Large-Scale Nearest Neighbor Classification with Statistical Guarantees

Guang Cheng

Big Data Theory Lab Department of Statistics Purdue University

Joint Work with Xingye Qiao and Jiexin Duan

July 3, 2018

A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

SUSY Data Set (UCI Machine Learning Repository)

- Size: 5,000,000 particles from the accelerator
- Predictor variables: 18 (properties of the particles)
- Goal Classification:

distinguish between a signal process and a background process



2 / 48

The *k*NN classifier predicts the class of $x \in \mathbb{R}^d$ to be the most frequent class of its *k* nearest neighbors (Euclidean distance).



Computational/Space Challenges for kNN Classifiers

- Time complexity: O(dN + Nlog(N))
 - dN: computing distances from the query point to all N observations in \mathbb{R}^d
 - *Nlog(N)*: sorting *N* distances and selecting the *k* nearest distances
- Space complexity: O(dN)
- Here, N is the size of the entire dataset



How to Conquer "Big Data?"



- E - N

A (1) > (1)

Train the total data at one time: oracle kNN



Figure: I'm Pac-Superman!

6 / 48



To build a Pac-superman for oracle-kNN, we need:

• Expensive super-computer



To build a Pac-superman for oracle-kNN, we need:

- Expensive super-computer
- Operation system, e.g. Linux



To build a Pac-superman for oracle-kNN, we need:

- Expensive super-computer
- Operation system, e.g. Linux
- Statistical software designed for big data, e.g. Spark, Hadoop



To build a Pac-superman for oracle-kNN, we need:

- Expensive super-computer
- Operation system, e.g. Linux
- Statistical software designed for big data, e.g. Spark, Hadoop
- Write complicated algorithms (e.g. MapReduce) with limited extendibility to other classification methods

Machine Learning Literature

- In machine learning area, there are some nearest neighbor algorithms for big data already
 - Muja&Lowe(2014) designed an algorithm to search nearest neighbors for big data
 - Comments: complicated; without statistical gurantee; lack of extendibility
- We hope to design an algorithm framework:
 - Easy to implement
 - Strong statistical guarantee
 - More generic and expandable (e.g. easy to be applied to other classifiers (decision trees, SVM, etc)

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

If We Don't Have A Supercomputer...

What can we do?



A (10) > (10)

What about splitting the data?



Divide & Conquer (D&C) Framework

Goal: deliver different insights from regression problems.



Big-kNN for SUSY Data

- Oracle-kNN: kNN trained by the entire data
- Number of subsets: $s = N^{\gamma}$, $\gamma = 0.1, 0.2, \dots 0.8$



Figure: Risk and Speedup for Big-kNN.

- Risk of classifier ϕ : $R(\phi) = \mathbb{P}(\phi(X) \neq Y)$
- Speedup: running time ratio between Oracle-kNN & Big-kNN

Construction of Pac-Batman (Big Data)



イロト イポト イヨト イヨト

æ

We need a *statistical* understanding of Big-*k*NN....

臣

▶ < 문 ▶ < 문 ▶</p>

< A

Definition: WNN

In a dataset with size n, the WNN classifier assigns a weight w_{ni} on the *i*-th neighbor of x:

$$\widehat{\phi}_{n}^{w_{n}}(x) = \mathbb{1}\left\{\sum_{i=1}^{n} w_{ni}\mathbb{1}\{Y_{(i)}=1\} \ge \frac{1}{2}\right\} \ s.t. \ \sum_{i=1}^{n} w_{ni} = 1,$$

When $w_{ni} = k^{-1} \mathbb{1}\{1 \le i \le k\}$, WNN reduces to kNN.

・ロト ・日本・ ・ ヨト・

15 / 48

Construction of Big-WNN

In a big data set with size N = sn, the Big-WNN is constructed as:

$$\widehat{\phi}_{n,s}^{Big}(x) = \mathbb{1}\left\{s^{-1}\sum_{j=1}^{s}\widehat{\phi}_{n}^{(j)}(x) \ge 1/2\right\},\$$

where $\widehat{\phi}_n^{(j)}(x)$ is the local WNN in *j*-th subset.

