



Web Scraping of Commodities for Consumer Price Index in the National Capital Region, Philippines

GLEN G. POLO

Price Statistics Division, Economic Sector Statistics Service
Philippine Statistics Authority

Meeting of the Groups Experts on Consumer Price Index
Geneva, Switzerland
07 to 09 June 2023

Outline

- I. Introduction**
- II. Methodology**
- III. Results**
- IV. Issues and Challenges**
- V. Ways Forward**

I. Objective of the Study

- To know whether prices collected from websites via web scraping can be used as substitute for the data collected via traditional survey in computing the 2012-based CPI for National Capital Region, Philippines.
- To be used as benchmark for the use of Big Data for official statistics

II. Methodology

Geographic Domain:
National Capital Region



II. Methodology



Frequency of Collection:
Daily (except in Saturdays and Sunday)

II. Methodology

Sample Outlets/Websites:



look good. feel great.



II. Methodology

Total No. of URLs Web Scraped: 1,354

Name of Online Stores	No. of URLs	Commodity Division Code									
		01	02	03	04	05	06	07	08	09	11
Total	1,354	402	15	94	38	231	74	5	8	233	254
Abensons	16					13				3	
Ace Hardware	15				4	11					
Ansons	12					11				1	
Lazada	552	155	11	39	14	87	16	2	3	107	118
National Bookstore	3	1	1				1				
PushKart	23									23	
Shopee	74	65	1			1					7
Watsons	539	151	2	43	14	84	16	3	5	96	125
Western Appliance	45						41				4
Wilcon	17					16				1	
Zalora	16				6	8				2	
Zagana	30	30									

Legend:

01 – Food and Non-Alcoholic Beverages
02 – Alcoholic Beverages and Tobacco
03 – Clothing and Footwear

04 – Housing, Water, Electricity, Gas and Other Fuels
05 – Furnishing, Household Equipment and Routine Household Maintenance

