# Secure K-Nearest Neighbor Search by Keywords

**Goduguchita S Rajitha Priya[1], Annavazulu Mrinalini [2]**

[1] Dept. of CSE, J.N.T.U Anantapur, Tirupathi, Andhra Pradesh, India, Email- reen100.nou@gmail.com

[2] Dept. of CSE, SVU, Tirupathi, Andhra Pradesh, India

## Abstract

Information Mining has wide applications in numerous territories, for example, saving money, prescription, experimental examination and among government organizations. Grouping is one of the regularly utilized undertakings as a part of information mining applications. For as far back as decade, because of the ascent of different security issues, numerous hypothetical and down to earth answers for the characterization issue have been proposed under various security models. In any case, with the late prevalence of distributed computing, clients now have the chance to outsource their information, in scrambled structure, and in addition the information mining undertakings to the cloud. Since the information on the cloud is in scrambled structure, existing security saving characterization procedures are not pertinent. In this paper, we concentrate on tackling the arrangement issue over scrambled information. Specifically, we propose a protected k-NN classifier over scrambled information in the cloud. The proposed convention ensures the classification of information, protection of client's info inquiry, and conceals the information access designs. To the best of our insight, our work is the first to build up a safe k-NN classifier over encoded information under the semi-legit model. Additionally, we experimentally break down the proficiency of our proposed convention utilizing a genuine dataset under various parameter settings.

*Keywords: Encryption, Classification, PPkNN, SMC, Cloud Storage, AES, Cryptosystem.*

## 1. Introduction

As of late, the distributed computing worldview [1] is reforming the associations' method for working their information especially in the way they store, get to and process information. As a developing registering worldview, distributed computing pulls in numerous associations to consider genuinely in regards to cloud potential as far as its cost-productivity, adaptability, and offload of managerial overhead. Frequently, associations designate their computational operations notwithstanding their information to the cloud. Regardless of colossal favorable circumstances that the cloud offers, protection and security issues in the cloud are avoiding organizations to use those favorable circumstances. At the point when information is profoundly touchy, the information should be encoded before outsourcing to the cloud. However, when information are scrambled, independent of the hidden encryption plan, performing any information mining errands turns out to be extremely testing while never decoding the information.

### 1.1 Existig System

Existing work on Privacy-Preserving Data Mining (either annoyance or secure multi-party calculation based methodology) can't take care of the DMED issue. Annoyed information don't have semantic security, so information bother strategies can't be utilized to scramble very delicate information. Likewise the bothered information don't produce exceptionally exact information mining results. Secure multi-party calculation (SMC) based methodology accepts information are conveyed and not scrambled at each taking an interest gathering.

Because of space constraints, here we quickly audit the current related work and give a few definitions as a foundation. It would be ideal if you allude to our specialized report for a more explained related work and background. At to begin with, it appears to be completely homomorphism cryptosystems (e.g., [6]) can take care of the DMED issue since it permits an outsider (that has the encoded information) to execute discretionary capacities over scrambled information while never unscrambling them. Nonetheless, we stretch that such procedures are extremely costly and their utilization in down to earth applications have yet to be investigated. For instance, it was appeared in [7] that notwithstanding for frail security parameters one "bootstrapping" operation of the homomorphism operation would take no less than 30 seconds on a superior machine. It is conceivable to utilize the current mystery sharing systems in SMC, for example, Shamir's plan [8], to build up a PPkNN convention. Nonetheless, our work is not quite the same as the mystery sharing based arrangement in the accompanying angle. Arrangements in light of the mystery sharing plans require no less than three gatherings while our work require just two gatherings. For instance, the developments taking into account Share mind [9], a surely understood SMC structure which depends on the mystery sharing plan, accept that the quantity of taking interest gatherings is three. Therefore, our work is orthogonal to Share mind and other mystery sharing based plans. The main drawbacks are:

➢ No characterization issues have been proposed under various security models.