Statistical Guarantees: Classification Accuracy & Stability

- Primary criterion: accuracy
 - Regret=Expected Risk-Bayes Risk= $\mathbb{E}_{\mathcal{D}}\left[R(\widehat{\phi}_n)\right] R(\phi^{\text{Bayes}})$
 - A small value of regret is preferred

Statistical Guarantees: Classification Accuracy & Stability

- Primary criterion: accuracy
 - Regret=Expected Risk-Bayes Risk= $\mathbb{E}_{\mathcal{D}} \left| R(\hat{\phi}_n) \right| R(\phi^{\text{Bayes}})$
 - A small value of regret is preferred
- Secondary criterion: stability (Sun et al, 2016, JASA)

Definition: Classification Instability (CIS)

Define classification instability of a classification procedure $\boldsymbol{\Psi}$ as

$$\operatorname{CIS}(\Psi) = \mathbb{E}_{\mathcal{D}_1, \mathcal{D}_2} \Big[\mathbb{P}_X \Big(\widehat{\phi}_{n1}(X) \neq \widehat{\phi}_{n2}(X) \Big) \Big],$$

where $\widehat{\phi}_{n1}$ and $\widehat{\phi}_{n2}$ are the classifiers trained from the same classification procedure Ψ based on \mathcal{D}_1 and \mathcal{D}_2 with the same distribution.

• A small value of CIS is preferred

・ロト ・ 日 ・ ・ 日 ・

Theorem (Regret of Big-WNN)

Under regularity assumptions, with s upper bounded by the subset size n, we have as $n, s \to \infty$,

Regret (Big-WNN)
$$\approx B_1 s^{-1} \sum_{i=1}^n w_{ni}^2 + B_2 \Big(\sum_{i=1}^n \frac{\alpha_i w_{ni}}{n^{2/d}} \Big)^2,$$

where $\alpha_i = i^{1+\frac{2}{d}} - (i-1)^{1+\frac{2}{d}}$, w_{ni} are the local weights. The constants B_1 and B_2 are distribution-dependent quantities.

・ロト ・ 日 ・ ・ 日 ・

Theorem (Regret of Big-WNN)

Under regularity assumptions, with s upper bounded by the subset size n, we have as $n, s \to \infty$,

Regret (Big-WNN)
$$\Rightarrow B_1 s^{-1} \sum_{i=1}^{n} w_{ni}^2 + B_2 \left(\sum_{i=1}^{n} \frac{\alpha_i w_{ni}}{n^{2/d}} \right)^2$$
,
where $\alpha_i = i^{1+\frac{2}{d}} + (i-1)^{1+\frac{2}{d}}$, w_{ni} are the local weights. The onstants B_1 and B_2 are distribution-dependent quantities.

Variance part

С

Theorem (Regret of Big-WNN)

Under regularity assumptions, with s upper bounded by the subset size n, we have as $n, s \rightarrow \infty$,



Asymptotic Regret Comparison



イロト イヨト イヨト イヨト

э

Theorem (Regret Ratio of Big-WNN and Oracle-WNN)

Given an Oracle-WNN classifier, we can always design a Big-WNN by adjusting its weights according to those in the oracle version s.t.

 $\frac{\textit{Regret (Big-WNN)}}{\textit{Regret (Oracle-WNN)}} \rightarrow Q,$

where $Q = (\frac{\pi}{2})^{\frac{4}{d+4}}$ and 1 < Q < 2.

20 / 48

Theorem (Regret Ratio of Big-WNN and Oracle-WNN)

Given an Oracle-WNN classifier, we can always design a Big-WNN by adjusting its weights according to those in the oracle version s.t.

 $\frac{Regret (Big-WNN)}{Regret (Oracle-WNN)} \to Q,$

where $Q = \left(\frac{\pi}{2}\right)^{\frac{4}{d+4}}$ and 1 < Q < 2.

• E.g., in Big-*k*NN case, we set $k = \lfloor \left(\frac{\pi}{2}\right)^{\frac{d}{d+4}} \frac{k^0}{s} \rfloor$ given the oracle k^0 . We need a multiplicative constant $\left(\frac{\pi}{2}\right)^{\frac{d}{d+4}}!$

・ロト ・日子・ ・日ト

20 / 48

Theorem (Regret Ratio of Big-WNN and Oracle-WNN)

Given an Oracle-WNN classifier, we can always design a Big-WNN by adjusting its weights according to those in the oracle version s.t.

 $\frac{Regret (Big-WNN)}{Regret (Oracle-WNN)} \to Q,$

where $Q = \left(\frac{\pi}{2}\right)^{\frac{4}{d+4}}$ and 1 < Q < 2.

- E.g., in Big-*k*NN case, we set $k = \lfloor \left(\frac{\pi}{2}\right)^{\frac{d}{d+4}} \frac{k^O}{s} \rfloor$ given the oracle k^O . We need a multiplicative constant $\left(\frac{\pi}{2}\right)^{\frac{d}{d+4}}!$
- We name Q the Majority Voting (MV) Constant.

