

# LDP two-sample chi-squared test

## 1 Setting

- $\mathbf{Y}_i \stackrel{iid}{\sim} \text{multi}(n_1, \mathbf{p}_Y), \mathbf{Z}_i \stackrel{iid}{\sim} \text{multi}(n_2, \mathbf{p}_Z)$  with  $k$  categories
- One-hot vector form i.e. random vectors with dependent Bernoulli random variable entries
- Allow for  $n_1 \neq n_2$

## 2 Generalized Randomized Response and two sample Pearson chi-square statistic

### 2.1 Privacy mechanism: Generalized Randomized Response

**Definition 2.1** (Generalized Randomized Response (Theorem 5.4. of Gaboardi and Rogers [1])).  
For a multinomial random vector  $\mathbf{Y}_i \stackrel{iid}{\sim} \text{multi}(n_1, \mathbf{p}_Y)$ , we define

$$\mathbb{P}(\mathcal{M}_{\text{GenRR}}(\mathbf{Y}_i) = \mathbf{y}' | \mathbf{Y}_i = \mathbf{y}) := \begin{cases} \frac{\exp(\alpha)}{\exp(\alpha) + k - 1} & \text{if } \mathbf{y}' = \mathbf{y} \\ \frac{1}{\exp(\alpha) + k - 1} & \text{if } \mathbf{y}' \neq \mathbf{y}. \end{cases}$$

Then  $\tilde{\mathbf{Y}}_i := \mathcal{M}_{\text{GenRR}}(\mathbf{Y}_i)$  is a multinomial random vector with probability vector

$$\tilde{\mathbf{p}}_Y := \mathbf{p}_Y \frac{\exp(\alpha)}{\exp(\alpha) + k - 1} + (1 - \mathbf{p}_Y) \frac{1}{\exp(\alpha) + k - 1}.$$

Since  $e^\alpha > 1$  for  $\alpha > 0$ , the probability of sending the original category is a little bit higher than sending the other category. Gaboardi and Rogers [1] constructs a private goodness-of-fit test based on a chi-square statistic evaluated on  $\tilde{\mathbf{Y}}_i$ 's. They demonstrate that the limiting distribution is chi-square distribution both under the null and alternative.

### 2.2 Two sample chi-square statistic

We extend the goodness-of-fit test by Gaboardi and Rogers [1] into two-sample testing by privatizing the raw samples  $\mathbf{Z}_i \stackrel{iid}{\sim} \text{multi}(n_2, \mathbf{p}_Z)$  into  $\tilde{\mathbf{Z}}_j := \mathcal{M}_{\text{GenRR}}(\mathbf{Z}_j)$ . Under the null,  $\mathcal{M}_{\text{GenRR}}(\mathbf{Y}_i)$  and  $\mathcal{M}_{\text{GenRR}}(\mathbf{Z}_j)$  follow multinomial distributions with the same probability vector. Therefore, the usual two-sample chi-square test statistic

$$T_\chi := \sum_{\ell=1}^k \frac{(n_2 \sum_{i=1}^{n_1} \tilde{\mathbf{Y}}_i(\ell) - n_1 \sum_{j=1}^{n_2} \tilde{\mathbf{Z}}_j(\ell))^2}{n_1 n_2 (n_1 + n_2) \sum_{j=1}^{n_2} (\tilde{\mathbf{Y}}_j(\ell) + \tilde{\mathbf{Z}}_j(\ell))}$$

converges to a chi-square distribution with degree of freedom  $k - 1$  and yields a valid test with size  $\gamma$ . This test statistic is from Van der Vaart's book Asymptotic Statistics, pp. 253.