06 – Health
07 – Transport
08 – Communication

09 – Recreation and Culture
11 – Restaurant and Miscellaneous Goods and Services

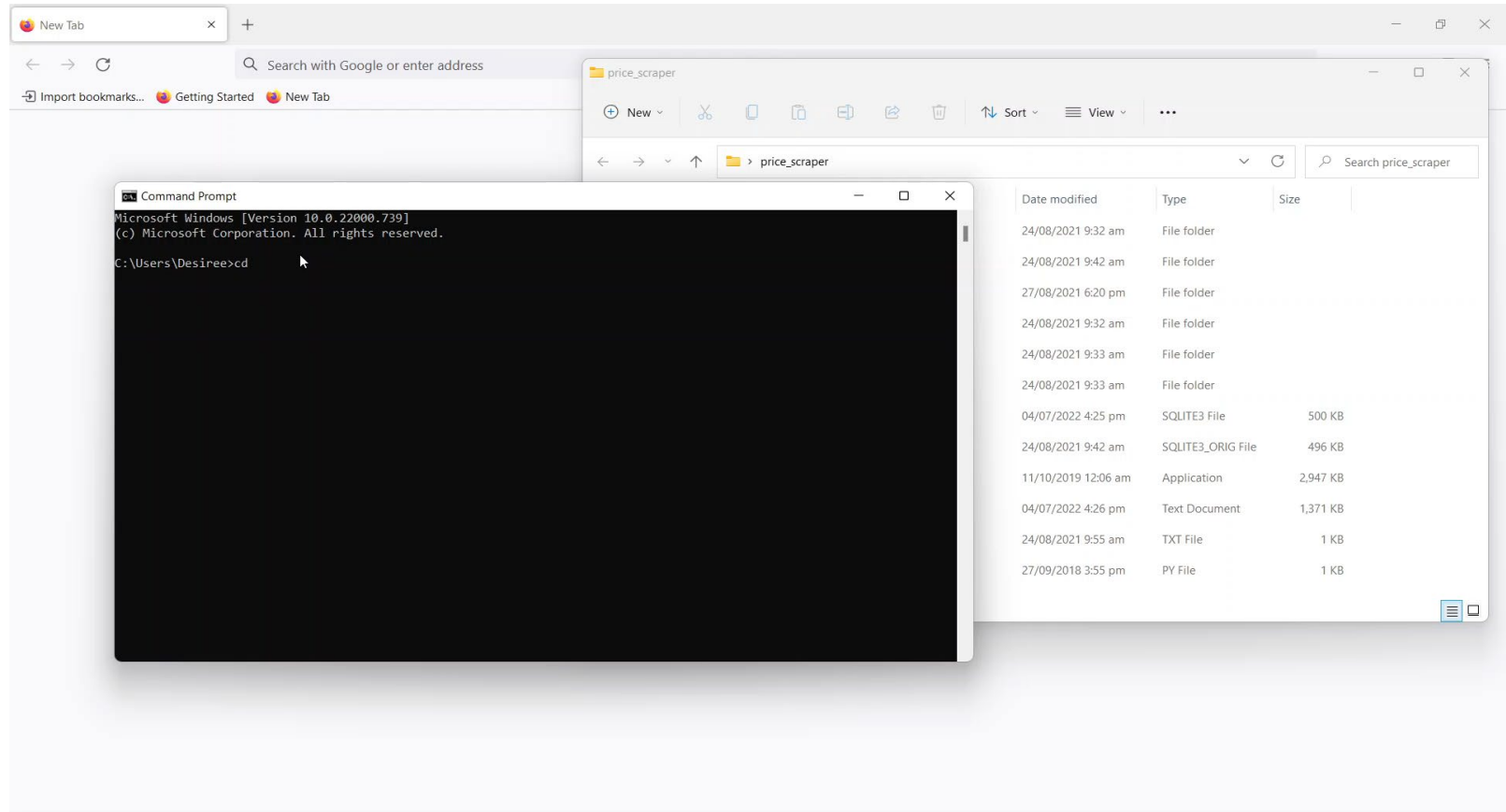
II. Methodology

Total No. of Commodities Web Scraped: 517

Division	No. of Commodities Web Scraped
01 - Food and Non-Alcoholic Beverages	183
02 - Alcoholic Beverages and Tobacco	10
03 - Clothing and Footwear	41
04 - Housing, Water, Electricity, Gas and Other Fuels	14
05 - Furnishing, Household Equipment, and Routine Household Maintenance	85
06 - Health	41
07 - Transport	2
08 - Communication	2
09 - Recreation and Culture	64
11 - Restaurants and Miscellaneous Goods and Services	75
Total	517

II. Methodology

- Web Scraping Application



II. Methodology

• Web Scraping Application: **Folders**

CPI WebScraper Home Single Product Multiple Products

[Add Product](#)

Single Products

Show entries Search:

Include	Name	Total Urls	Action
<input checked="" type="checkbox"/>	abensons	22	Delete
<input type="checkbox"/>	ansons	21	Delete
<input type="checkbox"/>	pushkart	82	Delete
<input type="checkbox"/>	acehardware	25	Delete
<input type="checkbox"/>	waltermart	79	Delete
<input type="checkbox"/>	watsons	61	Delete
<input type="checkbox"/>	wilcon	20	Delete
<input type="checkbox"/>	zalora	16	Delete
<input type="checkbox"/>	nationalbookstore	37	Delete
<input type="checkbox"/>	zagana	31	Delete

Showing 1 to 10 of 26 entries
[Previous](#)
[1](#)
[2](#)
[3](#)
[Next](#)

II. Methodology

- **Web Scraping Application:**
HTML Structure

HTML Structure and Listed Urls : [abensons](#)

HTML Structure

Folder:

Description Element:

Subdetails Element:

Price Element

Sale Price Element

[Save Changes](#) [Return](#)

List of Urls: [Add Url](#)

Show entries Search:

Url	Action
https://www.abenson.com/apple-ipad-mini-5-wi-fi-64gb-space-gray.html	Edit Delete
https://www.abenson.com/condura-ctd700mni.html	Edit Delete
https://www.abenson.com/dowell-di-583ns.html	Edit Delete
https://www.abenson.com/es-w600.html	Edit Delete
https://www.abenson.com/f-40dyp.html	Edit Delete
https://www.abenson.com/gs-600.html	Edit Delete
https://www.abenson.com/hi-89.html	Edit Delete
https://www.abenson.com/l1-ls-l2.html	Edit Delete
https://www.abenson.com/la-germania-e-726-w.html	Edit Delete
https://www.abenson.com/panasonic-na-s6518bsp.html	Edit Delete

II. Methodology

- **Web Scraping Application: Completion Prompt**

CPI WebScraper Home Single Product Multiple Products

Scraping Complete

Type	Action	Message
Single	Initialized	abensons -- Scraping Initialized
Single	Completed	abensons -- Scraping Complete

[Return to Homepage](#)

