in 29th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.
- [65]. Arya Mazumdar and Barna Saha for Clustering With Noisy Queries in January 22nd 2017.
- [66]. Antti Ukkonen Department of Computer Science, University of Helsinki, Helsinki, Finland for Crowdsourced Correlation Clustering With Relative Distance Comparisons in 2017.
- [67]. He Jiang And Xin Chen, Dalian University Of Technology, Tieke He And Zhenyu Chen, Nanjing University Xiaochen Li for Fuzzy Clustering of Crowdsourced Test Reports for Apps in ACM Trans. Internet Technol. 18, 2, Article 18 (February 2018).
- [68]. Maryam M Najafabadi, Flavio Villanustre, Taghi M Khoshgoftar, Naeem Seliya, Randall Wald and Edin Muharemagic on Deep learning applications and challenges in big data analytics, Journal of Big Data (2015) 2:1.
- [69]. Shayan Doroudi, Ece Kamar, Emma Brunskill, Eric Horvitz for Toward a Learning Science for Complex Crowdsourcing Tasks, CHI'16, May 07-12, 2016, San Jose, CA, USA, 2016, ACM.
- [70]. Mohammad Abu Alsheikh, Dusit Niyato, Shaowei Lin, Hwee-Pink Tan, and Zhu Han for Mobile Big Data Analytics Using Deep Learning and Apache Spark
- [71]. Chong Sun, Narasimhan Rampalli, Frank Yang, An Hai Doan “Chimera: Large-Scale Classification using Machine Learning, Rules, and Crowdsourcing” 40th International Conference on Very Large Data Bases, September 1st - 5th 2014, Hangzhou, China. Proceedings of the VLDB Endowment, Vol. 7, No. 13 pages: 1529-1540
- [72]. Matthew Lease “On Quality Control and Machine Learning in Crowdsourcing” Human Computation: Papers from the 2011 AAAI Workshop (WS-11-11) pages: 97-102
- [73]. STEVEN BURROWS, MARTIN POTTHAST, and BENNO STEIN, Bauhaus-Universität Weimar “Paraphrase Acquisition via Crowdsourcing and Machine Learning” ACM Transactions on Intelligent Systems and Technology, Vol. 4, No. 3, Article 43, Publication date: June 2013.
- [74]. STEVEN BURROWS And MARTIN POTTHAST And BENNO STEIN “Paraphrase Acquisition via Crowdsourcing and Machine Learning” ACM Transactions on Intelligent Systems and Technology, Vol. V, No. N, January 2012, Pages 1–22.
- [75]. Ece Kamar, Severin Hacker, Eric Horvitz “Combining Human and Machine Intelligence in Large-scale Crowdsourcing” Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2012)
- [76]. Qianqian Xu, Ming Yan, Chendi Huang, Jiechao Xiong, Qingming Huang, Yuan Yao “Exploring Outliers in Crowdsourced Ranking for QoE” MM'17, October 23-27, 2017, Mountain View, CA, USA
- [77]. Honglei Zhuang, Aditya Parameswaran, Dan Roth, and Jiawei Han “Debiasing Crowdsourced Batches” KDD. 2015 August ; 2015: 1593–1602

- [78]. YukinoBabaandHisashiKashima
“StatisticalQualityEstimationforGeneralCrowdsourcingTasks”19thACMSIGKDDConferenceonKnowledge
Discoveryand Data Mining (KDD) (Baba and Kashima 2013).
- [79]. Amna Basharat, I. BudakArpinar, KhaledRasheed“ Leveraging Crowdsourcing for the Thematic
Annotation of the Qur’an” WWW’16 Companion, April 11–15, 2016, Montréal, Québec, Canada. ACM
978-1-4503-4144-8/16/04.
- [80]. Matthew Brehmer, Bongshin Lee, Petra Isenberg, EunChoe “Visualizing Ranges over Time on Mobile
Phones: A Task-Based Crowdsourced Evaluation,16 Aug 2018
- [81]. Shixia Liu, Changjian Chen, Yafeng Lu, FangxinOuyang, Bin Wang, An Interactive Method to Improve
Crowdsourced Annotations submitted on IEEE Trans viscomput Graph , 2018 Aug 20.
- [82]. XiangjieKong ,Menglin Li, Tao Tang , KaiqiTian,Luis Moreira-Matias , Feng Xia, Shared Subway Shuttle
Bus Route Planning Based on Transport Data Analytics submitted on IEEE Transactions on Automation
Science and Engineering, vol. 15, no. 4, october 2018.
- [83]. YoonjungKima,Choong-kiKima, Dong KunLeeb, Hyun-woo Leea, Rogelio II. T. Andradac, Quantifying
nature-based tourism in protected areas in developing countries by using social big data Published on 2015
Feb 24.
- [84]. Mark Birkin, Spatial data analytics of mobility with consumer data Published on 1 August 2018.
- [85]. Min Chen, Jun Yang, Long Hu, M. ShamimHossain, GhulamMuhammad,Urban Healthcare Big Data
System based on Crowdsourced and Cloud-based Air Quality Indicators Published on 25 october 2018.
- [86]. Mohammad MasudurRahmanChanchal K. Roy, Effective Reformulation of Query for Code Search using
Crowdsourced Knowledge and Extra-Large Data Analytics published on 23 July 2018