



**15<sup>TH</sup> NATIONAL  
CONVENTION  
ON STATISTICS**

03-05 OCTOBER 2022

*Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority*



# **Mathematically Rigorous Approach to Anonymize the CBMS Micro-data using Differential Privacy**

**A. E. Abdulsamad**

Supervising Statistical Specialist  
Philippine Statistics Authority

**Statistical Quality Assurance and Data Privacy, Confidentiality, and Protection**

Crowne Plaza Galleria Manila  
08:30-10:00AM, 05 October 2022



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## Motivation

“

*Please be informed that all information shared are strictly confidential pursuant to Section 10 (Confidentiality of Information) of Republic Act (RA) No. 11315 or the CBMS Act and Section 8 (Confidentiality) of RA No. 10173 or the Data Privacy Act of 2012 and will not be used against you or to any of your household members for taxation, investigation, or law enforcement purposes.*

”

Confidentiality guarantee in the 2021 Pilot CBMS Household Profile Questionnaire



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

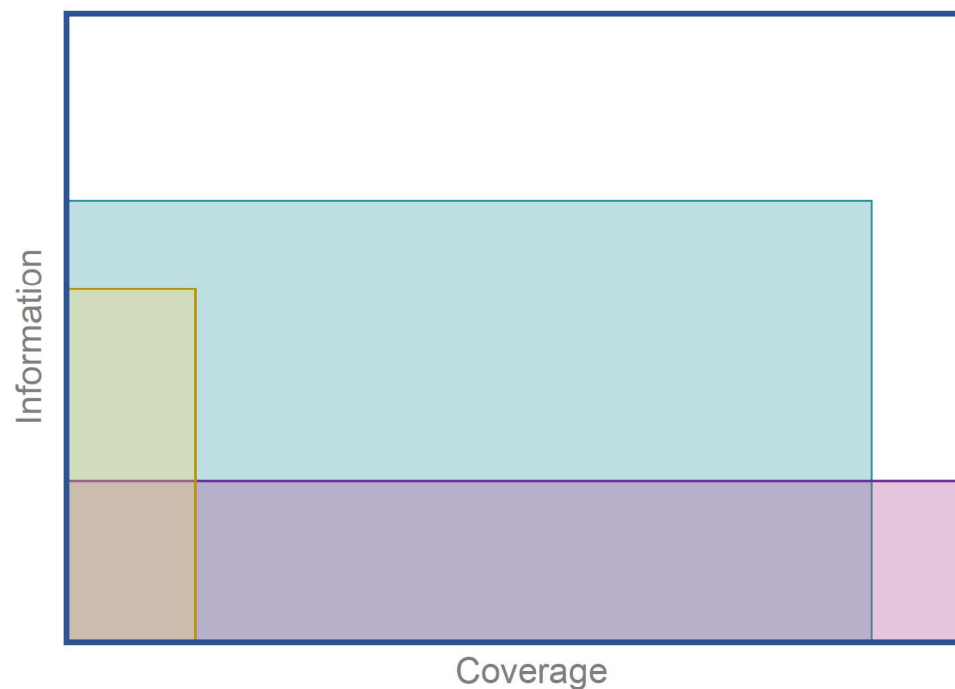
Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## Motivation

Why **data privacy** is  
crucial for **CBMS**?

- Sample surveys
- Census of population
- CBMS





# **15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS**

03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## **What we usually do to protect the privacy of our respondents**

### **De-identification**

Anonymize by removing personally identifiable information

### **Geo-swapping**

Randomly swap the data subject's geographic information

### **Release only summary or aggregate statistics**

Never release the microdata



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## Privacy guarantee

Aggregate statistics

Microdata



Privacy relatively guaranteed

Privacy is at risk



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

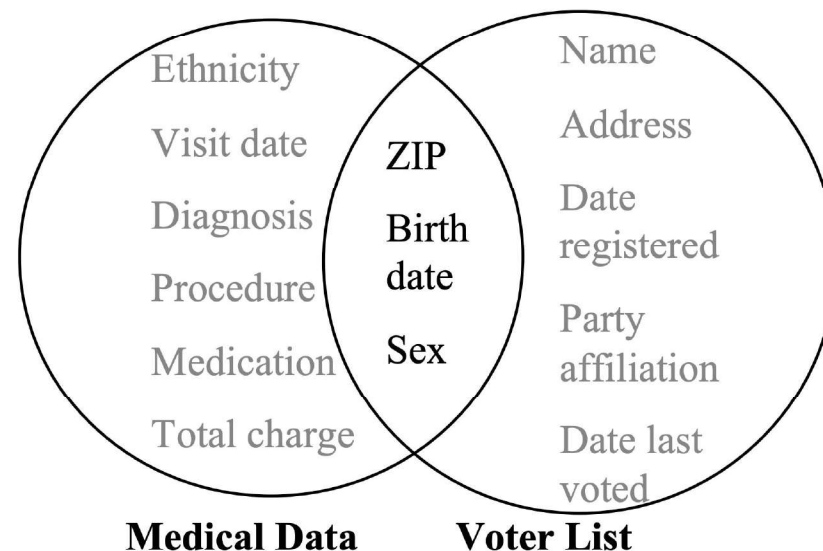
03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## Linkage attack