II. Methodology

- **Web Scraping Application:**
Sample Output

	A	B	C	D	E	F	G
1	Url	Description	Sub Details	Price	Sale Price		
2	https://wv	APPLE IPAD MINI 5 WI-FI	Item is discontinued.		23,990		
3	https://wv	CONDURA CTD700MNI	Item is discontinued.	19,997			
4	https://wv	DOWELL DI 583NS	SKU 161693	798			
5	https://wv	SHARP ES-W600	SKU 112944	3,997			
6	https://wv	PANASONIC F-40DYP	SKU 56136		1,748		
7	https://wv	HANABISHI GS 600	SKU 3776		648		
8	https://wv	HANABISHI HI-89	SKU 96012		698		
9	https://wv	PANASONIC NA-S6518BSP	SKU 161243	4,799			
10	https://wv	ASAHI RB-6004	SKU 118847		2,098		
11	https://wv	STANDARD SDS 12W	SKU 135746		1,298		
12	https://wv	STANDARD SGS 235S 2B	SKU 136929		1,998		
13	https://wv	SHARP SJ DTH55BS SL	Item is discontinued.	11,697			
14	https://wv	SONY KDL 32R307F	Item is discontinued.	14,499			
15	https://wv	CANON POWERSHOT SX620HS	SKU 144585	15,198			
16	https://wv	TEFAL RK104E	SKU 163548		3,895		
17	https://wv	TEFAL RK7405	SKU 161277		8,995		
18	https://wv	TEFAL RK8145	SKU 161278		10,995		
19	https://wv	LA GERMANIA E-726 W	SKU 170556	6,798			
20	https://wv	TEKNO TKX- 180	SKU 164815	648			
21	https://wv	TEKNO TKX-780	SKU 164814	1,278			
22	https://wv	KELVINATOR WKELH010EA	SKU 147117	18,498			
23							
24							
25							
26							
27							

single-abensons-742022

Ready Accessibility: Unavailable

II. Methodology

• Data Processing:

1. Validations are done daily: consistency checking, checking for the presence of web scraped prices, checking if the links are still active and if the price being collected is correct.
2. Computation of Average Prices, Indices, M-o-M Growth Rate, Y-o-Y Growth Rate follow the official CPI compilation.

II. Methodology

CPI Computation:

- Computation of Average Prices, Indices, M-o-M Growth Rate, Y-o-Y Growth Rate follow the official CPI compilation.
- Two types of CPI were computed and compared with the Official CPI:
 - Online – All prices used are collected from websites (web scraped)
 - Hybrid – combination of offline (traditional survey) and online (web scraped) prices.