IJISET - International Journal of Innovative Science, Engineering & Technology, Vol. 3 Issue 4, April 2016.

www.ijiset.com

ISSN 2348 – 7968

- ➢ Existing approach expect information are dispersed and not scrambled.
- ➢ Data don't have semantic security

## 2. Proposed System

We concentrate on taking care of the characterization issue over encoded information. Specifically, we propose a safe k-NN classifier over scrambled information in the cloud. The proposed convention ensures the classification of information, protection of client's info question, and shrouds the information access designs. To the best of our insight, our work is the first to build up a protected k-NN classifier over scrambled information under the semi-fair model. Likewise, we exactly break down the proficiency of our proposed convention utilizing a genuine dataset under various parameter settings. We proposed novel techniques to successfully take care of the DMED issue accepting that the encoded information are outsourced to a cloud. In particular, we concentrate on the order issue since it is a standout amongst the most widely recognized information mining assignments. Since every grouping system has their own leverage, to be solid, this paper focuses on executing the k closest neighbor arrangement technique over encoded information in the distributed computing environment.

### 2.1 Advantages

- ➢ Solving the arrangement issue over scrambled information
- ➢ protects the classification of information
- ➢ Secure k-NN classifier over scrambled information under the semi-legitimate model3. System Architecture

### 2.2 Contributions

In this paper, we propose a novel PPkNN convention, a safe k-NN classifier over semantically secure encoded information. In our convention, once the scrambled information is outsourced to the cloud, Alice does not take an interest in any calculations. Consequently, no data is uncovered to Alice. What's more, our convention meets the accompanying security prerequisites:

- ✓ Contents of D or any middle of the road results ought to not be uncovered to the cloud.
- ✓ Bob's question q ought not be uncovered to the cloud.
- ✓ $C_q$ ought to be uncovered just to Bob. Likewise, no other data ought to be uncovered to Bob.
- ✓ Data access examples, for example, the records relating to the k-closest neighbors of q, ought to

not be uncovered to Bob and the cloud (to forestall any surmising assaults).

We underscore that the middle results seen by the cloud in our convention are either recently created randomized encryptions or arbitrary numbers. In this manner, which information records relate to the k-closest neighbors and the yield class name are not known not cloud. Likewise, subsequent to sending his scrambled inquiry record to the cloud, Bob does not include in any calculations. Consequently, information access examples are further shielded from Bob (see Section 5 for more points of interest).
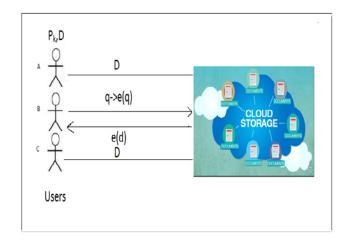
## 3. System Architecture

Fig. 1. Architecture

Fig.1 shows cloud storage and multiple users' relationships. And this architecture contains number of modules. They are Cloud Server, AES Cryptosystem.

### 3.1 Cloud server

A cloud server is a coherent server that is manufactured, facilitated and conveyed through a distributed computing stage over the Internet. Cloud servers have and show comparable capacities and usefulness to a run of the mill server however are gotten too remotely from a cloud administration supplier. A cloud server might likewise be known as a virtual server or virtual private disjoin..

### 3.2 AES Cryptosystem

- ➢ Stage 1 - Secure Retrieval of k-Nearest Neighbors (SRkNN): In this stage, Bob at first sends his inquiry q (in scrambled structure) to C1. After this, C1 and C2 include in an arrangement of sub-conventions to safely recover (in encoded

IJISET - International Journal of Innovative Science, Engineering & Technology, Vol. 3 Issue 4, April 2016.

www.ijiset.com

ISSN 2348 – 7968

structure) the class names relating to the k-closest neighbors of the data question q. Toward the end of this stride, encoded class names of k-closest neighbors are known just to C1.

➢ Stage 2 - Secure Computation of Majority Class (SCMCk):Following from Stage 1, C1 and C2 together figure the class name with a greater part voting among the k-closest neighbors of q. Toward the end of this stride, just Bob knows the class mark comparing to his information inquiry record q.

## 4. PPkNN Protocol

In this area, we propose a novel security saving k-NN arrangement convention, signified by PPkNN, which is built utilizing the conventions examined as a part of Section 3 as building squares. As specified before, we accept that Alice's database comprises of n records, denoted by $D = t_1, \ldots, t_n$, and m + 1 properties, where $t_{i,j}$ means the $j^{th}$ quality estimation of record $t_i$. At first, Alice scrambles her database quality shrewd, that is, she registers $E_{pk}(t_{i,j})$, for $1 \le i \le n$ and $1 \le j \le m + 1$, where section (m + 1) contains the class marks. Let the scrambled database be signified by $D'$. We expect that Alice outsources $D'$ and in addition the future grouping procedure to the cloud. Without loss of all inclusive statement, we expect that all trait values and their Euclidean separations lie in $[0, 2^l)$. In addition, let w mean the quantity of one of a kind class names in D.