Re-identification of individuals by directly matching or linking two (2) or more datasets with shared attributes



Linking to re-identify data (Sweeney, 2002)



## Reconstruction attack

- The U.S. Census Bureau reconstructed 100% of the 2010 Census micro-data records (308,745,538 persons)
- The reconstructed records matched the confidential data (2010 CEF) exactly (every single bit) for 46% of the population (142 million people) and allowing age +/- 1 year for 71% of the population (219 million people).

Reconstruction-abetted re-identification attacks and other traditional vulnerabilities (Abowd, n.d.)



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

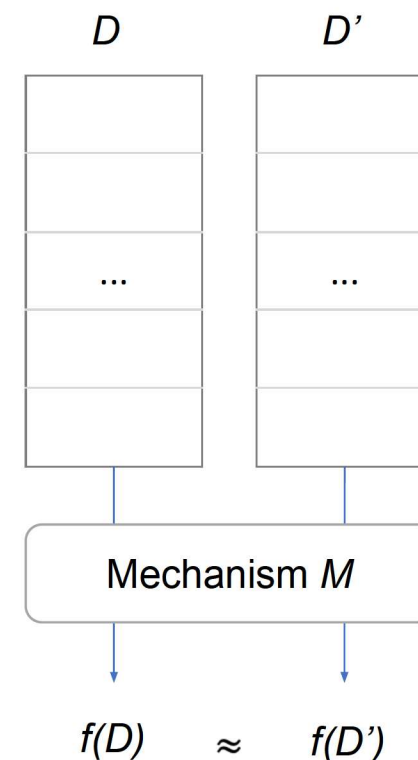
03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## Differential Privacy

- Gives a promise of a quantifiable and provable privacy guarantee at the same time preserve statistical properties of the original data
- Ensures that the same conclusions will be reached, independent of whether any individual opts in or opts out of the dataset.







# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## Differential Privacy

*Formal Definition:* A randomized mechanism  $M$  gives  $(\epsilon, \delta)$ -differential privacy for every set of outputs  $S$ , and for any neighbouring datasets of  $D$  and  $D'$ , if  $M$  satisfies:

$$P [M(D) \in S] \leq e^\epsilon \cdot P [M(D') \in S] + \delta.$$



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## Differential Privacy

- The parameter  $\epsilon$  is defined as the privacy-loss budget which controls the privacy guarantee level of mechanism  $M$ .
- For a particular output, the ratio on two probabilities is bounded by  $e^\epsilon$