### 3 Bit flip privatization and related test statisitc

#### 3.1 Bit flip privatization

We next consider another LDP algorithm  $\mathcal{M}_{bit} : \{\mathbf{e}_1, \dots, \mathbf{e}_k\} \rightarrow \{0, 1\}^k$ , which is the Algorithm 4 of Gaboardi and Rogers [1]. It flips each bit with some biased probability. The Algorithm 1 is

---

**Algorithm 1** Bit Flip Local Randomizer:  $\mathcal{M}_{bit}$

---

**Input:**  $\mathbf{y} \in \{\mathbf{e}_1, \dots, \mathbf{e}_k\}$ ,  $\alpha$ .

**for**  $\ell \in [k]$  **do**

    Set  $\tilde{\mathbf{y}}_\ell = \mathbf{y}_\ell$  with probability  $e^{\alpha/2}/(e^{\alpha/2} + 1)$ , otherwise  $\tilde{\mathbf{y}}_\ell = 1 - \mathbf{y}_\ell$

**end for**

**Output:**  $\tilde{\mathbf{y}}$

---

$\alpha$ -LDP (Theorem 5.5 of Gaboardi and Rogers [1]).)

#### 3.2 test statistic

We first review the one-sample test statistic of Gaboardi and Rogers [1] and expand it into two-sample statistic.

##### 3.2.1 Review of one-sample statistic

We first review how Gaboardi and Rogers [1] builds goodness-of-fit chi-square statistic for histogram of bit-flipped observations. We start with applying CLT to the bit-flipped obseervations.

**Lemma 3.1** (Applying CLT to bit-flipped observations, Lemma 5.7 of [1]). *When  $Y_i \stackrel{iid}{\sim} \text{multinomial}(\mathbf{p}, 1)$ , the mean vector and covariance matrix of the flipped observation are computed as follows:*

$$\tilde{\mathbf{p}} := \mathbb{E}(\mathcal{M}_{bit}(Y_1)) = \frac{(\exp(\alpha/2) - 1)\mathbf{p} + \mathbf{1}}{\exp(\alpha/2) + 1}, \text{ and} \quad (1)$$

$$\Sigma_{\tilde{\mathbf{p}}} := \text{Var}(\mathcal{M}_{bit}(Y_1)) = \left( \frac{\exp(\alpha/2) - 1}{\exp(\alpha/2) + 1} \right)^2 (\text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top) + \frac{\exp(\alpha/2)}{(\exp(\alpha/2) + 1)^2} I_d, \quad (2)$$

For any  $\alpha > 0$  and  $\mathbf{p} > 0$ ,  $\Sigma_{\tilde{\mathbf{p}}}$  is positive definite and one of its eigenvector is one-vector. Denote the histogram of flipped observations as

$$\tilde{\mathbf{H}} := \sum_{i=1}^n \mathcal{M}_{bit}(Y_i). \quad (3)$$

By the CLT for i.i.d random vectors, we get the following asymptotic distribution:

$$\sqrt{n}(\tilde{\mathbf{H}}/n - \tilde{\mathbf{p}}) \xrightarrow{d} N(0, \Sigma_{\tilde{\mathbf{p}}}) \quad (4)$$

In non-private chi-square test, we apply CLT and multiply by  $\text{diag}(\mathbf{p})^{-1/2}$  to turn the covariance matrix on the RHS into a projector matrix. Here in the private setting, we also need a scaling matrix to turn the covariance matrix into a projector matrix. Gaboardi and Rogers [1] proposes  $\tilde{\mathbf{p}}^{-1/2}\Pi$ , where  $\Pi := I_k - \frac{1}{k}\mathbf{1}\mathbf{1}^\top$ . The properties of  $\Pi$  are as follows:

1. It is symmetric idempotent (a projector matrix).
2. Its null space is  $\text{span}\{\mathbf{1}\}$ , so when multiplied to a symmetric matrix, it deletes an eigenvector  $\mathbf{1}$ .

$$\begin{aligned}\Pi x = 0 &\iff x = (1/k)\mathbf{1}\mathbf{1}^\top x \\ &\iff x = (1/k)\mathbf{1}(\mathbf{1}^\top x) = ((\mathbf{1}^\top x)/k)\mathbf{1} = c\mathbf{1}\end{aligned}$$

