# On Differentially Private Gaussian Hypothesis Testing

Kwassi H. Degue and Jerome Le Ny

*Abstract*—Data analysis for emerging systems such as syndromic surveillance or intelligent transportation systems requires testing statistical models based on privacy-sensitive data collected from individuals, e.g., medical records or location traces. In this paper, we design a differentially private hypothesis test based on the generalized likelihood ratio method to decide if data modeled as a sequence of independent and identically distributed Gaussian random variables has a given mean value. Analytic formulas for decision thresholds and for the test's receiver operating characteristic curve show explicitly the performance impact of the privacy constraint. We then apply the algorithm to the design of a differentially private anomaly (or fault) detector and study its performance for the analysis of a syndromic surveillance dataset from the Centers for Disease Control and Prevention in the United States.

## I. INTRODUCTION

Hypothesis testing provides a structured approach to compare and select models with reasonable confidence based on observed data. This data increasingly contains highly privacy-sensitive information collected from individuals, such as medical records, wages and location traces. Data analysis methods must take this fact into account and provide guarantees to survey participants that a published result does not significantly increase the likelihood of privacy breaches [1], [2].

Many recent approaches to privacy-preserving data analysis are based on the notion of *differential privacy* [3], [4]. When the result of a data mining computation based on private information from individuals is released, differential privacy promises to these individuals that whether or not they choose to provide their data will not significantly change an adversary's ability to deduct new knowledge about them, and in exchange they might obtain important benefits from the system [5], [6].

For categorical data following a multinomial distribution, differentially private mechanisms have been studied for various classical hypothesis tests in [7]–[13] for example. In particular, [9] and [13] compute detection thresholds providing target levels of significance by using Monte

Carlo (MC) approaches. However, many applications rely on numerical data such as sensor measurements instead of categorical ones. Ghassemi et al. [14] describe a differentially private mechanism to train an anomaly detector. To the best of our knowledge only [15] considers the design of a statistical test (a binary classifier) for normally distributed data. A differentially private sequential anomaly detector based on the cumulative sums (CUSUM) algorithm is proposed in [16]. CUSUM algorithms, however, can respond slowly when it comes to detecting large shifts in the mean of a stochastic process [17]. Moreover, [16] considers only scalar processes with a trivial state-space representation.

The first contribution of this paper lies in designing a differentially private statistical test to decide if the mean of a sequence of independent and identically distributed (iid) Gaussian random variables has a given value. We use the generalized likelihood ratio method [18], where the value of the a priori unknown mean is replaced by its maximum likelihood estimate, here the empirical mean. Instead of adding directly the privacy-preserving noise to the data, we perturb the empirical mean used at decision times, which reduces the impact of the privacy-preserving noise. In addition, in the Gaussian model analytical formulas can be used to determine detection thresholds and performance, instead of MC approaches as in [9], [13], The second contribution lies in designing an anomaly detection algorithm that preserves differential privacy for individuals' measurements data, which are assumed to originate from a general linear time-invariant Gaussian dynamic model.

In Section II, we present the problem statement and we provide some background on differential privacy and hypothesis testing. We apply these results in Section III to design the differentially private hypothesis test, whose performance is analyzed in Section IV. We provide analytic formulas to set appropriate detection thresholds and determine the detector's receiver operating characteristic (ROC) curves. In Section V, we extend the discussion to sequential detection and Section VI presents an application of the results to an epidemic outbreak detection scenario.

*Notation:* We fix a generic probability triple $(\Omega, \mathcal{F}, \mathbb{P})$, where $\mathcal{F}$ stands for a $\sigma$-algebra on the sample space $\Omega$ and $\mathbb{P}$ is a probability measure defined on $\mathcal{F}$. The $\ell_p$-norm of a vector $x \in \mathbb{R}^k$ is denoted by $|x|_p := (\sum_{i=1}^{k} |x_i|^p)^{1/p}$,

K. H. Degue and J. Le Ny are with the department of Electrical Engineering, Polytechnique Montreal and GERAD, QC H3T-1J4, Montreal, Canada `kwassi-holali.degue,jerome.le-ny@polymtl.ca`

for $p \in [1, \infty]$. To indicate that a random vector $X$ is distributed according to a normal (or Gaussian) distribution with mean $\mu$ and covariance matrix $\Sigma$, we use the notation $X \sim \mathcal{N}(\mu, \Sigma)$.