In our issue setting, we accept the presence of two non-conniving semi-fair cloud administration suppliers, indicated by $C_1$ and $C_2$, which together shape a combined cloud. Under this setting, Alice outsources her encoded database $D'$ to $C_1$ and the mystery key sk to $C_2$. Here it is workable for the information proprietor Alice to supplant $C_2$ with her private server. Be that as it may, if Alice has a private server, we can contend that there is no requirement for information outsourcing from Alice's perspective. The principle reason for utilizing $C_2$ can be persuaded by the accompanying two reasons. (i) With constrained registering asset and specialized ability, it is to the greatest advantage of Alice to totally outsource its information administration and operational errands to a cloud. For case, Alice might need to get to her information and investigative results utilizing an advanced mobile phone or any gadget with exceptionally constrained registering ability. (ii) Suppose Bob needs to keep his info inquiry and access designs private from Alice. For this situation, if Alice utilizes a private server, then she needs to perform calculations accepted by $C_2$ under which the very reason for outsourcing the encoded information to $C_1$ is

discredited. All in all, whether Alice utilizes a private server or cloud administration supplier C2 really relies on upon her assets. Specifically to our issue setting, we like to utilize $C_2$ as this keeps away from the aforementioned disservices (i.e., if there should be an occurrence of Alice utilizing a private server) through and through. In our answer, in the wake of outsourcing scrambled information to the cloud, Alice does not partake in any future computations. The objective of the PPkNN convention is to arrange clients' question records utilizing D′ as a part of a security protecting manner. Consider an approved client Bob who needs to group his inquiry record $q = q_1, \ldots, q_m$ in view of D′ in $C_1$. The proposed PPkNN convention mostly comprises of the accompanying two stages:

Stage 1 - Secure Retrieval of k-Nearest Neighbors (SRkNN): In this stage, Bob at first sends his question q (in scrambled structure) to $C_1$. After this, $C_1$ and $C_2$ include in an arrangement of sub-conventions to safely recover (in scrambled structure) the class marks relating to the k-closest neighbors of the information inquiry q. Toward the end of this stride, scrambled class marks of k-closest neighbors are known just to $C_1$.

Stage 2 - Secure Computation of Majority Class (SCMCk): Following from Stage 1, $C_1$ and $C_2$ mutually process the class mark with a lion's share voting among the k-closest neighbors of q. toward the end of this stride; just Bob knows the class mark relating to his info question record q.

## 5. Literature Survey

### 4.1 Studies about Fully Homomorphic Encryption Using Ideal Lattices

We propose a completely homomorphic encryption plan i.e., a plan that permits one to assess circuits over encoded information without having the capacity to unscramble. Our answer comes in three stages. To begin with, we give a general result that, to build an encryption plan that allows assessment of subjective circuits, it suffices to develop an encryption plan that can assess (marginally expanded variants of) its own unscrambling circuit; we call a plan that can assess its (enlarged) decoding circuit bootstrappable. Next, we depict an open key encryption plan utilizing perfect cross sections that are verging on bootstrappable. Cross section based cryptosystems regularly have unscrambling calculations with low circuit unpredictability, frequently ruled by an inward item calculation that is in NC1. Additionally, perfect grids give both added substance and multiplicative homeomorphisms (modulo an open key perfect in a polynomial ring that is

IJISET - International Journal of Innovative Science, Engineering & Technology, Vol. 3 Issue 4, April 2016.
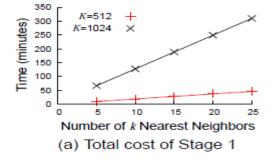
www.ijiset.com

ISSN 2348 – 7968

spoken to as a cross section), as expected to assess general circuits. Sadly, our underlying plan is not exactly bootstrappable – i.e., the profundity that the plan can accurately evaluate can be logarithmic in the grid measurement, much the same as the profundity of the decoding circuit, however the last is more prominent than the previous. In the last step, we demonstrate to adjust the plan to diminish the profundity of the unscrambling circuit, and in this way acquire a bootstrappable encryption plan, without lessening the profundity that the plan can assess. Stomach muscle strictly, we achieve this by empowering the encrypted to begin the unscrambling process, leaving less work for the decrypter, much like the server leaves less work for the decrypter in a server-helped cryptosystem.
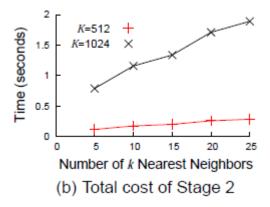
## 4.2 Studies about Building Castles out of Mud: Practical Access Pattern Privacy and Correctness on Entrusted Storage