・ロット 全部 マイロッ

MV constant monotonically decreases to one as d increases



For example:

- d=1, Q=1.44
- d=2, Q=1.35
- d=5, Q=1.22
- d=10, Q=1.14
- d=20, Q=1.08
- d=50, Q=1.03
- d=100, Q=1.02

How to Compensate Statistical Accuracy Loss, i.e., Q > 1

Why is there statistical accuracy loss?



Statistical Accuracy Loss due to Majority Voting

Accuracy loss due to the transformation from (continuous) percentage to (discrete) 0-1 label



Majority Voting/Accuracy Loss Once in Oracle Classifier

Only one majority voting in oracle classifier



Figure: I'm Pac-Superman!

Majority Voting/Accuracy Loss Twice in D&C Framework



A (10) > (10)

문 🕨 문



Is it possible to apply majority voting once in D&C?



A Continuous Version of D&C Framework



< ロ > < 同 > < 回 > < 回 > < 回 > <

æ

Only Loss Once in the Continuous Version



≣) ≣

Construction of Pac-**Bat**man (Continuous Version)



< ロ > < 同 > < 回 > < 回 > < 回 > <

æ

Continuous Version of Big-WNN (C-Big-WNN)

Definition: C-Big-WNN

In a data set with size N = sn, the C-Big-WNN is constructed as:

$$\widehat{\phi}_{n,s}^{CBig}(x) = \mathbb{1}\left\{\frac{1}{s}\sum_{j=1}^{s}\widehat{S}_{n}^{(j)}(x) \ge 1/2\right\},\$$

where $\widehat{S}_{n}^{(j)}(x) = \sum_{i=1}^{n} w_{ni,j} Y_{(i),j}(x)$ is the weighted average for nearest neighbors of query point x in *j*-th subset.

Continuous Version of Big-WNN (C-Big-WNN)

Definition: C-Big-WNN

In a data set with size N = sn, the C-Big-WNN is constructed as:

$$\widehat{\phi}_{n,s}^{CBig}(x) = \mathbb{1}\left\{\frac{1}{s}\sum_{j=1}^{s}\widehat{S}_{n}^{(j)}(x) \ge 1/2\right\},\$$

where $\widehat{S}_{n}^{(j)}(x) = \sum_{i=1}^{n} w_{ni,j} Y_{(i),j}(x)$ is the weighted average for nearest neighbors of query point x in *j*-th subset.

Step 1: Percentage calculation

イロト イポト イヨト イヨト

30 / 48

Continuous Version of Big-WNN (C-Big-WNN)

Definition: C-Big-WNN

In a data set with size N = sn, the C-Big-WNN is constructed as:

$$\widehat{\phi}_{n,s}^{CBig}(x) = \mathbb{1}\left\{\frac{1}{s}\sum_{j=1}^{s}\widehat{S}_{n}^{(j)}(x) \geq 1/2\right\},$$

where $\widehat{S}_{n}^{(j)}(x) = \sum_{i=1}^{n} w_{ni,j} Y_{(i),j}(x)$ is the weighted average for nearest neighbors of query point x in *j*-th subset.

Step 1: Percentage calculation Step 2: Majority voting

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

30 / 48

C-Big-kNN for SUSY Data



Figure: Risk and Speedup for Big-*k*NN and C-Big-*k*NN.

Э

Asymptotic Regret Comparison



・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

æ

Theorem (Regret Ratio between C-Big-WNN and Oracle-WNN)

Given an Oracle-WNN, we can always design a C-Big-WNN by adjusting its weights according to those in the oracle version s.t.

 $\frac{\textit{Regret (C-Big-WNN)}}{\textit{Regret (Oracle-WNN)}} \rightarrow 1,$

・ロト ・日子・ ・日子

Theorem (Regret Ratio between C-Big-WNN and Oracle-WNN)

Given an Oracle-WNN, we can always design a C-Big-WNN by adjusting its weights according to those in the oracle version s.t.

 $\frac{\textit{Regret (C-Big-WNN)}}{\textit{Regret (Oracle-WNN)}} \rightarrow 1,$

E.g., in the C-Big-kNN case, we can set k = \[\frac{k^O}{s} \] given the oracle k^O. There is no multiplicative constant!

Theorem (Regret Ratio between C-Big-WNN and Oracle-WNN)

Given an Oracle-WNN, we can always design a C-Big-WNN by adjusting its weights according to those in the oracle version s.t.