## II. PROBLEM STATEMENT

### A. Gaussian Hypothesis Testing Problem

Consider a sequence $r := \{r_i\}_{i=1}^{i=n}$ of $n$ iid scalar Gaussian random variables with unknown mean value $\theta$ and known standard deviation $\sigma$. In other words, the probability density of each sample $r_i$ can be written as follows

$$p_\theta(r_i) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(r_i-\theta)^2}{2\sigma^2}}.$$

We aim at designing a statistical test to determine if the mean value $\theta$ is equal to zero or not, in other words, a rule mapping the values in the sequence $r$ to a decision between the following two hypotheses

$$\begin{cases} H_0 : & \theta = \theta_0 = 0 \\ H_1 : & \theta = \theta_1 \neq 0, \end{cases} \tag{1}$$

with $\theta_1$ unknown. The discussion generalizes to any known value of $\theta_0$ by redefining $\tilde{\theta} := \theta - \theta_0$. For a given test, denote by $P_I$ the probability of incorrectly rejecting $H_0$ (type I error), and by $P_{II}$ the probability of incorrectly accepting $H_0$ (type II error).

Define the log-likelihood ratio $l(r)$ for the data from $r_1$ to $r_n$ as follows

$$l(r) = \sum_{i=1}^{n} s_i(r_i), \text{ with } s_i(r_i) = \ln \frac{p_{\theta_1}(r_i)}{p_{\theta_0}(r_i)}. \tag{2}$$

Here, in the Gaussian case with $\theta_0 = 0$, we have

$$l(r) = \frac{1}{2\sigma^2} \sum_{i=1}^{n} \left[ r_i^2 - (r_i - \theta_1)^2 \right],$$

$$l(r) = \frac{\theta_1}{\sigma^2} \sum_{i=1}^{n} \left( r_i - \frac{\theta_1}{2} \right). \tag{3}$$

If $\theta_1$ were known, $l$ could be computed explicitly from the data and the Neyman-Pearson lemma [18] shows that optimal decision rules $d$ (i.e., *most powerful tests*) to minimize the probability of type II error given an acceptable probability of type I error are of the form

$$d(r) = \begin{cases} 0 & \text{if } l(r) < h : H_0 \text{ is chosen} \\ 1 & \text{if } l(r) > h : H_1 \text{ is chosen}, \end{cases} \tag{4}$$

where the threshold $h$ is adequately chosen to satisfy a given bound on the probability of type I error. In the case where $l = h$, one might have to randomize in (4), but this case occurs with probability 0 for Gaussian random variables.

Since in our case the value of $\theta_1$ is unknown but enters the expression (3), the decision rule (4) cannot be directly used. In this case, a potentially suboptimal but very common test is the Generalized Likelihood Ratio (GLR) test [19], which consists here in replacing $\theta_1$ in the numerator of (2) by the maximum likelihood (ML) estimate of $\theta_1$ computed from $r$, assuming that $r_i \sim \mathcal{N}(\theta_1, \sigma^2)$. In the Gaussian case the ML estimate of $\theta_1$ is just the sample mean $\hat{\theta}_1 = \frac{1}{n} \sum_{i=1}^{n} r_i$. Replacing $\theta_1$ by $\hat{\theta}_1$ gives

$$\hat{l}(r) = \frac{1}{2\sigma^2 n} \left( \sum_{i=1}^{n} r_i \right)^2 \tag{5}$$

and the GLR test reads

$$d(r) = \begin{cases} 0 & \text{if } \hat{l}(r) < h : H_0 \text{ is chosen} \\ 1 & \text{if } \hat{l}(r) > h : H_1 \text{ is chosen}. \end{cases} \tag{6}$$

The decision rule (6) depends on the data $r$, which in this paper is assumed to be privacy-sensitive. As explained in the next subsection, we aim to publish the result of statistical tests for $H_0$ under a differential privacy constraint. This requires modifying the decision rule (6), which is done in Section III.