We present another viable instrument for remote information stockpiling with effective access design protection and rightness. A capacity customer can convey this system to issue scrambled peruses, composes, and embeds to a possibly inquisitive and malignant capacity administration supplier, without uncovering data or access designs. The supplier can't set up any relationship between's progressive gets to, or even to recognize a read and a compose. In addition, the customer is given with solid rightness affirmations to its operations – illegal supplier conduct does not go undetected. We assembled a first reasonable framework request of extent speedier than existing usage – that can execute more than a few inquiries for every second on 1Tbyte+ databases with full computational security and accuracy.

## 6. Simulated Result

### 5.1 Performance of PPkNN

(a) Total cost of Stage 1

(b) Total cost of Stage 2
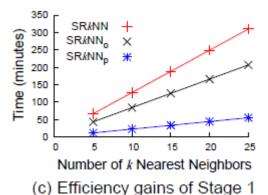
(c) Efficiency gains of Stage 1

Fig. 2. Computation costs of PPkNN for varying number of k nearest neighbors and encryption key size K.

In simulated result, we initially assessed the calculation expenses of Stage 1 in PPkNN for fluctuating number of k-closest neighbors. Also, the Paillier encryption key size K is either 512 on the other hand 1024 bits. The outcomes are appeared in Figure 2(a). For K=512 bits, the calculation expense of Stage 1 fluctuates from 9.98 to 46.16 minutes when k is changed from 5 to 25, separately. Then again, when K=1024 bits, the calculation expense of Stage 1 fluctuates from 66.97 to 309.98 minutes when k is changed from 5 to 25, respectively. In either case, we watched that the expense of Stage 1 becomes directly with k. For any given k, we distinguished that the expense of Stage 1 increments by just about a variable of 7 at whatever point k is multiplied. E.g., when k=10, Stage 1 took 19.06 and 127.72 minutes to produce the encoded class names of the 10 closest neighbors under K=512 and 1024 bits, separately. In addition, when k=5, we watch that around 66.29% of expense in Stage 1 is accounted because of SMINn which is started k times in PPkNN (once in each iteration).Also, the expense caused because of SMINn increments from 66.29% to 71.66% when k is expanded from 5 to 25.We now assess the calculation expenses of

IJISET - International Journal of Innovative Science, Engineering & Technology, Vol. 3 Issue 4, April 2016.

www.ijiset.com

Stage 2 for changing k and K. As appeared in Figure 2(b), for K=512 bits, the calculation time for Stage 2 to produce the last class name comparing to the data question shifts from 0.118 to 0.285 seconds when k is changed from 5 to 25. Then again, for K=1024 bits, Stage 2 took 0.789 and 1.89 seconds when k = 5 and 25, separately. The low calculation expenses of Stage 2 were because of SMAXw which brings about altogether less calculation than SMINn in Stage 1.This further legitimizes our hypothetical examination in Section 5.5. Note that, in our dataset, w=4 and n=1728. Like in Stage 1, for any given k, the calculation time of Stage 2 increments by very nearly a variable of 7 at whatever point K is multiplied. E.g., when k=10, the calculation time of Stage 2 differs from 0.175 to 1.158 seconds when the encryption key size K is changed from 512 to 1024 bits. As appeared in Figure 2(b), a comparable examination can be watched for different estimations of k and K. It is clear that the calculation expense of Stage 1 is essentially higher than that of Stage 2 in PPkNN. In particular, we watched that the calculation time of Stage 1 represents no less than 99% of the aggregate time in PPkNN. For instance, when k = 10 and K=512 bits, the calculation expenses of Stage 1 and 2 are 19.06 minutes and 0.175 seconds, individually. Under this situation, expense of Stage 1 is 99.98% of the aggregate expense of PPkNN. We likewise watched that the aggregate calculation time of PPkNN becomes directly with n and k.6. Conclusion and Future Work

This paper initially performed client study to investigate the wellbeing seeker needs. This gives the bits of knowledge of group based wellbeing administrations. It then displayed a scantily associated profound learning plot that can surmise the conceivable sicknesses given the inquiries of wellbeing seekers. This plan is built by means of option mark mining and pre preparing in an incremental way. It licenses unsupervised element gaining from other extensive variety of infection sorts. In this way, it is generalizable and adaptable as thought about to past infection surmising utilizing shallow learning approaches, which are typically prepared on healing facility created persistent records with organized fields. Traditional profound learning structures are thickly associated and the hub number in each concealed layers are dully balanced. In contract, our model is scantily associated with enhanced learning proficiency; what's more, the quantity of shrouded hubs is consequently decided.