By multiplying  $\Sigma_{\tilde{\mathbf{p}}}^{-1/2}\Pi$  to the LHS vector of (4), we get

$$\sqrt{n}\Sigma_{\tilde{\mathbf{p}}}^{-1/2}\Pi(\tilde{\mathbf{H}}/n - \tilde{\mathbf{p}}) \xrightarrow{d} N(0, \Sigma_{\tilde{\mathbf{p}}}^{-1/2}\Pi\Sigma_{\mathbf{p}}\Pi\Sigma_{\tilde{\mathbf{p}}}^{-1/2}). \quad (5)$$

The next lemma specifies the property of the covariance matrix  $\Sigma_{\tilde{\mathbf{p}}}^{-1/2}\Pi\Sigma_{\mathbf{p}}\Pi\Sigma_{\tilde{\mathbf{p}}}^{-1/2}$ .

**Lemma 3.2.** *Let  $\Sigma \in \mathbb{R}^{k \times k}$  be a symmetric positive definite matrix one of whose eigenvector is  $\mathbf{1}$ . Then we can diagonalize as  $\Sigma = BDB^\top$ . Let  $\Pi := I_k - \frac{1}{k}\mathbf{1}\mathbf{1}^\top$ . Then the following matrix is the identity matrix except one of the entries on the diagonal is zero.*

$$\Sigma^{-1/2}\Pi\Sigma\Pi\Sigma^{-1/2} = \Sigma^{-1/2}\Pi BDB^\top\Pi\Sigma^{-1/2} \quad (6)$$

Now we invoke the following classical theorem to derive an asymptotic chi-square distribution with degree of freedom  $k - 1$ .

**Theorem 3.1** (Ferguson (1996)). *If  $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$  and  $\Sigma$  is a projection matrix of rank  $\nu$  and  $\Sigma\boldsymbol{\mu} = \boldsymbol{\mu}$  then  $\mathbf{X}^\top \mathbf{X} \sim \chi_\nu^2(\boldsymbol{\mu}^\top \boldsymbol{\mu})$ .*

### 3.3 Extension to two sample test statistic

We extend the previous result to two-sample setting. For simplicity, we follow the equal sample size setting of Gaboardi and Rogers [1]. Now we have two collections of raw data  $\{\mathbf{Y}_i\}_{i \in [n]}$  and  $\{\mathbf{Z}_j\}_{j \in [n]}$  generated from multinomial distributions with probability vectors  $\mathbf{p}_1$  and  $\mathbf{p}_2$ , respectively. According to Lemma 3.1, flipped observations of these samples have following moments:

$$\begin{aligned}\tilde{\mathbf{p}}_Y &:= \mathbb{E}(\mathcal{M}_{bit}(Y_1)) = \frac{(\exp(\alpha/2) - 1)\mathbf{p}_Y + 1}{\exp(\alpha/2) + 1}, \\ \Sigma_{\tilde{\mathbf{p}}_Y} &:= \text{Var}(\mathcal{M}_{bit}(Y_1)) = \left( \frac{\exp(\alpha/2) - 1}{\exp(\alpha/2) + 1} \right)^2 (\text{diag}(\mathbf{p}_Y) - \mathbf{p}_Y\mathbf{p}_Y^\top) + \frac{\exp(\alpha/2)}{(\exp(\alpha/2) + 1)^2} I_d \\ \tilde{\mathbf{p}}_Z &:= \mathbb{E}(\mathcal{M}_{bit}(Z_1)) = \frac{(\exp(\alpha/2) - 1)\mathbf{p}_Z + 1}{\exp(\alpha/2) + 1}, \\ \Sigma_{\tilde{\mathbf{p}}_Z} &:= \text{Var}(\mathcal{M}_{bit}(Z_1)) = \left( \frac{\exp(\alpha/2) - 1}{\exp(\alpha/2) + 1} \right)^2 (\text{diag}(\mathbf{p}_Z) - \mathbf{p}_Z\mathbf{p}_Z^\top) + \frac{\exp(\alpha/2)}{(\exp(\alpha/2) + 1)^2} I_d\end{aligned}$$