### B. Differentially Privacy

Given a space $\mathcal{H}$ of datasets, we define a *mechanism $M$* as a random map from $\mathcal{H}$ to some measurable output space $\mathcal{O}$. The output of a differentially private mechanism [5] should have similar distributions for inputs that we want to make hard to distinguish. Formally, this requires defining a symmetric binary relation Adj on $\mathcal{H}$, called adjacency, to specify which inputs are in some sense considered close. Typically, two adjacent inputs differ by the data of a single individual. In this paper, $\mathcal{H}$ is the space $\mathbb{R}^n$ of sequences $r$ and we consider the following adjacency relation

$$\text{Adj}(r, r') \text{ iff } |r - r'|_1 \leq \rho, \tag{7}$$

with $\rho \in \mathbb{R}_+$ a given positive number. Therefore, with the interpretation of adjacent sequences above, we assume that a single participant contributes additively to possibly each $r_i$ but in such a way that its overall influence on the whole sequence is bounded in 1-norm by $\rho$. We now provide the formal definition of differential privacy [3], [4].

**Definition 1.** *Consider $\mathcal{H}$, a space equipped with a symmetric binary relation denoted Adj, and $(\mathcal{O}, \mathcal{M})$ a measurable space, where $\mathcal{M}$ stands for a given $\sigma$-algebra over $\mathcal{O}$. Let $\epsilon, \delta \geq 0$. A randomized mechanism $M$ from $\mathcal{H}$ to $\mathcal{O}$ is $(\epsilon, \delta)$-differentially private (for Adj) if for all $r, r' \in \mathcal{H}$ such that Adj$(r, r')$, for all sets $S$ in $\mathcal{M}$,*

$$\mathbb{P}(M(r) \in S) \leq e^\epsilon \mathbb{P}(M(r') \in S) + \delta. \tag{8}$$

*If $\delta = 0$, the mechanism is said to be $\epsilon$-differentially private.*

Note that (8) expresses the fact that the distributions of the random variables $M(r)$ and $M(r')$ are close for $r$ and $r'$ adjacent. Next, we describe a basic mechanism, the Gaussian mechanism [20], which can be used to publish differentially private numerical outputs by adding an appropriate amount of Gaussian noise to a non-private output.

**Definition 2.** *Consider $\mathcal{H}$ equipped with an adjacency relation Adj and let $\mathcal{O}$ be a vector space with norm $\| \cdot \|_{\mathcal{O}}$. The sensitivity of a query $q : \mathcal{H} \mapsto \mathcal{O}$ is defined as $\triangle_{\mathcal{O}} q := \sup_{\{r,r':Adj(r,r')\}} \|q(r) - q(r')\|_{\mathcal{O}}$. In particular, for $\mathcal{O} = \mathbb{R}^k$ (where $k = +\infty$ is a possibility) equipped with the p-norm for $p \in [1, \infty]$, this defines the $\ell_p$-sensitivity, denoted $\triangle_p q$.*

In the following Theorem, we consider the $\mathcal{Q}$-function defined by $\mathcal{Q}(x) := \frac{1}{\sqrt{2\pi}} \int_x^{\infty} \exp(-\frac{u^2}{2}) du$ and we let $\kappa_{\delta,\epsilon} = \frac{1}{2\epsilon}(\mu + \sqrt{\mu^2 + 2\epsilon})$, with $\mu = Q^{-1}(\delta)$.

**Theorem 1.** *[20], [21] Let $q : \mathcal{H} \to \mathbb{R}^k$ be a query. The Gaussian mechanism $M_q$ defined by $M_q(r) = q(r) + \zeta$, with $\zeta \sim \mathcal{N}(0, \sigma_\zeta^2 I_k)$ and $\sigma_\zeta \geq \kappa_{\delta,\epsilon}(\Delta_2 q)$, is $(\epsilon, \delta)$-differentially private.*

When $q(r) = r$, this leads to the so-called input perturbation mechanism, where the raw participants' data is perturbed directly. An additional crucial property of differential privacy is its "resilience to post-processing", i.e., transforming a result that is differentially private does not weaken the guarantee, as long as the transformation does not involve the sensitive data [5], [21, Theorem 1].

Our goal is to design a statistical test for $H_0$ publishing a decision that is differentially private with respect to the adjacency relation (7) on the sequences $r = \{r_i\}_{i=1}^n \in \mathbb{R}^n$.