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



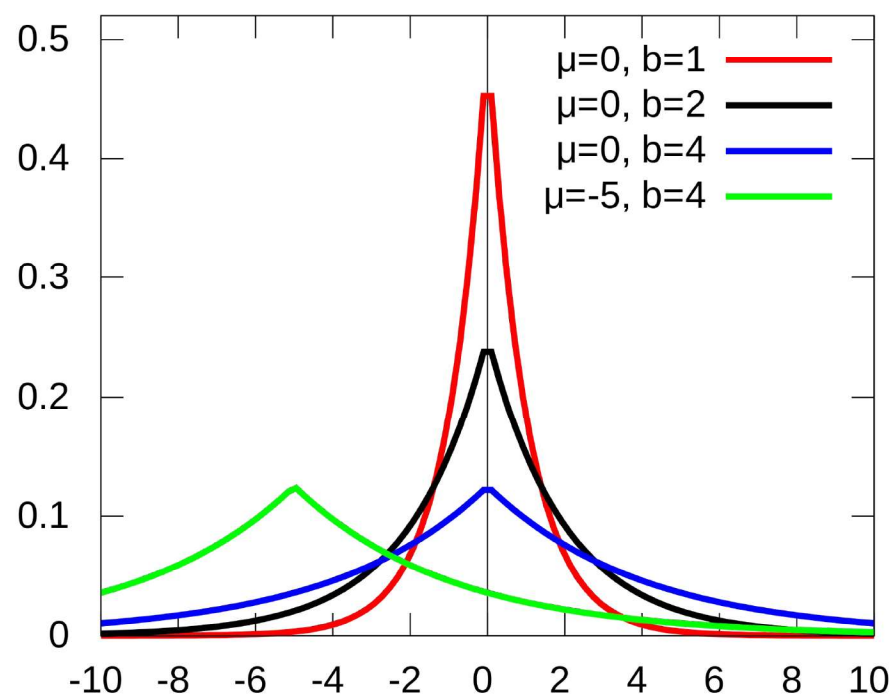
## Laplace mechanisms

$$M(D) = f(D) + \text{Lap}\left(\frac{\Delta f}{\epsilon}\right)$$

$$\Delta f = \max ||f(D) - f(D')||_1$$

### Other mechanisms:

- Gaussian
- Exponential



# Properties of Differential Privacy

## Protection against arbitrary risks

Moving beyond protection against re-identification

## Automatic neutralization of linkage attacks

Including all those attempted with all past, present, and future datasets and other forms and sources of auxiliary information

## Quantification of privacy loss

Decide the level of privacy and accuracy of the output

# Properties of Differential Privacy

## Composition

Analysis and control of cumulative privacy loss over multiple computations

## Group privacy

Analysis and control of privacy loss incurred by groups, such as families

## Closure under post-processing

Data analyst cannot increase privacy loss





# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

## Original Dataset $D$

Histogram query  
using the Laplace  
mechanism

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



Age Group	Count	$\epsilon = 0.01$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$	$\epsilon = 0.5$	$\epsilon = 1.0$	$\epsilon = 5.0$
< 5	2,193	2,138	2,198	2,192	2,191	2,189	2,196	2,192	2,193
5-9	2,710	2,796	2,713	2,708	2,713	2,710	2,714	2,708	2,710
10-14	2,739	2,719	2,784	2,733	2,739	2,748	2,739	2,738	2,739
15-19	2,596	2,741	2,600	2,590	2,599	2,600	2,596	2,596	2,596
20-24	2,440	2,653	2,445	2,436	2,443	2,444	2,438	2,440	2,440
25-29	2,153	1,913	2,154	2,153	2,156	2,150	2,157	2,155	2,153
30-34	1,870	1,876	1,872	1,867	1,870	1,867	1,869	1,870	1,870
35-39	1,659	1,813	1,654	1,665	1,661	1,660	1,655	1,660	1,659
40-44	1,344	1,355	1,332	1,332	1,345	1,343	1,349	1,344	1,344
45-49	1,195	1,186	1,221	1,194	1,197	1,196	1,196	1,195	1,195
50-54	1,118	1,363	1,134	1,123	1,096	1,117	1,115	1,119	1,118
55-59	1,045	1,035	1,050	1,038	1,045	1,043	1,045	1,046	1,045
60-64	801	845	810	802	798	803	806	804	801
65-69	665	681	635	661	664	661	665	665	665
70-74	411	253	411	404	412	411	411	411	411
75-79	250	411	257	254	249	250	251	250	250
80-84	130	59	122	138	125	131	129	126	130
85-89	74	0	69	73	72	73	73	73	74
90-94	14	0	6	16	15	14	12	15	14
95-99	1	241	0	0	0	7	0	0	1
100+	1	152	0	0	4	1	8	0	1
<b>Total</b>	25,409	26,230	25,467	25,379	25,394	25,418	25,424	25,407	25,409



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

## Neighboring Dataset $D'$

Histogram query  
using the Laplace  
mechanism

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



Age Group	Count	$\epsilon = 0.01$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$	$\epsilon = 0.5$	$\epsilon = 1.0$	$\epsilon = 5.0$
< 5	2,193	2,174	2,209	2,193	2,192	2,194	2,191	2,194	2,193
5-9	2,710	2,647	2,704	2,707	2,710	2,708	2,712	2,711	2,710
10-14	2,739	2,769	2,735	2,735	2,732	2,739	2,740	2,742	2,739
15-19	2,596	2,496	2,632	2,598	2,594	2,598	2,602	2,594	2,596
20-24	2,440	2,570	2,443	2,428	2,439	2,437	2,440	2,438	2,440
25-29	2,153	2,221	2,174	2,156	2,150	2,152	2,149	2,154	2,153
30-34	1,870	1,911	1,869	1,868	1,874	1,870	1,871	1,871	1,870
35-39	1,659	1,686	1,657	1,658	1,655	1,662	1,660	1,655	1,659
40-44	1,344	1,314	1,348	1,349	1,347	1,349	1,341	1,345	1,344
45-49	1,194	1,200	1,182	1,203	1,194	1,198	1,194	1,195	1,194
50-54	1,118	1,256	1,120	1,115	1,119	1,119	1,116	1,118	1,117
55-59	1,045	1,063	1,038	1,058	1,041	1,051	1,041	1,045	1,046
60-64	801	915	827	804	802	801	801	800	800
65-69	665	618	667	667	663	665	661	661	665
70-74	411	465	411	400	413	410	409	411	411
75-79	250	187	248	249	249	249	246	250	250
80-84	130	151	144	130	139	122	131	130	130
85-89	74	70	71	75	83	74	73	74	74
90-94	14	0	8	17	16	17	11	15	14
95-99	1	15	0	10	0	0	0	0	1
100+	1	176	0	2	0	0	4	1	1
Total	25,408	25,904	25,487	25,422	25,412	25,415	25,393	25,404	25,407



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

## Original Dataset $D$

Histogram query  
using the **Laplace**  
mechanism

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



Age Group	$\epsilon = 0.01$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$	$\epsilon = 0.5$	$\epsilon = 1.0$	$\epsilon = 5.0$
< 5	55	-5	1	2	4	-3	1	0
5-9	-86	-3	2	-3	0	-4	2	0
10-14	20	-45	6	0	-9	0	1	0
15-19	-145	-4	6	-3	-4	0	0	0
20-24	-213	-5	4	-3	-4	2	0	0
25-29	240	-1	0	-3	3	-4	-2	0
30-34	-6	-2	3	0	3	1	0	0
35-39	-154	5	-6	-2	-1	4	-1	0
40-44	-11	12	12	-1	1	-5	0	0
45-49	9	-26	1	-2	-1	-1	0	0
50-54	-245	-16	-5	22	1	3	-1	0
55-59	10	-5	7	0	2	0	-1	0
60-64	-44	-9	-1	3	-2	-5	-3	0
65-69	-16	30	4	1	4	0	0	0
70-74	158	0	7	-1	0	0	0	0
75-79	-161	-7	-4	1	0	-1	0	0
80-84	71	8	-8	5	-1	1	4	0
85-89	74	5	1	2	1	1	1	0
90-94	14	8	-2	-1	0	2	-1	0
95-99	-240	1	1	1	-6	1	1	0
100+	-151	1	1	-3	0	-7	1	0
Total	-821	-58	30	15	-9	-15	2	0



# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

## Neighboring Dataset $D'$

Histogram query  
using the **Laplace**  
mechanism

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



Age Group	$\epsilon = 0.01$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$	$\epsilon = 0.5$	$\epsilon = 1.0$	$\epsilon = 5.0$
< 5	19	-16	0	1	-1	2	-1	0
5-9	63	6	3	0	2	-2	-1	0
10-14	-30	4	4	7	0	-1	-3	0
15-19	100	-36	-2	2	-2	-6	2	0
20-24	-130	-3	12	1	3	0	2	0
25-29	-68	-21	-3	3	1	4	-1	0
30-34	-41	1	2	-4	0	-1	-1	0
35-39	-27	2	1	4	-3	-1	4	0
40-44	30	-4	-5	-3	-5	3	-1	0
45-49	-6	12	-9	0	-4	0	-1	0
50-54	-138	-2	3	-1	-1	2	0	1
55-59	-18	7	-13	4	-6	4	0	-1
60-64	-114	-26	-3	-1	0	0	1	1
65-69	47	-2	-2	2	0	4	4	0
70-74	-54	0	11	-2	1	2	0	0
75-79	63	2	1	1	1	4	0	0
80-84	-21	-14	0	-9	8	-1	0	0
85-89	4	3	-1	-9	0	1	0	0
90-94	14	6	-3	-2	-3	3	-1	0
95-99	-14	1	-9	1	1	1	1	0
100+	-175	1	-1	1	1	-3	0	0
Total	-496	-79	-14	-4	-7	15	4	1





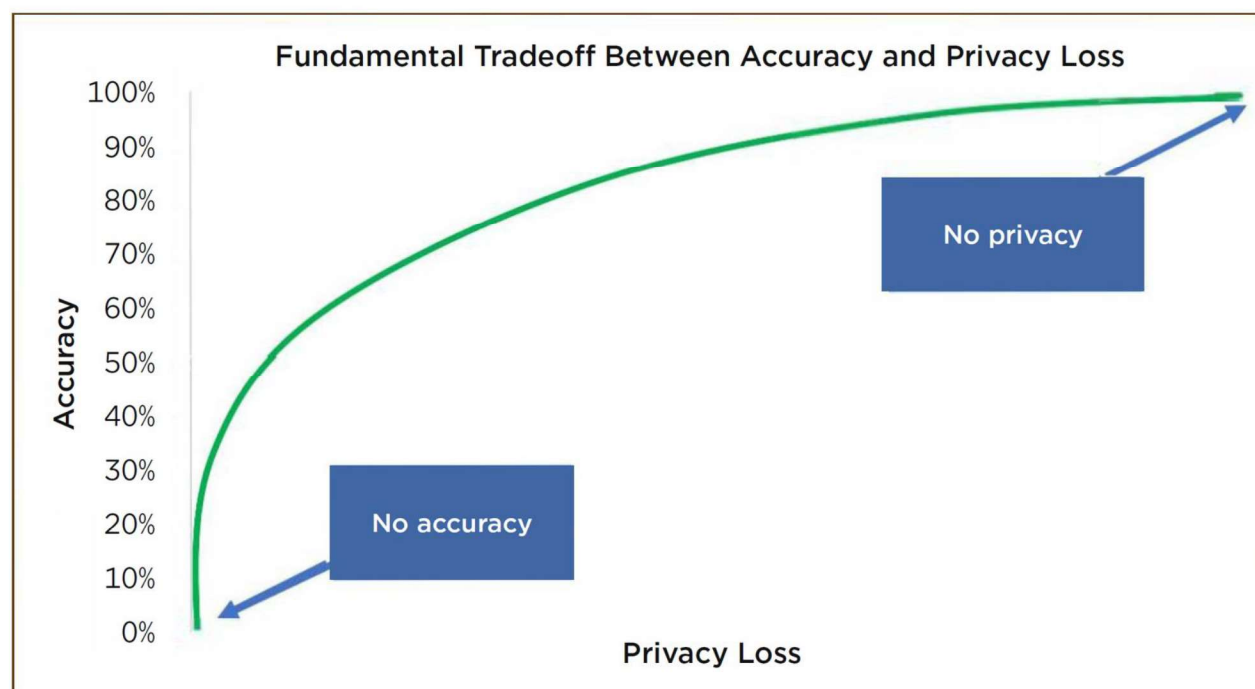
# 15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS

03-05 OCTOBER 2022

Organized by the Philippine Statistical System  
Spearheaded by the Philippine Statistics Authority



## Accuracy-privacy tradeoff



Source: US Census Bureau



## Discussion

### **Appropriate algorithm for implementing differential privacy**

Design algorithms that fits our need; determination of privacy budget; accuracy specification for pre-determined tabulations

### **Differentially private microdata**

Generate rich synthetic data with differential privacy

### **Explore interactive data release with differential privacy**

## References

Abowd, J. (2019). "Balancing Privacy and Accuracy: New Opportunity for Disclosure Avoidance Analysis." Research Matters (blog). US Census Bureau, October 29, 2019.  
[https://www.census.gov/newsroom/blogs/research-matters/2019/10/balancing\\_privacyan.html](https://www.census.gov/newsroom/blogs/research-matters/2019/10/balancing_privacyan.html).

Abowd, J. (n.d.). Tweakitorial: Reconstruction-abetted re-identification attacks and other traditional vulnerabilities. Retrieved October 4, 2022 from <https://blogs.cornell.edu/abowd/special-materials/245-2/>

Dwork, C. and Roth, A. (2014). The Algorithmic Foundations of Differential Privacy. Found. Trends Theor. Comput. Sci. 9, 3–4 (August 2014), 211–407. <https://doi.org/10.1561/04000000042>

Garfinkel, S., Abowd, J., and Martindale, C. (2018). Understanding Database Reconstruction Attacks on Public Data: These attacks on statistical databases are no longer a theoretical danger. <https://dl.acm.org/doi/pdf/10.1145/3291276.3295691>

Sweeney, L. (2002). k-anonymity: a model for protecting privacy. International Journal on Uncertainty, Fuzziness and Knowledge-based Systems, 10 (5), 2002; 557-570.



# **15<sup>TH</sup> NATIONAL CONVENTION ON STATISTICS**

03-05 OCTOBER 2022



*Organized by the Philippine Statistical System Spearheaded by the Philippine  
Statistics Authority*



## Thank you!



<http://www.psa.gov.ph/ncs>



<http://openstat.psa.gov.ph>



<https://twitter.com/PSAgovph>



<https://www.facebook.com/PhilippineStatisticsAuthority>