 $\frac{\textit{Regret}(\textit{C-Big-WNN})}{\textit{Regret}(\textit{Oracle-WNN})} \rightarrow 1,$

- E.g., in the C-Big-kNN case, we can set k = [k^O/s] given the oracle k^O. There is no multiplicative constant!
- Note that this theorem can be directly applied to the optimally weighted NN (Samworth, 2012, *AoS*), i.e., OWNN, whose weights are optimized to lead to the minimal regret.

・ コ ト ・ 白 ト ・ 山 ト ・

 It is time consuming to apply OWNN in each subset although it leads to the minimal regret

- It is time consuming to apply OWNN in each subset although it leads to the minimal regret
- Rather, it is more easier and faster to implement kNN in each subset (with accuracy loss, though)

- It is time consuming to apply OWNN in each subset although it leads to the minimal regret
- Rather, it is more easier and faster to implement kNN in each subset (with accuracy loss, though)
- What is the statistical price (in terms of accuracy and stability) in order to trade accuracy with speed?

Theorem (Statistical Price)

We can design an optimal C-Big-kNN by setting its weights as $k = \lfloor \frac{k^{O,opt}}{s} \rfloor$ s.t.

$$rac{Regret(Optimal \ C-Big-kNN)}{Regret(Oracle-OWNN)}
ightarrow Q',$$

$$\frac{CIS(Optimal \ C-Big-kNN)}{CIS(Oracle-OWNN)} \to \sqrt{Q'},$$

where $Q' = 2^{-4/d+4} (\frac{d+4}{d+2})^{(2d+4)/(d+4)}$ and 1 < Q' < 2. Here, $k^{O,opt}$ minimizes the regret of Oracle-kNN.

Interestingly, both Q and Q' only depend on data dimension!

・ロン ・四シ ・ヨン ・ヨン

Q' converges to one as d grows in a unimodal way



The worst dimension is $d^* = 4 \Longrightarrow Q^* = 1.089$.

Consider the classification problem for Big-*k*NN and C-Big-*k*NN:

- Sample size: N = 27,000
- Dimensions: d = 6, 8
- $P_0 \sim N(0_d, \mathbb{I}_d)$ and $P_1 \sim N(rac{2}{\sqrt{d}} \mathbb{1}_d, \mathbb{I}_d)$
- Prior class probability: $\pi_1 = Pr(Y = 1) = 1/3$
- Number of neighbors in Oracle-kNN: $k^O = N^{0.7}$
- Number of subsamples in D&C: $s = N^{\gamma}$, $\gamma = 0.1, 0.2, \dots 0.8$
- Number of neighbors in Big-*k*NN: $k^d = \lfloor (\frac{\pi}{2})^{\frac{d}{d+4}} \frac{k^O}{s} \rfloor$
- Number of neighbors in C-Big-kNN: $k^{c} = \lfloor \frac{k^{O}}{s} \rfloor$

Simulation Analysis – Empirical Risk



Figure: Empirical Risk (Testing Error). Left/Right: d = 6/8.

イロト イポト イヨト イヨト

臣

Simulation Analysis – Running Time



Figure: Running Time. Left/Right: d = 6/8.

Simulation Analysis – Empirical Regret Ratio



Figure: Empirical Ratio of Regret. Left/Right: d = 6/8.

 Q: Regret(Big-kNN)/Regret(Oracle-kNN) or Regret(C-Big-kNN)/Regret(Oracle-kNN)

Simulation Analysis – Empirical CIS



Figure: Empirical CIS. Left/Right: d = 6/8.

イロト イポト イヨト イヨト

æ

Data	Size	Dim	Big- <i>k</i> NN	C-Big- <i>k</i> NN	Oracle- <i>k</i> NN	Speedup
htru2	17898	8	3.72	3.36	3.34	21.27
gisette	6000	5000	12.86	10.77	10.78	12.83
musk1	476	166	36.43	33.87	35.71	3.6
musk2	6598	166	10.43	9.98	10.14	14.68
occup	20560	6	5.09	4.87	5.19	20.31
credit	30000	24	21.85	21.37	21.5	22.44
SUSY	5000000	18	30.66	30.08	30.01	68.45

Table: Test error (Risk): Big-kNN compared to Oracle-kNN in real datasets. Best performance is shown in bold-face. The speedup factor is defined as computing time of Oracle-kNN divided by the time of the slower Big-kNN method. Oracle $k = N^{0.7}$, number of subsets $s = N^{0.3}$.



・ロト ・ 日 ・ ・ ヨ ・ ・ モ ・ ・

2