Our ebb and flow model can't distinguish discriminant highlights for every particular infection. Later on, we will give careful consideration on that.

# 7. Conclusion and Future Work

To ensure client protection, different security saving grouping methods has been proposed over the previous decade. The current procedures are not pertinent to outsourced database situations where the information lives in encoded structure on an outsider server. This paper proposed a novel security safeguarding k-NN arrangement convention over encoded information in the cloud. Our convention ensures the privacy of the information, client's info question, and shrouds the information access designs. We additionally assessed the execution of our convention under various parameter settings.

Since enhancing the effectiveness of SMINn is a critical initial step for enhancing the execution of our PPkNN convention, we plan to examine option and more productive answers for the SMINn issue in our future work. Additionally, we will explore and extend our examination to other grouping calculations.

## References

[[1] P. Mell and T. Grance, "The nist definition of cloud computing (draft)," NIST special publication, vol. 800, p. 145, 2011.

[2] S. De Capitani di Vimercati, S. Foresti, and P. Samarati, "Managing and accessing data in the cloud: Privacy risks and approaches," in CRiSIS, pp. 1 –9, 2012.

[3] P. Williams, R. Sion, and B. Carbunar, "Building castles out of mud: practical access pattern privacy and correctness on untrusted storage," in ACM CCS, pp. 139–148, 2008.

[4] P. Paillier, "Public key cryptosystems based on composite degree residuosity classes," in Eurocrypt, pp. 223–238, 1999.

[5] B. K. Samanthula, Y. Elmehdwi, and W. Jiang, "k-nearest neighbor classification over semantically secure encrypted relational data." eprint arXiv:1403.5001, 2014.

[6] C. Gentry, "Fully homomorphic encryption using ideal lattices," in ACM STOC, pp. 169–178, 2009.

[7] C. Gentry and S. Halevi, "Implementing gentry's fullyhomomorphic encryption scheme," in EUROCRYPT, pp. 129–148, Springer, 2011.

[8] A. Shamir, "How to share a secret," Commun. ACM, vol. 22, pp. 612–613, Nov. 1979.

[9] D. Bogdanov, S. Laur, and J. Willemson, "Sharemind: A framework for fast privacy-preserving computations," in ESORICS, pp. 192–206, Springer, 2008.

[10] R. Agrawal and R. Srikant, "Privacy-preserving data mining," in ACM Sigmod Record, vol. 29, pp. 439–450, ACM, 2000.

[11] Y. Lindell and B. Pinkas, "Privacy preserving data mining," in Advances in Cryptology (CRYPTO), pp. 36–54, Springer, 2000.

[12] P. Zhang, Y. Tong, S. Tang, and D. Yang, "Privacy preserving naive bayes classification," ADMA, pp. 744–752, 2005.

[13] A. Evfimievski, R. Srikant, R. Agrawal, and J. Gehrke, "Privacy preserving mining of association rules," Information Systems, vol. 29, no. 4, pp. 343–364, 2004.

[14] R. J. Bayardo and R. Agrawal, "Data privacy through optimal k-anonymization," in IEEE ICDE, pp. 217–228, 2005.

[15] H. Hu, J. Xu, C. Ren, and B. Choi, "Processing private queries over untrusted data cloud through privacy homomorphism," in IEEE ICDE, pp. 601–612, 2011.

**Goduguchita S Rajitha** received the B.Tech Degree in Computer Science and Engineering from S.V.College of Engineering, University of JNTUA in 2014.She is currently working towards the Master's Degree in Computer Science, in AITS University of JNTUA. She interest lies in the areas of Web Development Platforms, SQL, and Cloud Computing Technology.

**Annavazulu Mrinalini** received M.Tech in SVU. Currently he is an Assistant Professor in the Department of Computer Science at AITS-Tirupati.