We also denote the histograms of flipped observations as

$$\tilde{\mathbf{H}}_{\mathbf{Y}} := \sum_{i=1}^n \mathcal{M}_{bit}(Y_i) \text{ and } \tilde{\mathbf{H}}_{\mathbf{Z}} := \sum_{i=1}^n \mathcal{M}_{bit}(Z_i). \quad (7)$$

Then we have the following asymptotic distribution:

$$\sqrt{n} \left( \left( \frac{\tilde{\mathbf{H}}_{\mathbf{Y}}}{n} - \frac{\tilde{\mathbf{H}}_{\mathbf{Z}}}{n} \right) - (\tilde{\mathbf{p}}_{\mathbf{Y}} - \tilde{\mathbf{p}}_{\mathbf{Z}}) \right) \xrightarrow{d} N(0, \Sigma_{\tilde{\mathbf{p}}_{\mathbf{Y}}} + \Sigma_{\tilde{\mathbf{p}}_{\mathbf{Z}}}). \quad (8)$$

According to Lemma 3.1,  $\Sigma_{\tilde{\mathbf{p}}_{\mathbf{Y}}}$  and  $\Sigma_{\tilde{\mathbf{p}}_{\mathbf{Z}}}$  are positive-definite. Since the set of symmetric positive-definite matrices is closed under nonnegative linear combination,  $\Sigma_{\tilde{\mathbf{p}}_{\mathbf{Y}}} + \Sigma_{\tilde{\mathbf{p}}_{\mathbf{Z}}}$  is also symmetric positive definite. Lemma 3.1 also implies that both of  $\Sigma_{\tilde{\mathbf{p}}_{\mathbf{Y}}}$  and  $\Sigma_{\tilde{\mathbf{p}}_{\mathbf{Z}}}$  have eigenvector  $\mathbf{1}$ . Therefore,  $\Sigma_{\tilde{\mathbf{p}}_{\mathbf{Y}}} + \Sigma_{\tilde{\mathbf{p}}_{\mathbf{Z}}}$  also has eigenvector  $\mathbf{1}$ . Thus we can invoke Lemma 3.2 to modify the asymptotic distribution (8). Let us denote  $\tilde{\Sigma} := \Sigma_{\tilde{\mathbf{p}}_{\mathbf{Y}}} + \Sigma_{\tilde{\mathbf{p}}_{\mathbf{Z}}}$ . Then we have the following asymptotic distribution:

$$\sqrt{n} \tilde{\Sigma}^{-1/2} \Pi \left( \left( \frac{\tilde{\mathbf{H}}_{\mathbf{Y}}}{n} - \frac{\tilde{\mathbf{H}}_{\mathbf{Z}}}{n} \right) - (\tilde{\mathbf{p}}_{\mathbf{Y}} - \tilde{\mathbf{p}}_{\mathbf{Z}}) \right) \xrightarrow{d} N(0, \tilde{\Sigma}^{-1/2} \Pi \tilde{\Sigma} \Pi \tilde{\Sigma}^{-1/2}). \quad (9)$$