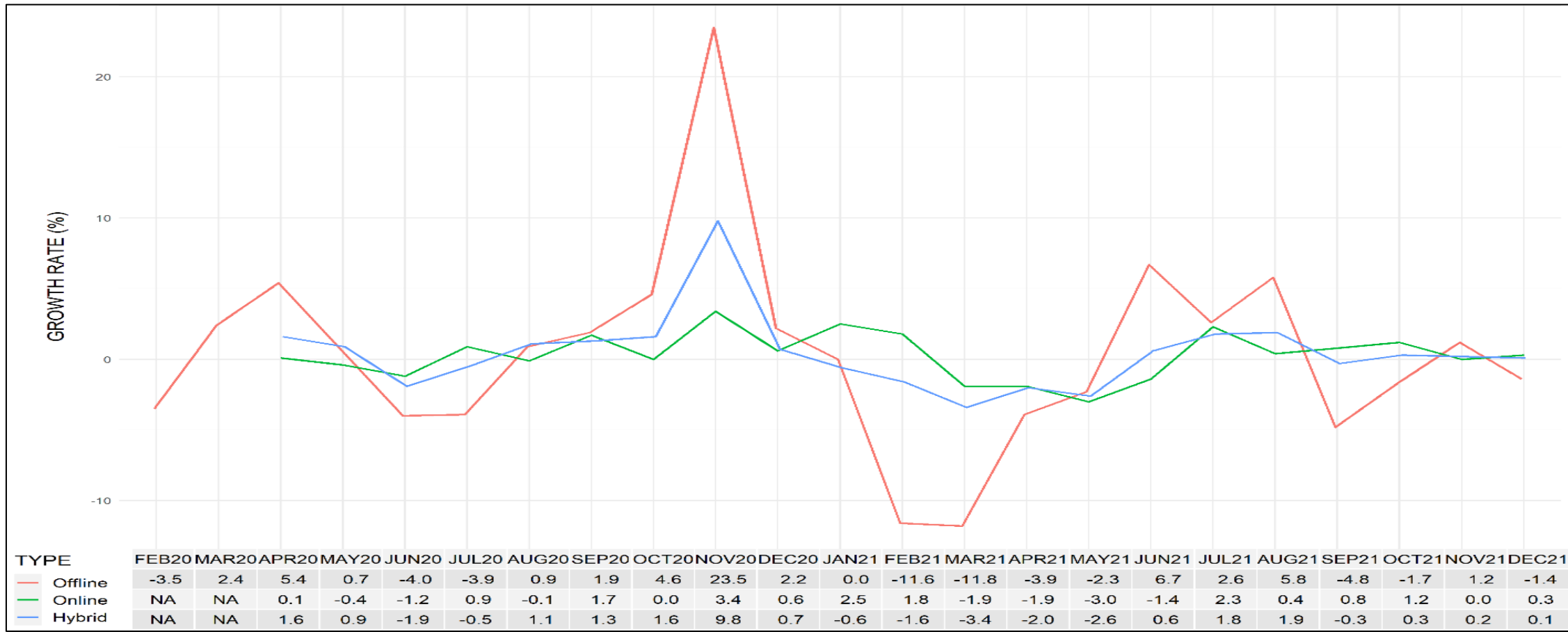
III. Results

Year-on-Year: Fish and Seafood



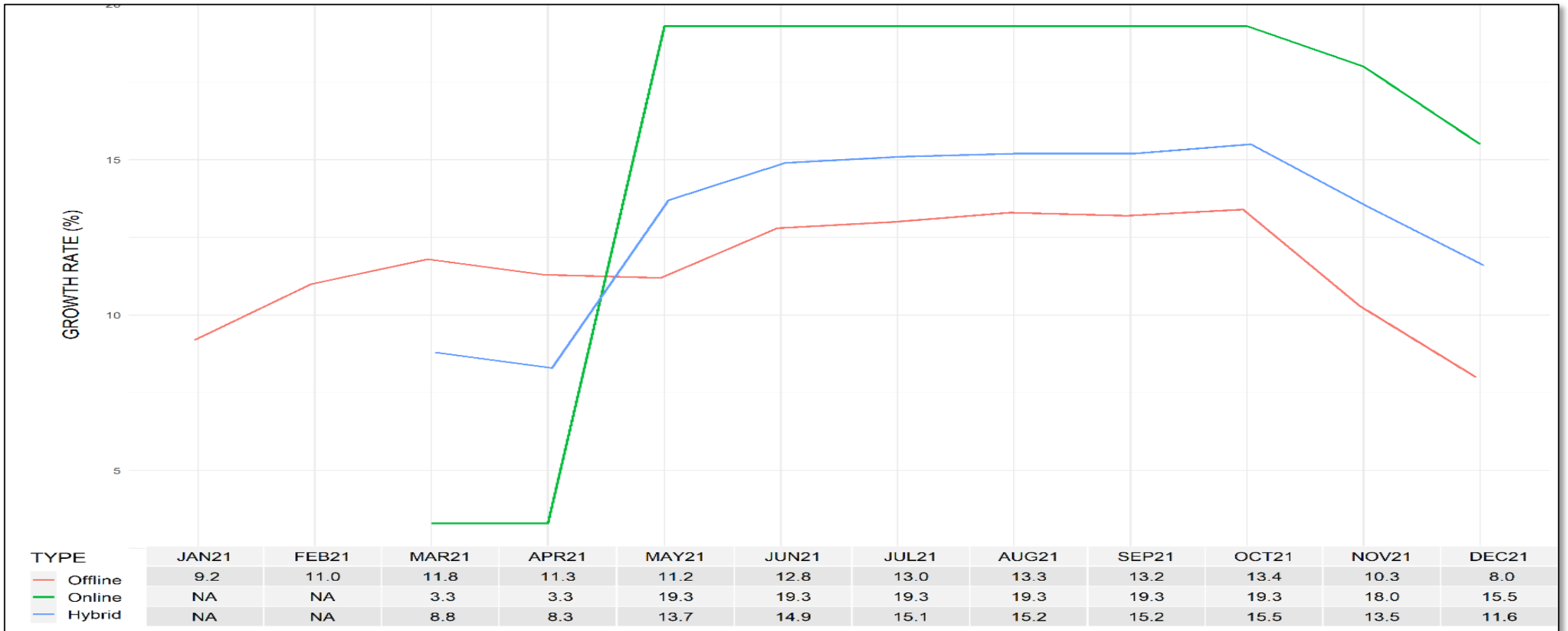
III. Results

Year-on-Year: Vegetables



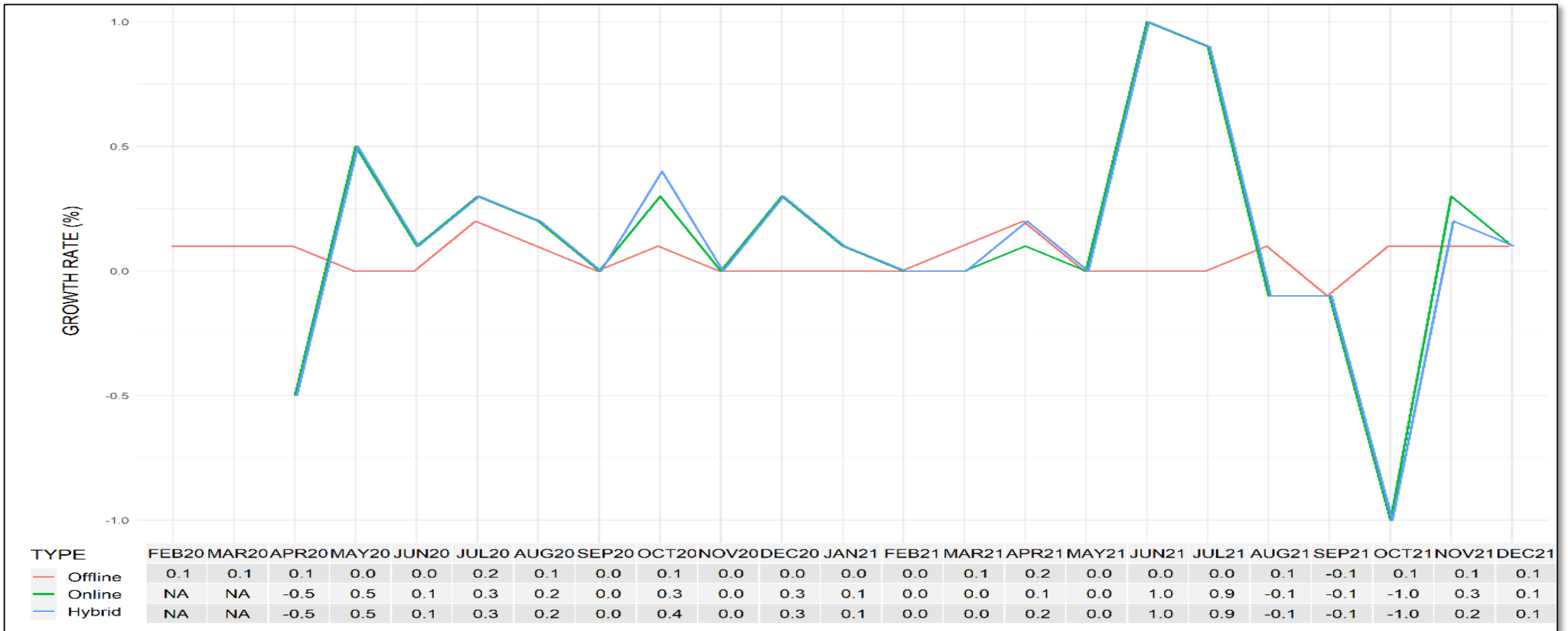
III. Results

Year-on-Year: Tobacco



III. Results

Year-on-Year: Garments



IV. Issues and Challenges

- 1. Websites selected for scraping are not CPI sample outlets. Chosen based on availability of commodities listed in the market basket**
- 2. Not all web scraped commodities have exactly similar specifications with those from the market-basket.**
- 3. Not all of the subclass (5-digit level PCOICOP) and class (4-digit level PCOICOP) have complete commodities.**
- 4. There is an issue with legality and ethics.**

V. Ways Forward

1. **Start the web scraping simultaneous with price collection for the new CPI series**
2. **Collect prices from the websites of the CPI sample outlets**

Authors:

Divina Gracia L. Del Prado, Deputy National Statistician

Elena G. Varona (ret.), Chief, Price Statistics Division (PSD)

Glen G. Polo, Officer-in-Charge, PSD

Desiree R. Robles, Senior Statistical Specialist

Rosario S. Lodovice, Statistical Specialist II

Jo Louise L. Buhay, Statistical Specialist I

THANK YOU!



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



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



<https://twitter.com/PSAgovph>



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