## III. DIFFERENTIALLY PRIVATE GAUSSIAN HYPOTHESIS TESTING

Given the form of the GLR test, by the resilience to post-processing property, basing the decision rule (6) on a differentially private version of $\tilde{l}(r)$ provides a differentially private test.

*Proposition 1.* A mechanism that publishes $\tilde{M}(r) = \frac{1}{n} \sum_{i=1}^n r_i + \zeta$ with $\zeta \sim \mathcal{N}(0, \sigma_\zeta^2)$ and $\sigma_\zeta \geq \kappa_{\delta,\epsilon} \frac{\rho}{n}$ is $(\epsilon, \delta)$-differentially private for the adjacency relation (7).

*Proof.* We have immediately, for two adjacent sequences $r$ and $r'$

$$\left| \frac{1}{n} \sum_{i=1}^n r_i - \frac{1}{n} \sum_{i=1}^n r_i' \right| \leq \frac{1}{n} \sum_{i=1}^n |r_i - r_i'| \leq \frac{\rho}{n}.$$

The averaging query being scalar in this case, its $\ell_p$ sensitivity with respect to (7) is the same for any $p$ and bounded by $\rho/n$. The result follows from Theorem 1. $\square$

Consider now the following statistical test for $H_0$. First, compute

$$\tilde{l}(r) = \frac{n}{2\sigma^2} \left[ \left( \frac{1}{n} \sum_{i=1}^n r_i \right) + \zeta_n \right]^2, \qquad (9)$$

with $\zeta_n \sim \mathcal{N}(0, \sigma_{\zeta_n}^2)$ and $\sigma_{\zeta_n} = \kappa_{\delta,\epsilon} \rho/n$. Then publish

$$\tilde{d}(r) = \begin{cases} 0 & \text{if } \tilde{l}(r) < \tilde{h} : H_0 \text{ is chosen} \\ 1 & \text{if } \tilde{l}(r) > \tilde{h} : H_1 \text{ is chosen,} \end{cases} \qquad (10)$$

where $\tilde{h}$ is an appropriately chosen threshold, set independently of the input data $r$.

The following Corollary is then a consequence of Proposition 1 and resilience to post-processing.

*Corollary* 1. The test (10) is $(\epsilon, \delta)$-differentially private.

We skip the proofs of our results in next sections due to space limitations.

## IV. PERFORMANCE ANALYSIS

We are now interested in characterizing the privacy-utility trade-off of the privacy-preserving GLR test (10). This relies on the following calculation.

*Proposition* 2. For the random variable $\tilde{l}(r)$ defined in (9), we have

$$\tilde{l}(r) = \left( \frac{1}{2} + \frac{\kappa_{\delta,\epsilon}^2 \rho^2}{2\sigma^2 n} \right) \Xi, \qquad (11)$$

with $\Xi \sim \chi_1^2 \left( \frac{\theta^2}{\frac{\sigma^2}{n} + \frac{\kappa_{\delta,\epsilon}^2 \rho^2}{n^2}} \right)$, where $\chi_1^2(\lambda)$ represents a non-central chi-squared distribution with one degree of freedom and noncentrality parameter $\lambda$.

From Proposition 2, we see in particular that the impact of the privacy-preserving noise decreases with the number $n$ of samples in the sequence $r$ as $O(1/n^2)$, i.e., whereas the impact of the intrinsic noise scales as $\sigma^2/n$. Hence, the effect of the privacy-preserving noise becomes negligible asymptotically as $n$ increases.

We need the following definitions for the next proposition. Let $\mathcal{F}(\cdot; k, \lambda^2)$ be tail distribution (complement of the cumulative distribution function) of the noncentral chi-squared distribution with $k$ degrees of freedom and non-centrality parameter $\lambda^2$ [18, Chapter 2]. Let $\mathcal{F}^{-1}(\cdot; k, \lambda^2)$ denote the inverse of $\mathcal{F}(\cdot; k, \lambda^2)$, defined on $[0, 1)$.

*Proposition* 3. The decision rule (10) achieves a probability $P_I$ of type I error when the threshold is chosen as

$$\tilde{h} = \left( \frac{1}{2} + \frac{\kappa_{\delta,\epsilon}^2 \rho^2}{2\sigma^2 n} \right) \mathcal{F}^{-1}(P_I; 1, 0). \qquad (12)$$

The corresponding probability of correct detection is

$$P_D = 1 - P_{II}$$

$$= \mathcal{F}\left( \mathcal{F}^{-1}(P_I; 1, 0); 1, \frac{\theta_1^2}{\frac{\sigma^2}{n} + \frac{\kappa_{\delta,\epsilon}^2 \rho^2}{n^2}} \right). \qquad (13)$$

Note again the vanishing impact of the privacy-preserving noise as $n$ increases in (13), the term $\kappa_{\delta,\epsilon}^2 \rho^2 / n^2$ becoming negligible with respect to $\sigma^2/n$.

We can compare the probability of detection as a function of the sequence length $n$ when using the differentially private test (10) and when using input perturbation. Recall that the latter consists in directly perturbing each element $r_i$ of the input sequence to $\tilde{r}_i = r_i + \zeta_i$, with $\{\zeta_i\}_{1 \le i \le n}$ a sequence of iid Gaussian random variables with standard deviation $\kappa_{\delta,\epsilon}\rho$, so that $\tilde{r}$ becomes a differentially private sequence for which we can now use the standard GLR decision rule by resilience to post-processing. For input perturbation, the probability of correct detection for a given value of $P_I$ reads

$$P_D = \mathcal{F}\left( \mathcal{F}^{-1}(P_I; 1, 0); 1, \frac{n\theta_1^2}{\sigma^2 + \kappa_{\delta,\epsilon}^2 \rho^2} \right)$$

instead of (13). Clearly, with input perturbation, the privacy-preserving noise plays the same role as the intrinsic noise from the point of view of the detector's performance, and so in this case the impact of the privacy requirement does not vanish as $n$ increases.

### A. Receiver Operating Characteristic Curves

A common performance metric for detection algorithms is the receiver operating characteristic (ROC) curve, which characterizes the trade-off between false alarms and true detections [18]. Consider a scenario when the scalar $n = 1000$, the parameter $\rho = 500$ and the standard deviation $\sigma = 0.5$. We obtain the ROC curves of Fig.1 by plotting the probability of detection $P_D$ against the probability of false-alarm $P_I$ using (13), for different values of the privacy parameter $\epsilon$ when we fix $\delta = 0.05$. The overall accuracy of the algorithm is higher when the curve is closer to the upper left corner. Hence, we notice the worst accuracy performance of the test (10) in the high-privacy regime, i.e., as $\epsilon$ becomes small, which shows the trade-off between the detector's accuracy and the privacy requirements.
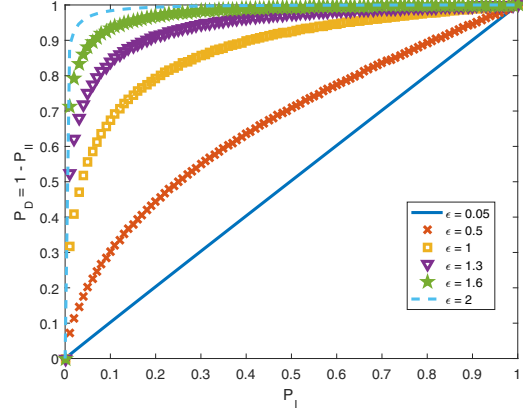


Fig. 1. ROC curve for different values of $\epsilon$ when $\delta = 0.05$.

## V. Application to Sequential Detection

Consider the following linear time-invariant system

$$x_{i+1} = Ax_i + Bu_i + Ef_i + w_i,$$
$$y_i = Cx_i + Du_i + Ff_i + v_i, \qquad (14)$$

where $x_i \in \mathbb{R}^n$ is the state vector, $u_i \in \mathbb{R}^m$ stands for a known input signal, $y_i \in \mathbb{R}$ represents the output signal available for measurement, and the process noise $w_i \sim \mathcal{N}(0, W)$ and measurement noise $v_i \sim \mathcal{N}(0, V)$ are independent sequences of iid zero-mean Gaussian random variables with covariances $W \succ 0$, $V > 0$. The initial condition $x_0$ is a Gaussian random vector independent of the noise processes $w_i$ and $v_i$. A monitoring system aims at detecting the vector of unknown input signals $f_i \in \mathbb{R}^{n_f}$, which represents additive faults or anomalies. A steady-state Kalman filter to estimate the state can be written

$$K = \Sigma C^{\mathrm{T}} (C \Sigma C^{\mathrm{T}} + V)^{-1},$$
$$\Sigma = A \Sigma A^{\mathrm{T}} + W - A \Sigma C^{\mathrm{T}} (C \Sigma C^{\mathrm{T}} + V)^{-1} C \Sigma A^{\mathrm{T}},$$
$$\hat{x}_i = \hat{x}_{i|i-1} + K(y_i - C \hat{x}_{i|i-1} - Du_i),$$
$$\hat{x}_{i+1|i} = A \hat{x}_i + Bu_i.$$

The innovation process is defined as follows

$$r_i = y_i - C \hat{x}_i - Du_i. \qquad (15)$$

In the absence of fault ($f_i = 0$) and in steady state, the innovation sequence is a white Gaussian process with $r_i \sim \mathcal{N}(0, \sigma^2)$, $\sigma^2 = C \Sigma C^{\mathrm{T}} + V > 0$ [22]. When a fault occurs, i.e. $f_i \ne 0$, the mean of $r_i$ is no longer equal to zero.

Consequently, given a privacy sensitive dataset consisting of the sequence $\{y_i\}$, we can apply the differentially private algorithm of Section III to the innovation process to discriminate between the hypotheses $f_i = 0$ and $f_i \ne 0$ in the model (14). For sequential detection however, one typically

wants to detect a change between these two regimes (the occurrence of a fault) quickly, while still controlling the rate of false alarms. One basic way of achieving this is to execute the test over the last $k$ samples only instead of the whole sequence, for a fixed value of $k$. At the end of the $p^{th}$ block of length $k$, one now sets

$$\tilde{l}_p(r) = \frac{k}{2\sigma^2} \left( \frac{1}{k} \sum_{i=k(p-1)+1}^{i=kp} r_i + \zeta_k \right)^2, \qquad (16)$$

and $n$ should be replaced by $k$ in the results of Section IV. An alarm is then raised at time $t_a$, which defines the following stopping rule

$$t_a = k\, p^* = k\, \min\{p : \tilde{d}_p(r) = 1\}. \qquad (17)$$

In words, an alarm is raised after the first sample of size $k$ for which the decision rule (10) chooses $H_1$. Increasing $k$ improves the accuracy of the detection, but has the negative impact of delaying the decision [19]. We describe the final differentially private sequential detection algorithm in Algorithm 1.

---

**Algorithm 1: Differentially private sequential detection**
**Given:**
1. A model (14) and a set of measurements data $\{y_i\}_{i=1}^{i=n}$.
2. A required target probability of type I error $P_I$.
3. Required differential privacy parameters $(\epsilon, \delta)$.
**Initialization at the $p$th sample:**
1. Select $k$.
2. Acquire $k$ data samples.
**At each $p$th sample:**
S1. Compute the residual sequence $\{r_i\}_{i=k(p-1)+1}^{i=kp}$ by using (15).
S3. Compute the ML estimate of the mean value $\hat{\theta}_1 = \frac{1}{k} \sum_{i=k(p-1)+1}^{i=kp} r_i$.
S4. Determine the amount of differential privacy noise $\zeta_k$ by using Corollary 1.
S5. Compute the decision function by using (12) with $n = k$.
S7. Take the decision by using (10).
**Result:** A sequence of decisions $\{\tilde{d}_p(r)\}_p$.

---

## VI. NUMERIC EXPERIMENTS

Let us apply the results of Section V to an epidemic outbreak detection scenario. Consider an abstract scenario in which the sequence $r$ consists of the data of the Centers for Disease Control and Prevention[1] (CDC) in the United States from the 15th week (from April 7 to April 14) to
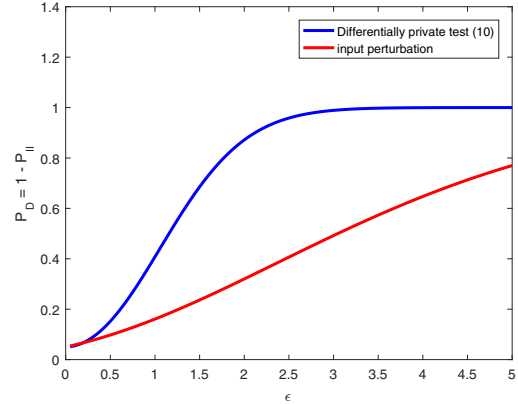
[1] https://www.cdc.gov/flu



Fig. 2. Probability of detection of the test (10), as a function of the privacy parameter $\epsilon$ when $k = 7$.

the 22th week (from June 3 to June 9) of 2018 [23]. In fact, the CDC collects data from U.S. clinical laboratories and hospital emergency departments, which consist of the number of positive tests reported to them for Influenza A(H1N1)pdm09, A(H3N2) and variant (H3N2v), and other Influenza A virus infections. We assume that the standard deviation $\sigma = 1$ and since the CDC's data consists of weekly reports, we set $k = 7$. Note that our settings may not be optimal (Fig. 3). We set the of probability of false-alarm $P_{FA} = 0.05$ and $\rho = 20$. The differentially private test (10) has better performance than the input perturbation mechanism (Fig. 2), as explained in Section IV, and the performance degrades as $\epsilon$ becomes smaller (high level of privacy). The alarm is raised at the first test time, i.e., at $t_a = 7$ days.

Assume now that the CDC executes the differentially private test (10) every 4 weeks instead of weekly, which means $k = 28$ (Fig. 3). As in the previous case, the differentially private test (10) shows better performance than the standard input perturbation mechanism (Fig. 4), but the probabilities of detection of these two mechanisms are much better when $k = 28$ than when $k = 7$. Indeed, the convergence of the probability of detection to 1 is much quicker when $k = 28$ than when $k = 7$, due to the increased sample size. However, the first opportunity to raise the alarm only occurs at $t_a = 28$ days. Accordingly, the monitoring system needs to trade-off detection accuracy and speed by appropriately selecting the sample size $k$.

## VII. CONCLUSION

This paper addresses the problem of hypothesis testing for the mean of a Gaussian sequence, under a differential privacy constraint. We study the case in which the data consist of iid Gaussian random variables and provide analytic
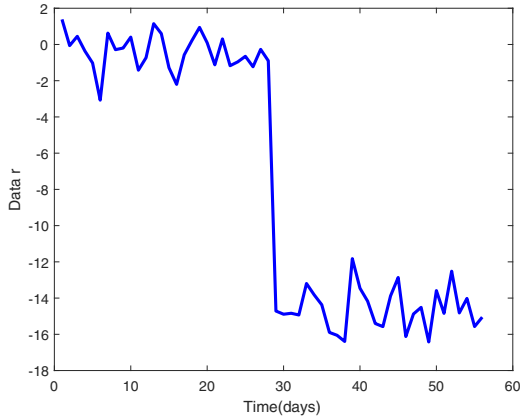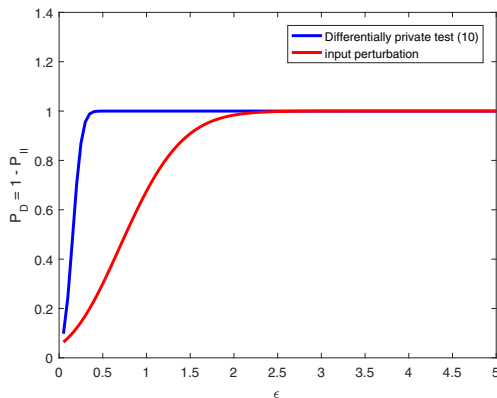
Fig. 3.    Data r when $k = 28$.



Fig. 4.    Probability of detection of the test (10), as a function of the privacy parameter $\epsilon$ when $k = 28$.

formulas to set the detection threshold and for the error probabilities. We apply the results to design a differentially private fault detector for linear dynamic processes with Gaussian noise, and propose a version of the algorithm when decisions are taken sequentially. We test the efficiency of the scheme on a dataset of the CDC for the problem of detecting an epidemic outbreak.

## REFERENCES

[1] L. Sweeney, "*k*-anonymity: a model for protecting privacy," *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, vol. 10, pp. 557–570, 2002.

[2] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam, "*l*-diversity: Privacy beyond *k*-anonymity," *ACM Trans. Knowl. Discov. Data*, vol. 1, no. 1, March 2007.

[3] C. Dwork, "Differential privacy," in *Proceedings of the 33rd International Colloquium on Automata, Languages and Programming (ICALP)*, ser. Lecture Notes in Computer Science, vol. 4052. Springer-Verlag, 2006.

[4] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in *Proceedings of the Third Theory of Cryptography Conference*, 2006, pp. 265–284.

[5] C. Dwork and A. Roth, *The Algorithmic Foundations of Differential Privacy*, t. e. o. k. now, Ed.   Foundations and Trends in Theoretical Computer Science, 2014, vol. 9, no. 3-4.

[6] K. H. Degue and J. Le Ny, "On differentially private Kalman filtering," in *Proceedings of the 5th IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, Montreal, Canada, Nov. 2017.

[7] D. Vu and A. Slavković, "Differential privacy for clinical trial data: Preliminary evaluations," in *Proceedings of the 2009 IEEE International Conference on Data Mining Workshops*, Miami, FL, USA, Dec. 2009.

[8] C. Dwork, W. Suy, and L. Zhang, "Private false discovery rate control," *arXiv preprint arXiv:1511.03803*, Nov. 2015.

[9] M. Gaboardi, H. W. Lim, R. Rogers, and S. P. Vadhan, "Differentially Private Chi-Squared Hypothesis Testing: Goodness of Fit and Independence Testing," in *Proceedings of the 33rd International Conference on Machine Learning*, New York City, NY, USA, Jun. 2016, pp. 2111–2120.

[10] C. Uhler, A. Slavković, and S. E. Fienberg, "Privacy-Preserving Data Sharing for Genome-Wide Association Studies," *Journal of Privacy and Confidentiality*, vol. 5, no. 1, pp. 137–166, 2013.

[11] S. Simmons and B. Berger, "Realizing privacy preserving genome-wide association studies," *Bioinformatics*, vol. 32, no. 9, pp. 1293–1300, Jan. 2016.

[12] R. Rogers and D. Kifer, "A new class of private chi-square hypothesis tests," in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, Fort Lauderdale, Florida, USA, Apr. 2017.

[13] Y. Wang, J. Lee, and D. Kifer, "Revisiting differentially private hypothesis tests for categorical data," *ArXiv e-prints*, Mar. 2017.

[14] M. Ghassemi, A. D. Sarwate, and R. Wright, "Differentially private online active learning with applications to anomaly detection," in *Proceedings of the 9th ACM Workshop on Artificial Intelligence and Security*, October 2016.

[15] X. Tong, B. Xi, M. Kantarcioglu, and A. Inan, "Gaussian mixture models for classification and hypothesis tests under differential privacy," in *31st Annual IFIP WG 11.3 Conference on Data and Applications Security and Privacy (DBSec'17)*, Philadelphia, PA, USA, Jul. 2017.

[16] L. Fan and L. Xiong, "Differentially private anomaly detection with a case study on epidemic outbreak detection," in *Proceedings of the 13th International Conference on Data Mining Workshops*, Dallas, TX, USA, Dec. 2013.

[17] J. S. Oakland, *Statistical Process Control*, 6th ed.    Oxford: Butterworth-Heinemann, 2008.

[18] S. M. Kay, *Fundamentals of Statistical Processing, Volume 2: Detection Theory*, ser. Prentice Hall Signal Processing Series.  Upper Saddle River: Prentice-Hall PTR, 2009.

[19] M. Basseville and I. V. Nikiforov, *Detection of Abrupt Changes: Theory and Application*.    Englewood Cliffs, NJ, USA: Prentice-Hall, 1993.

[20] C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor, "Our data, ourselves: Privacy via distributed noise generation," *Advances in Cryptology-EUROCRYPT*, vol. 4004, pp. 486–503, 2006.

[21] J. Le Ny and G. Pappas, "Differential private filtering," *IEEE Transactions on Automatic Control*, vol. 59, no. 2, pp. 341–354, February 2014.

[22] S. X. Ding, *Model-based Fault Diagnosis Techniques. Design Schemes, Algorithms and Tools*.  Berlin: Springer Heidelberg, 2008.

[23] "Weekly U.S. Influenza surveillance report," Centers for Disease Control and Prevention (CDC), Atlanta, GA, USA, Tech. Rep., Jul. 2018.