Since  $\|\mathcal{M}_{bit}(Y_i)\|_2 \leq \sqrt{k}$  and  $\|\mathcal{M}_{bit}(Z_j)\|_2 \leq \sqrt{k}$ , the sample covariance matrices

$$\begin{aligned} \hat{\Sigma}_{\tilde{\mathbf{p}}_{\mathbf{Y}}} &:= \frac{1}{n} \sum_{i=1}^n (\mathcal{M}_{bit}(Y_i) - \tilde{\mathbf{H}}_{\mathbf{Y}}/n)(\mathcal{M}_{bit}(Y_i) - \tilde{\mathbf{H}}_{\mathbf{Y}}/n)^\top \\ \hat{\Sigma}_{\tilde{\mathbf{p}}_{\mathbf{Z}}} &:= \frac{1}{n} \sum_{j=1}^n (\mathcal{M}_{bit}(Z_j) - \tilde{\mathbf{H}}_{\mathbf{Z}}/n)(\mathcal{M}_{bit}(Z_j) - \tilde{\mathbf{H}}_{\mathbf{Z}}/n)^\top \end{aligned}$$

converge in probability to  $\Sigma_{\tilde{\mathbf{p}}_{\mathbf{Y}}}$  and  $\Sigma_{\tilde{\mathbf{p}}_{\mathbf{Z}}}$ , respectively (Corollary 6.20 of Wainwright [2]). Let us denote  $\hat{\Sigma} := \hat{\Sigma}_{\tilde{\mathbf{p}}_{\mathbf{Y}}} + \hat{\Sigma}_{\tilde{\mathbf{p}}_{\mathbf{Z}}}$ . Since matrix inversion and matrix square root is continuous mapping on the space of positive symmetric definite matrices, we have

$$\hat{\Sigma}^{-1/2} \tilde{\Sigma}^{-1/2} \xrightarrow{p} I_k. \quad (10)$$

Therefore, by the Slutsky's theorem, we have

$$\sqrt{n} \hat{\Sigma}^{-1/2} \Pi \left( \left( \frac{\tilde{\mathbf{H}}_{\mathbf{Y}}}{n} - \frac{\tilde{\mathbf{H}}_{\mathbf{Z}}}{n} \right) - (\tilde{\mathbf{p}}_{\mathbf{Y}} - \tilde{\mathbf{p}}_{\mathbf{Z}}) \right) \xrightarrow{d} N(0, \tilde{\Sigma}^{-1/2} \Pi \tilde{\Sigma} \Pi \tilde{\Sigma}^{-1/2}). \quad (11)$$

Under the Under the null hypothesis of  $\mathbf{p}_{\mathbf{Y}} = \mathbf{p}_{\mathbf{Z}}$ , we have  $\tilde{\mathbf{p}}_{\mathbf{Y}} - \tilde{\mathbf{p}}_{\mathbf{Z}} = 0$ . Therefore, we have

$$\sqrt{n} \hat{\Sigma}^{-1/2} \Pi \left( \frac{\tilde{\mathbf{H}}_{\mathbf{Y}}}{n} - \frac{\tilde{\mathbf{H}}_{\mathbf{Z}}}{n} \right) \xrightarrow{d} N(0, \tilde{\Sigma}^{-1/2} \Pi \tilde{\Sigma} \Pi \tilde{\Sigma}^{-1/2}). \quad (12)$$

Finally, we invoke Theorem 3.1 to obtain the following asymptotic null distribution

$$n \left( \frac{\tilde{\mathbf{H}}_{\mathbf{Y}}}{n} - \frac{\tilde{\mathbf{H}}_{\mathbf{Z}}}{n} \right)^\top \Pi \hat{\Sigma}^{-1} \Pi \left( \frac{\tilde{\mathbf{H}}_{\mathbf{Y}}}{n} - \frac{\tilde{\mathbf{H}}_{\mathbf{Z}}}{n} \right) \xrightarrow{d} \chi_{(k-1)}^2, \quad (13)$$

and we define the lefthand side of (13) as our test statistic.

## References

- [1] Gaboardi, M. and Rogers, R. (2018). Local private hypothesis testing: Chi-square tests. *Proceedings of the 35th International Conference on Machine Learning*, 80:1626–1635.
- [2] Wainwright, M. J. (2019). *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge.