

Driving Growth & Revenue: A Predictive Analysis for Blinkit

Under the guidance of

Monika Bhattacharya

by

- Sahil Singh
- Atharv Koltharkar
- Atreyee Roy
- Rohan Karmakar
- Pritam Naskar
- Sourav Ghosh



Datasets

1

Blinkit Customer Orders

Transactional data of customer purchases including products, time, and payment.

4

Blinkit Customer Registration

Data on registered users like ID and location.

2

Blinkit Customer Feedback

Customer ratings and reviews for orders and services.

5

Blinkit Marketing Performance

Insights from marketing campaigns and conversions.

3

Blinkit Delivery Performance

Metrics on delivery time, success rate, and delays.

6

Blinkit Inventory Data

Data on stock levels and product availability.

Snapshots of Dataset for Revenue Prediction

order_id	customer_order_date	promised_delivery_time	actual_delivery_time	delivery_status	order_total	payment_method	delivery_partner_id	store_id
1961864118	30065862	17-07-2024 08:34	17-07-2024 08:52	17-07-2024 08:47	On Time	3197.07	Cash	63230 4771
1549769649	9573071	28-05-2024 13:14	28-05-2024 13:25	28-05-2024 13:27	On Time	976.55	Cash	14983 7534
9185164487	45477575	23-09-2024 13:07	23-09-2024 13:25	23-09-2024 13:29	On Time	839.05	UPI	39859 9886
9644738826	88067569	24-11-2023 16:16	24-11-2023 16:34	24-11-2023 16:33	On Time	440.23	Card	61497 7917
5427684290	83298567	20-11-2023 05:00	20-11-2023 05:17	20-11-2023 05:18	On Time	2526.68	Cash	84315 2741
3265154092	43367112	18-03-2023 16:29	18-03-2023 16:49	18-03-2023 16:48	On Time	3161.43	UPI	554 3442
4898355547	13284996	16-04-2023 18:50	16-04-2023 19:01	16-04-2023 19:02	On Time	956.4	Card	14630 1318
6568151549	88866835	31-03-2024 06:26	31-03-2024 06:37	31-03-2024 06:39	On Time	905.47	Cash	67714 115
6006693867	24496983	13-07-2023 23:49	14-07-2023 00:02	14-07-2023 00:05	On Time	1371.17	Card	91362 9021
374186990	52215833	09-08-2023 01:17	09-08-2023 01:37	09-08-2023 01:44	Slightly Delayed	1601.19	Wallet	77203 7955

Blinkit Orders Data

product_id	date	stock_received	damaged_stock
153019	17-03-2023	4	2
848226	17-03-2023	4	2
965755	17-03-2023	1	0
39154	17-03-2023	4	0
34186	17-03-2023	3	2
478531	17-03-2023	4	2
812607	17-03-2023	4	0
818990	17-03-2023	4	0
14145	17-03-2023	0	2

Blinkit Inventory Data

Snapshots of Dataset for User Growth Prediction

campaign_id	campaign_name	date	target_audience	channel	impressions	clicks	conversions	spend	revenue_generated	roas
548299	New User Discount	05-11-2024	Premium	App	3130	163	78	1431.85	4777.75	3.6
390914	Weekend Special	05-11-2024	Inactive	App	3925	494	45	4506.34	6238.11	2.98
834385	Festival Offer	05-11-2024	Inactive	Email	7012	370	78	4524.23	2621	2.95
241523	Flash Sale	05-11-2024	Inactive	SMS	1115	579	86	3622.79	2955	2.84
595111	Membership Drive	05-11-2024	New Users	Email	7172	795	54	2888.99	8951.81	2.22
176344	Category Promotion	05-11-2024	Inactive	App	4333	421	40	2714.78	6697.88	3.89
875646	App Push Notification	05-11-2024	Premium	Email	4752	937	21	3584.76	9152.99	2.14
5516	Email Campaign	05-11-2024	Premium	Social Media	1559	418	68	1457.88	2156.09	3.92
989534	Referral Program	05-11-2024	New Users	SMS	3838	911	45	4808.64	8074.17	1.7
243346	New User Discount	04-11-2024	Inactive	SMS	6289	834	14	3746.35	4474.02	2.56

Blinkit Marketing Data

customer_id	customer_name	phone	area	pincode	registration_date	customer_segment	total_orders	avg_order_value
97475543	Niharika Nagi	9.12988E+11	Udupi	321865	13-05-2023	Premium	13	451.92
22077605	Megha Sachar	9.15123E+11	Aligarh	149394	18-06-2024	Inactive	4	825.48
47822591	Hema Bahri	9.10034E+11	Begusarai	621411	25-09-2024	Regular	17	1969.81
79726146	Zaitra Vig	9.16264E+11	Kozhikode	826054	04-10-2023	New	4	220.09
57102800	Januja Verma	9.17294E+11	Ichalkaranji	730539	22-03-2024	Inactive	14	578.14
54748429	Darsh More	9.12582E+11	Visakhapatnam	883122	22-04-2024	New	2	669.35
49152878	Patrick Sandhu	9.12413E+11	Gwalior	649817	13-09-2024	Premium	4	973.62
16379942	Harshil Kuruvilla	9.13807E+11	Orai	332997	10-04-2024	Inactive	15	1370.56
11071601	Ojas Ahuja	9.12865E+11	Buxar	528426	04-07-2023	Inactive	17	1950.54

Blinkit Customer Registration Data

Problem Statement 1

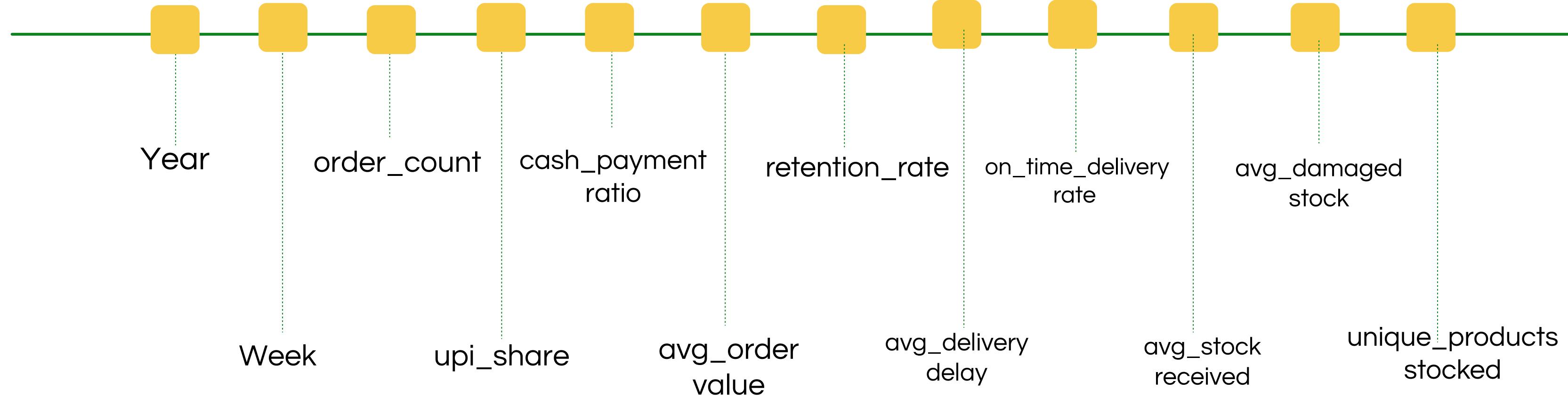
Predicting Weekly Revenue for
Blinkit using Linear Regression
Model.



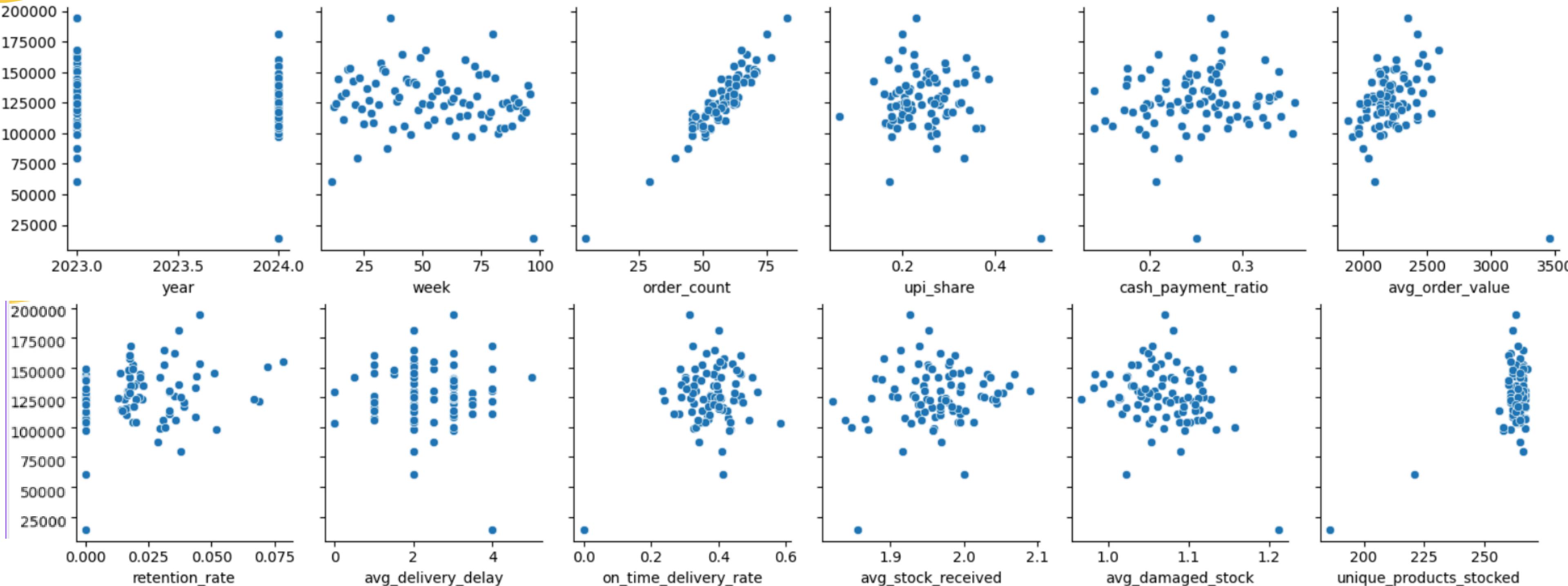
Data Preparation

Response Variable: Weekly Revenue

Predictors:



Exploratory Data Analysis



Initial Model Fitting

Multiple Linear Regression Model of Weekly Revenue

$$\begin{aligned} \text{Weekly Revenue}_t = & \beta_0 + \beta_1 \cdot \text{year}_t + \beta_2 \cdot \text{week}_t + \beta_3 \cdot \text{order_count}_{t-1} + \beta_4 \cdot \text{UPI_share}_{t-1} + \beta_5 \cdot \text{cash_payment_ratio}_{t-1} + \beta_6 \cdot \text{avg_order_value}_{t-1} + \beta_7 \cdot \text{retention_rate}_{t-1} \\ & + \beta_8 \cdot \text{avg_delivery_delay}_{t-1} + \beta_9 \cdot \text{on_time_delivery_rate}_{t-1} + \beta_{10} \cdot \text{avg_stock_received}_{t-1} + \beta_{11} \cdot \text{avg_damaged_stock}_{t-1} + \beta_{12} \cdot \text{unique_products_stocked}_{t-1} + \varepsilon \end{aligned}$$

- We have taken all the predictors except year and week at time (t-1) (lag 1).
- Year and week are at time t
- Response Variable is at time t.
- This setup helps us predicting the future weeks' revenue

Results of Multiple Linear Regression

OLS Regression Results						
Dep. Variable:	weekly_revenue	R-squared:	0.984			
Model:	OLS	Adj. R-squared:	0.982			
Method:	Least Squares	F-statistic:	466.6			
Date:	Tue, 15 Apr 2025	Prob (F-statistic):	6.42e-64			
Time:	04:01:22	Log-Likelihood:	-823.03			
No. Observations:	87	AIC:	1668.			
Df Residuals:	76	BIC:	1695.			
Df Model:	10					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-2.066e+05	3.12e+04	-6.631	0.000	-2.69e+05	-1.45e+05
order_count	2240.2883	45.565	49.167	0.000	2149.537	2331.039
upi_share	-1.931e+04	6440.256	-2.998	0.004	-3.21e+04	-6481.893
cash_payment_ratio	-1.172e+04	7988.522	-1.467	0.147	-2.76e+04	4193.979
avg_order_value	49.2226	2.349	20.955	0.000	44.544	53.901
retention_rate	4703.7493	2.02e+04	0.233	0.817	-3.56e+04	4.5e+04
avg_delivery_delay	1401.5824	672.440	2.084	0.040	62.302	2740.863
on_time_delivery_rate	3.126e+04	1.02e+04	3.078	0.003	1.1e+04	5.15e+04
avg_stock_received	4531.0909	7553.632	0.600	0.550	-1.05e+04	1.96e+04
avg_damaged_stock	-1.327e+04	1.03e+04	-1.295	0.199	-3.37e+04	7143.840
unique_products_stocked	355.1363	55.556	6.392	0.000	244.487	465.786
Omnibus:	35.348	Durbin-Watson:	1.599			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	150.826			
Skew:	1.145	Prob(JB):	1.77e-33			
Kurtosis:	9.031	Cond. No.	2.05e+05			

Ols Summary on whole data

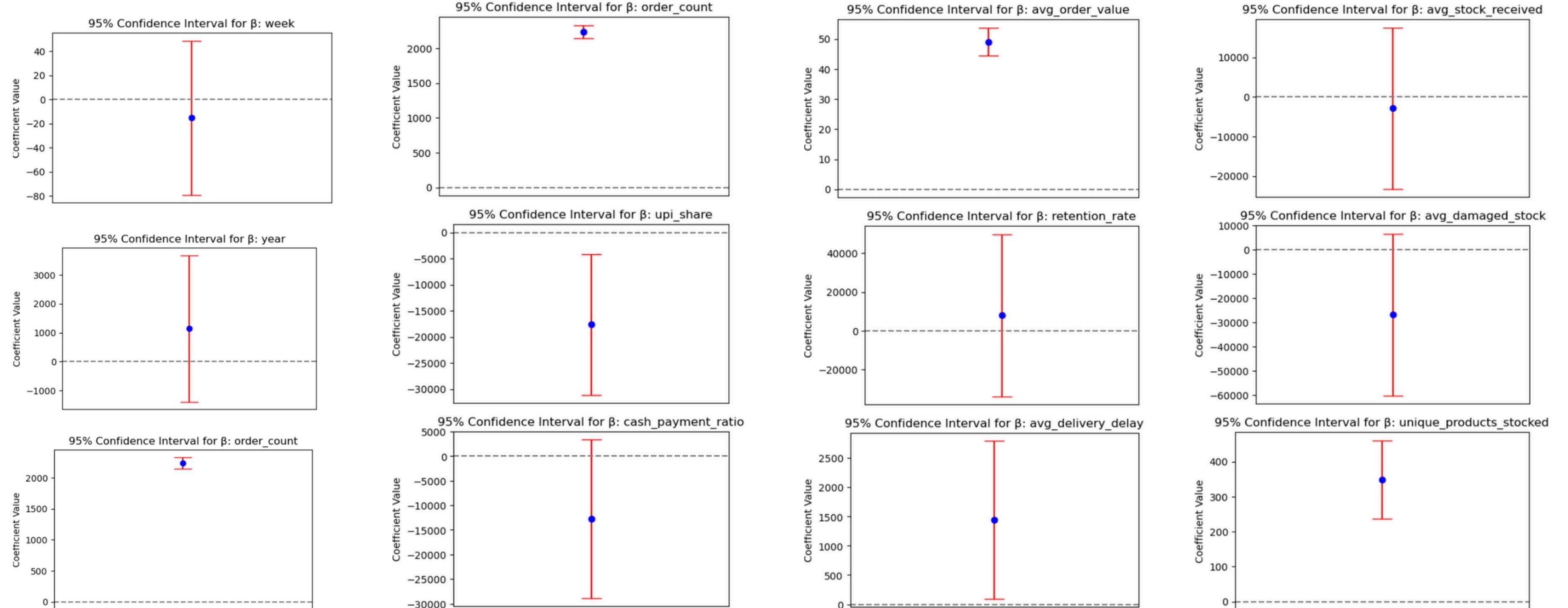
MSE: 57782808.61920831

R² Score: 0.8594537204958749

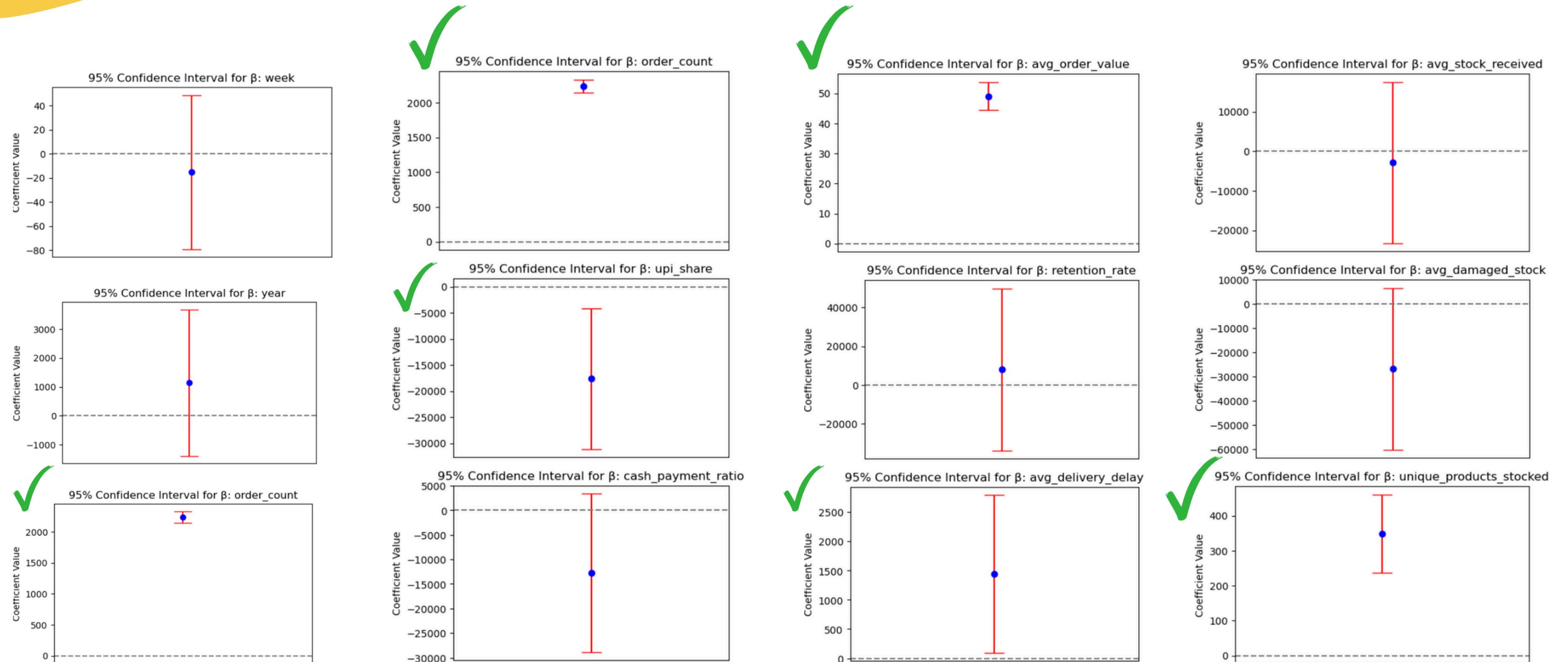
OLS results on test data

- **Model is highly significant** (F-statistic = 466.6, p < 0.001).
- **Significant predictors:** order_count, upi_share, avg_order_value, avg_delivery_delay, on_time_delivery_rate, unique_products_stocked.
- **High t-values** for order_count, avg_order_value, unique_products_stocked → strong influence on revenue.
- **Some evidence of non-normality in residuals** (JB test p-value ≪ 0.05, skew = 1.145, kurtosis = 9.03)

Confidence Interval plot for Regression Coefficients



Confidence Interval plot for Regression Coefficients



Confidence Interval Interpretation

	0	1
const	-7.580200e+06	2.593320e+06
year	-1.395030e+03	3.685439e+03
week	-7.902365e+01	4.877302e+01
order_count	2.151234e+03	2.334529e+03
UPI_Share	-3.108088e+04	-4.170874e+03
cash_payment_ratio	-2.889606e+04	3.426900e+03
avg_order_value	4.442847e+01	5.384725e+01
retention_rate	-3.375265e+04	4.988081e+04
avg_delivery_delay	9.343226e+01	2.797834e+03
on_time_delivery_rate	1.069981e+04	5.178959e+04
avg_stock_received	-2.325137e+04	1.758931e+04
avg_damaged_stock	-6.007317e+04	6.669881e+03
unique_products_stocked	2.366695e+02	4.613975e+02

- We visualized 95% confidence intervals for each regression coefficient to assess the **magnitude, uncertainty, and significance** of predictors.
- **Statistically significant predictors** (interval does not include 0):
`order_count`, `UPI_Share`, `avg_order_value`, `avg_delivery_delay`,
`on_time_delivery_rate`, `unique_products_stocked`
- These features show a reliable association with the target variable.
- Predictors like `year`, `week`, `retention_rate`, and `cash_payment_ratio` have intervals including 0, indicating **no statistically significant effect** at 95% confidence.
- Wider intervals (e.g., `retention_rate`, `avg_damaged_stock`) suggest higher uncertainty in estimating their effects.

Regression Diagnostics

1. Box Cox Transformation

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0 \\ \ln(y), & \text{if } \lambda = 0 \end{cases}$$

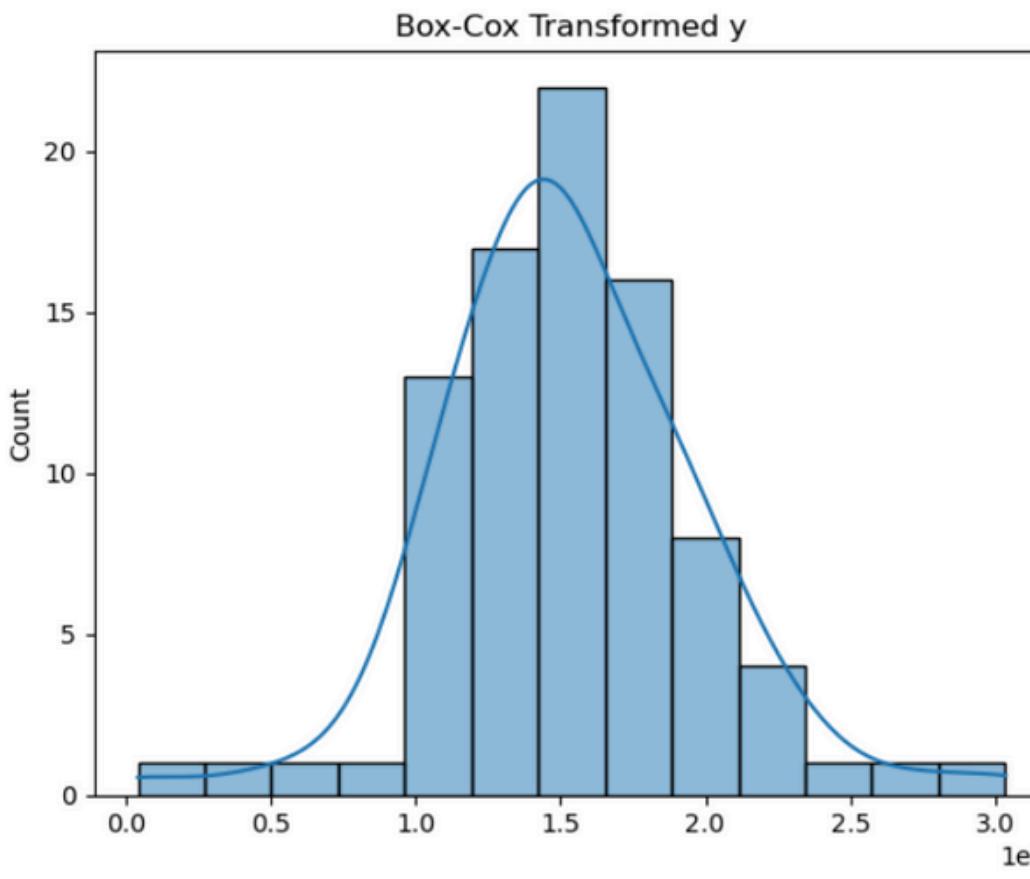
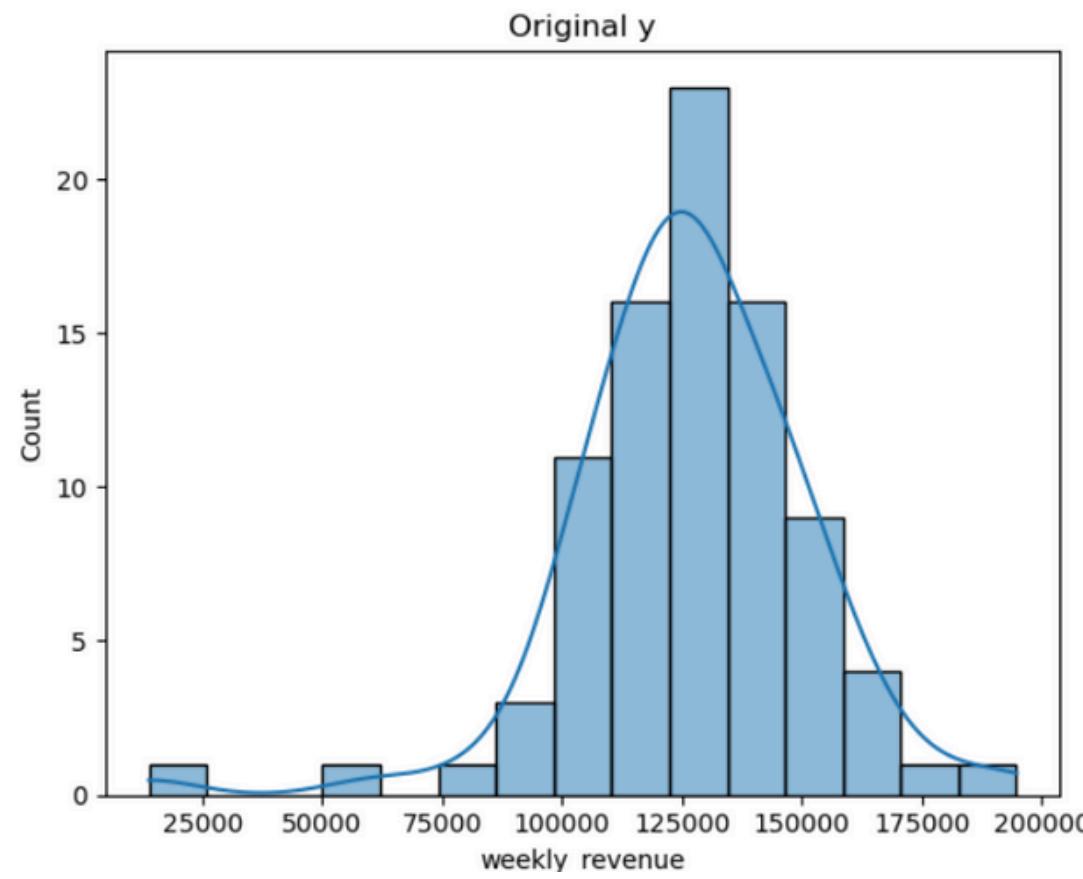
- y : The response variable (must be positive).
- λ : Transformation parameter to be estimated.

We applied the **Box-Cox transformation** to the response variable to stabilize variance and satisfy the assumptions of linear regression—specifically normality and homoscedasticity of residuals. **Using maximum likelihood estimation (MLE), we found the optimal transformation parameter λ .** The value of λ was estimated to lie between -2 and 2, which is a commonly accepted practical range. After applying the transformation, the residuals showed improved normality and more constant variance, enhancing model performance and interpretability.

Regression Diagnostics

Box Cox Transformation on Response Variable

Since the MSE is very large we applied Box Cox Transformation to scale down the the Response Variable. Then we plotted the residual vs Predicted Values before and after the Transformation.

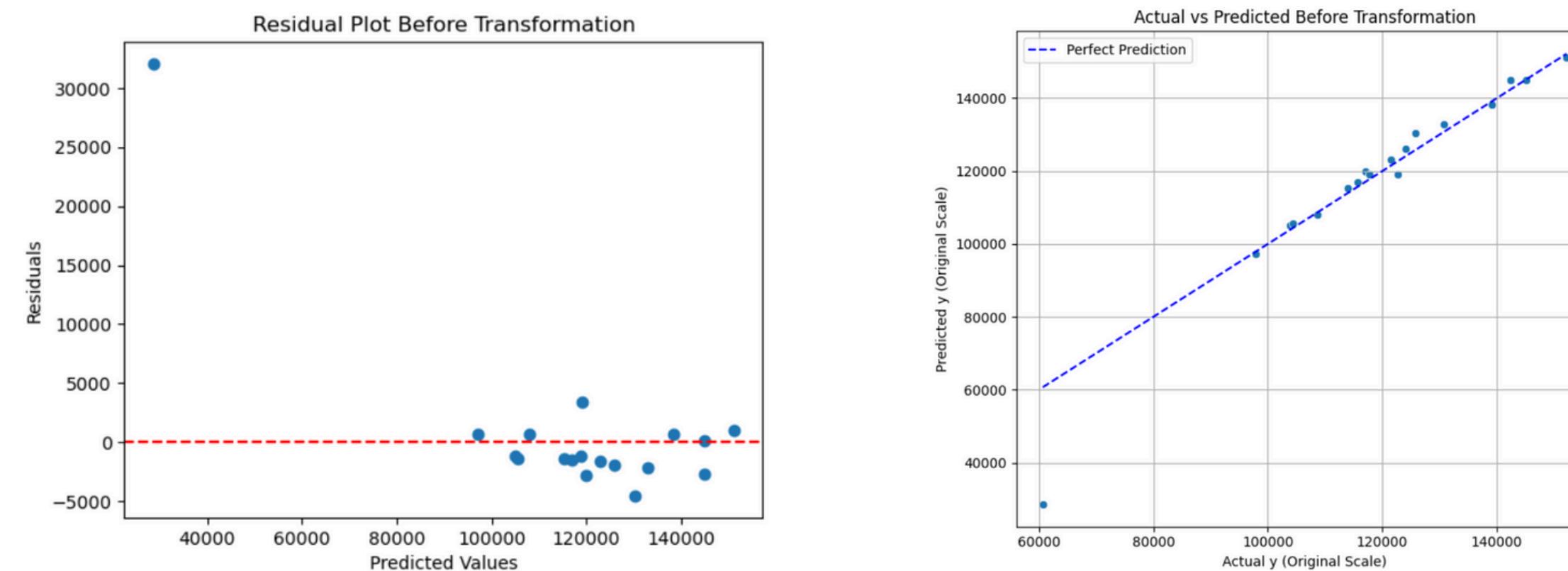


Box-Cox transformation was applied to correct negative skewness in `weekly_revenue`

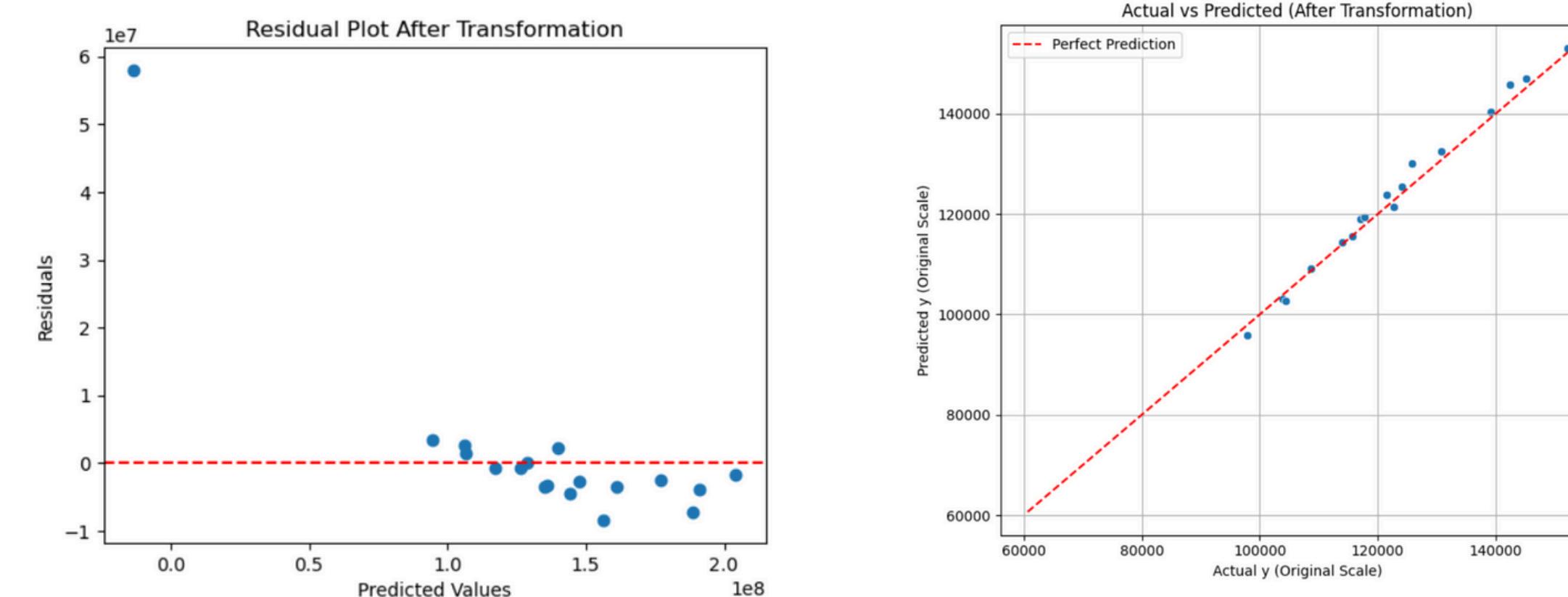
It **improved normality of the response**, aligning better with regression assumptions.

Regression Diagnostics

Before Transformation:

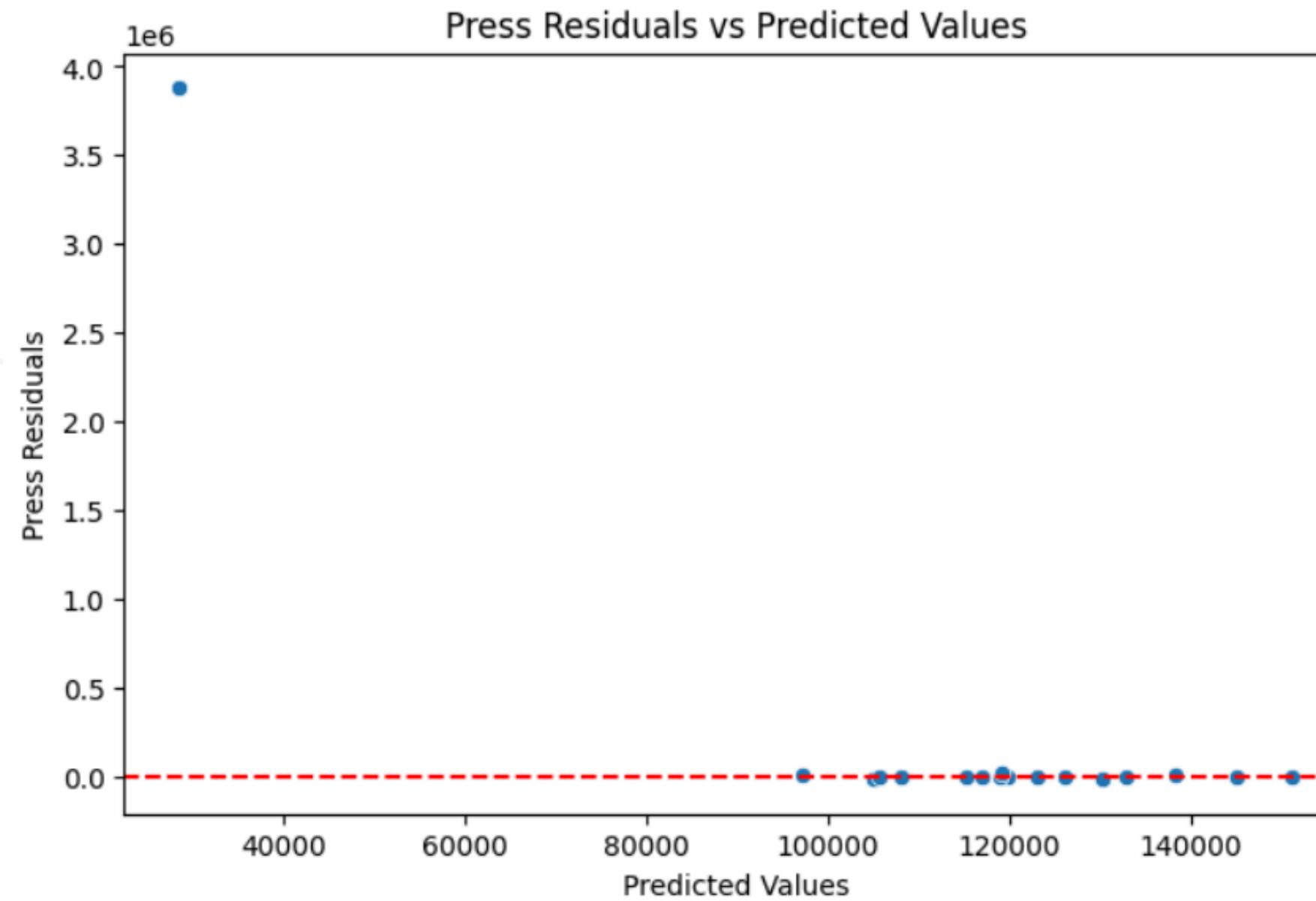


After Transformation:



Regression Diagnostics

2. Press Residuals



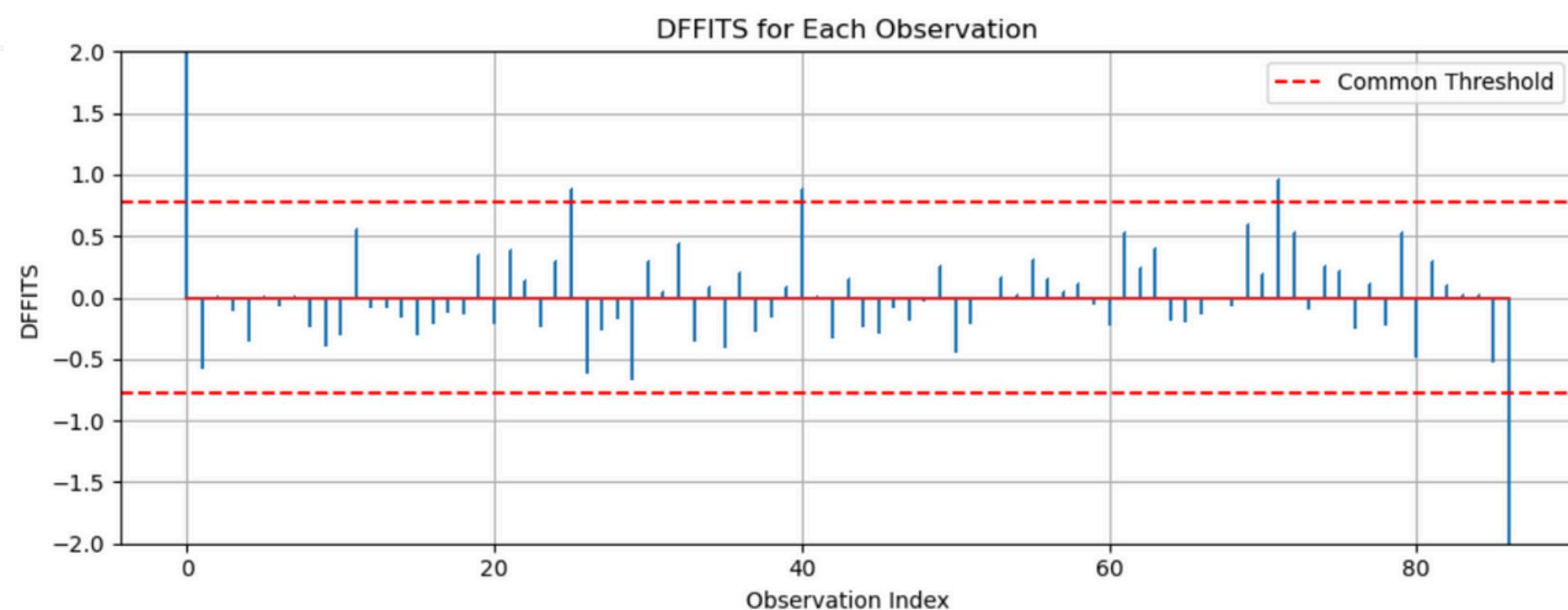
- **Most PRESS residuals are close to zero,** indicating good predictive performance for the majority of observations.
- **One clear outlier** suggests a data point with poor prediction accuracy that may need further investigation.

Outlier Diagnostics

1. Cook's D : Threshold: $D > 4/n$

Influential points (Cook's D > 0.0460): [0 25 40 71 86]

2. DFFITS:



Influential points based on DFFITS ($>|0.7731|$): [0 25 40 71 86]

Five **influential points** were identified using both **Cook's** Distance and **DFFITS**, indicating they have a significant impact on the model's predictions.

Outlier Diagnostics

3. DFBETAS:

The DFBETA for a specific observation i and coefficient j is calculated as:

$$\text{DFBETA}_{ij} = \frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\text{SE}(\hat{\beta}_j)}$$

const	0.848088
year	-0.759078
week	-1.510661
order_count	-0.355948
UPI_Share	-1.996221
Cash_Payment_Ratio	-1.251524
Avg_Order_Value	-4.930107
Retention_Rate	-0.864768
Avg_Delivery_Delay	1.484956
On_Time_Delivery_Rate	1.666815
Avg_Stock_Received	-0.337448
Avg_Damaged_Stock	-2.350873
Unique_Products_Stocked	-8.918328
Name:	0, dtype: float64

- Using DFFITS, we identified **5 influential observations** with values greater than the cutoff ($|DFFITS| > 0.7731$):
[0, 25, 40, 71, 86]
- We further validated these using DFBETAS, which quantify how much each regression coefficient changes when an observation is removed.
- For **observation 0**, multiple features showed large standardized changes:
 - Strong impact on:
unique_products_stocked (-8.91), avg_order_value (-4.93), upi_share (-1.99), avg_damaged_stock (-2.35)
 - These DFBETAS exceed the usual threshold of $|DFBETA| > 2 / \sqrt{n}$, suggesting significant influence on these coefficients.

Is there any way to improve the Model?

If yes, then how?



Making the Model Better

1. Checking Multi Collinearity :

	Feature	VIF
0	order_count	1.755791
1	upi_share	1.414567
2	cash_payment_ratio	1.307299
3	avg_order_value	1.759307
4	retention_rate	1.109829
5	avg_delivery_delay	3.314978
6	on_time_delivery_rate	4.446875
7	avg_stock_received	1.273674
8	avg_damaged_stock	1.496455
9	unique_products_stocked	2.287527

Component 1: Condition Index = 1.00
Component 2: Condition Index = 1.35
Component 3: Condition Index = 1.49
Component 4: Condition Index = 1.63
Component 5: Condition Index = 1.70
Component 6: Condition Index = 4.66
Component 7: Condition Index = 3.22
Component 8: Condition Index = 2.60
Component 9: Condition Index = 2.13
Component 10: Condition Index = 2.04

From the VIF Table and Condition Index Table it is evident that there is no multicollinearity among the predictors. That is, the predictors are independent of each other.

Making the Model Better

2. Feature Selection

- We used Recursive Feature Elimination that recursively removes the least important features to find the best k no of features for your model.
- Here we choose k=5 and get top 5 important features or key influenced features for our data

	Feature	Ranking	Selected
0	order_count	1	True
3	avg_order_value	1	True
5	avg_delivery_delay	1	True
6	on_time_delivery_rate	1	True
9	unique_products_stocked	1	True
1	upi_share	2	False
8	avg_damaged_stock	3	False
7	avg_stock_received	4	False
2	cash_payment_ratio	5	False
4	retention_rate	6	False

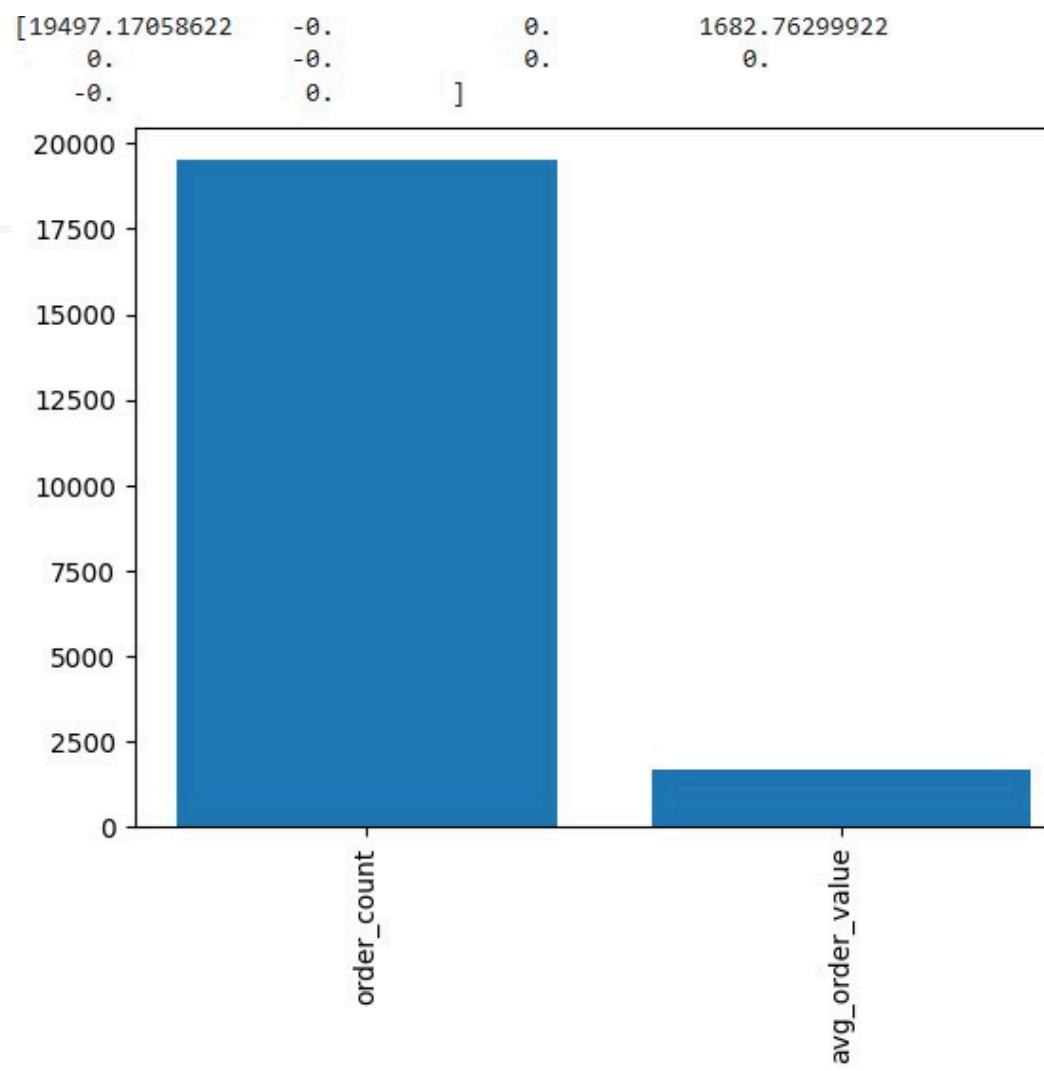
So we select 5 features that is,

- order count
- avg order value
- avg delivery delay
- on time delivery rate
- unique products stocked

Making the Model Better

3. Validating with Lasso Regression

- Lasso Regression applies L1 regularization, shrinking less important feature coefficients to zero.
- We used it to verify if predictors marked less significant by RFE drop to zero beyond a certain alpha threshold.



Lasso regression retained only the important predictors:

Among the features tested, `order_count` maintained a non-zero coefficient, indicating it has a strong predictive contribution.

Less significant predictors dropped to zero:

Features like `avg_order_value` were shrunk to zero, confirming the RFE (Recursive Feature Elimination) results and validating that these predictors have low influence in the presence of regularization.

APPLYING XG BOOST AND RANDOM FOREST & BACK PROPAGATION

XG_BOOST:

R2_SCORE: 0.81

RANDOM FOREST:

R2_SCORE: 0.81

REASON: In both cases they outperform linear regression model because they both are ensamble techniques

APPLYING XG BOOST AND RANDOM FOREST & BACK PROPAGATION

BACK PROPAGATION:

R2_SCORE: 0.85

REASON: Neural Network is built through multiple hidden layers where the starting layers catch premetive patterns while the ending layers catch complex patterns and give the solution that's why back propagation outperform every other techniques

Problem Statement 2

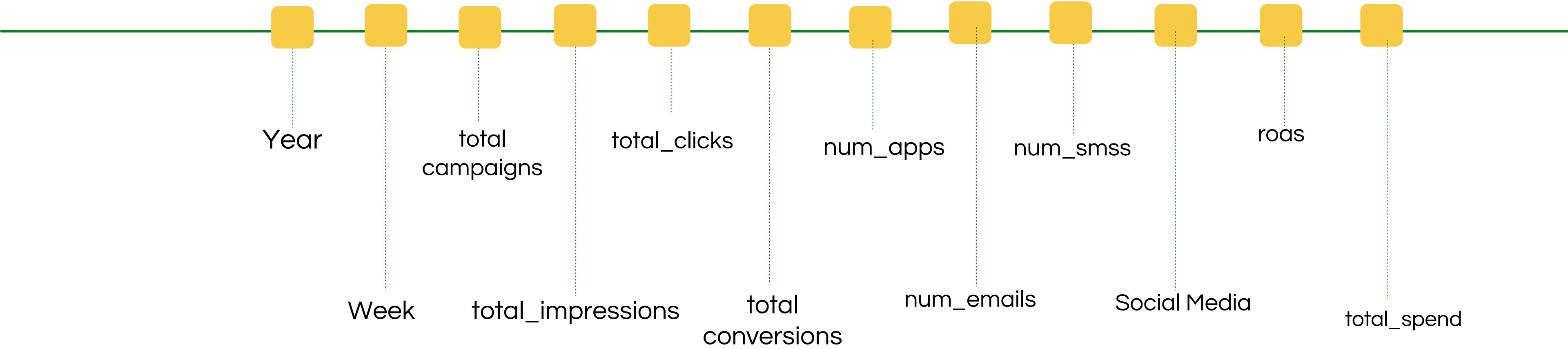
Predicting User Growth for Blinkit
using Linear Regression Model.



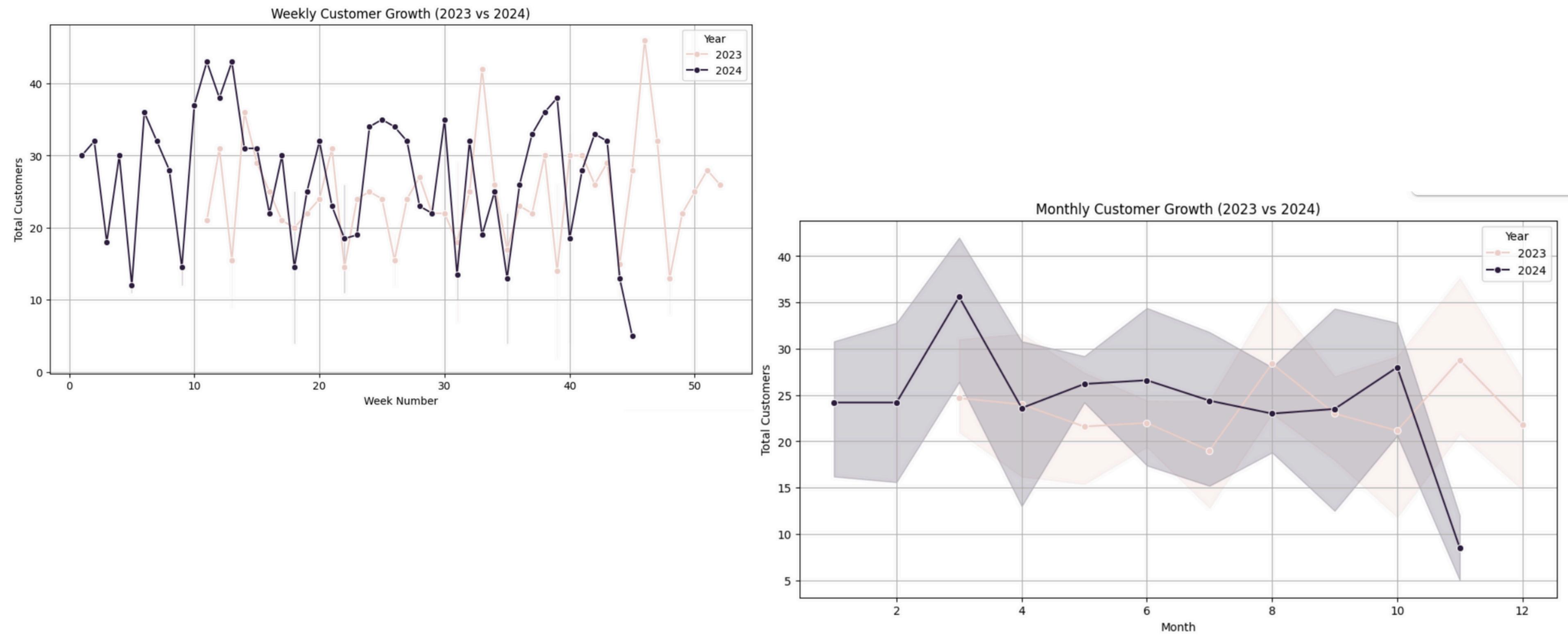
Data Preparation

Response Variable: Weekly Newly Registered Users

Predictors:



Exploratory Data Analysis



Initial Model Fitting

Multiple Linear Regression Model of Weekly New Users

$$\begin{aligned} \text{Weekly_New_Users}_t = & \beta_0 + \beta_1 \cdot \text{Year}_t + \beta_2 \cdot \text{Week}_t + \beta_3 \cdot \text{total_campaigns}_t + \beta_4 \cdot \text{total_impressions}_t + \beta_5 \cdot \text{total_clicks}_t + \beta_6 \cdot \text{total_conversions}_t \\ & + \beta_7 \cdot \text{num_apps}_t + \beta_8 \cdot \text{num_emails}_t + \beta_9 \cdot \text{num_smss}_t + \beta_{10} \cdot \text{SocialMedia}_t + \beta_{11} \cdot \text{roast}_t + \beta_{12} \cdot \text{total_spend}_t + \varepsilon_t \end{aligned}$$

- Unlike the weekly revenue model, we have taken all the predictors at time t.
- Response Variable is at time t.

Results of Multiple Linear Regression

OLS Regression Results							
Dep. Variable:	total_customers	R-squared:	0.704				
Model:	OLS	Adj. R-squared:	0.664				
Method:	Least Squares	F-statistic:	17.83				
Date:	Mon, 07 Apr 2025	Prob (F-statistic):	6.31e-19				
Time:	18:01:52	Log-Likelihood:	-314.22				
No. Observations:	103	AIC:	654.4				
Df Residuals:	90	BIC:	688.7				
Df Model:	12						
Covariance Type:	nonrobust						
	coef	std err	t	P> t	[0.025	0.975]	
const	-3927.5420	2487.237	-1.579	0.118	-8868.873	1013.789	
year	1.9397	1.229	1.579	0.118	-0.501	4.381	
month	0.3952	1.663	0.238	0.813	-2.908	3.699	
week	-0.1062	0.384	-0.277	0.783	-0.869	0.657	
total_campaigns	0.1548	0.299	0.519	0.605	-0.438	0.748	
total_impressions	2.857e-05	3.15e-05	0.907	0.367	-3.4e-05	9.11e-05	
total_clicks	-4.306e-05	0.000	-0.137	0.891	-0.001	0.001	
total_conversions	0.0007	0.003	0.201	0.841	-0.006	0.007	
total_spend	3.235e-05	8.9e-05	0.363	0.717	-0.000	0.000	
num_apps	0.0004	0.195	0.002	0.999	-0.388	0.389	
num_emails	0.1214	0.172	0.706	0.482	-0.220	0.463	
num_smss	-0.2187	0.164	-1.337	0.185	-0.544	0.106	
Social Media	0.2518	0.151	1.672	0.098	-0.047	0.551	
roas	1.4488	4.227	0.343	0.733	-6.949	9.847	
	Omnibus:	4.601	Durbin-Watson:	1.519			
Prob(Omnibus):		0.100	Jarque-Bera (JB):	3.980			
Skew:		0.392	Prob(JB):	0.137			
Kurtosis:		3.560	Cond. No.	2.97e+19			

Ols Summary on whole data

- The model explains 70% of the variance in the dependent variable ($R^2 = 0.70$), indicating **a good fit**
- The overall **model is statistically significant ($F = 17.83, p < 0.001$)**, suggesting meaningful predictive power.
- However, most individual predictors are not significant ($p > 0.05$), possibly due to **multicollinearity** or irrelevant features.
- Residuals show mild skewness (0.39) and slightly heavy tails (kurtosis = 3.56)**, indicating **near-normal error distribution**.

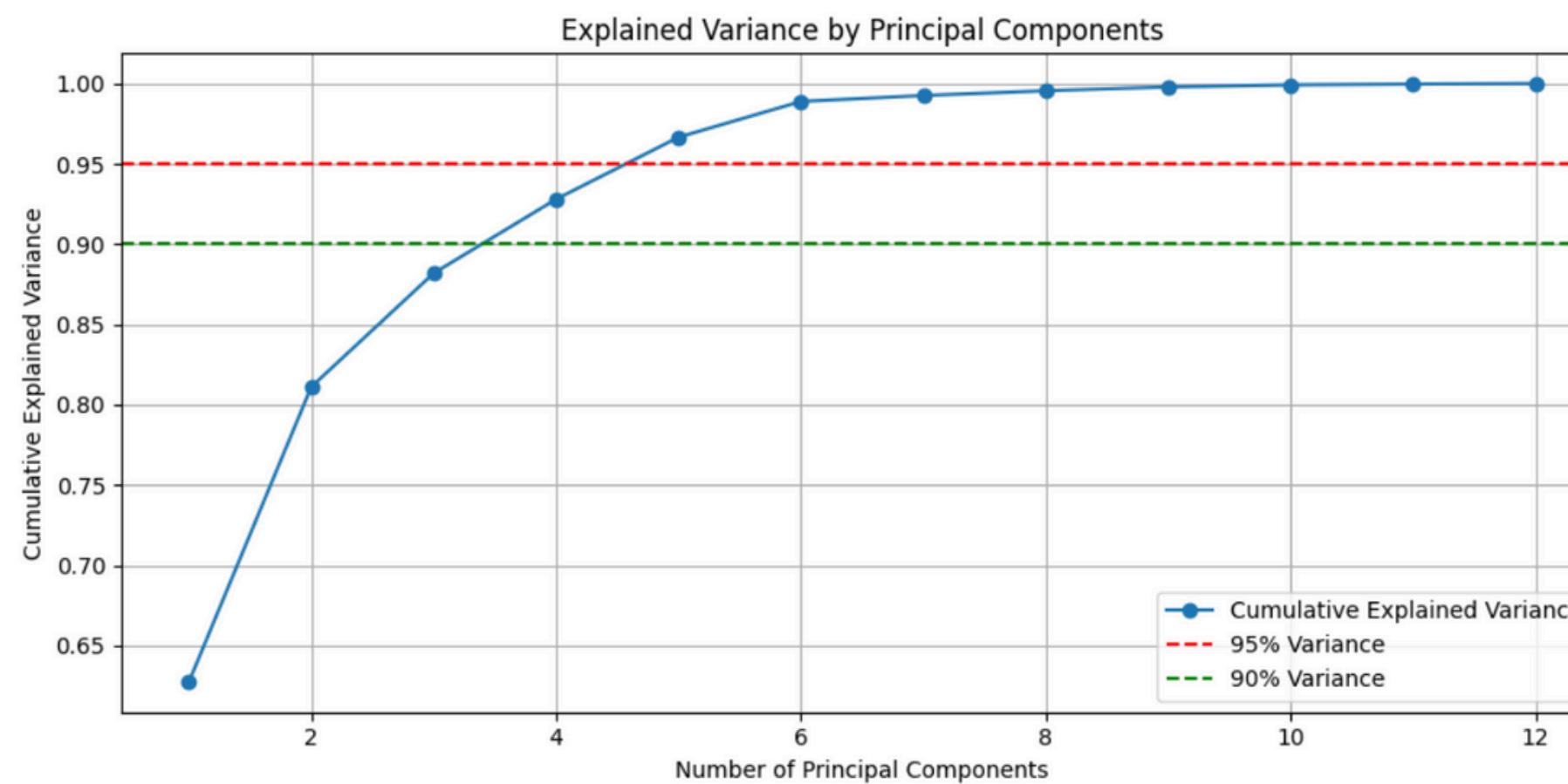
Checking Multicollinearity

	VIF	feature
0	1.298355	year
1	88.248234	month
2	88.413694	week
3	inf	total_campaigns
4	31.572135	total_impressions
5	30.454283	total_clicks
6	32.461501	total_conversions
7	74.201285	total_spend
8	inf	num_apps
9	inf	num_emails
10	inf	num_smss
11	inf	Social_Media
12	1.628076	roas

TOO high Multicollinearity !!

- There is high multicollinearity among the predictors, due to which, even if the model is significant, the predictors turn out to be insignificant.
- We need to tackle Multicollinearity by using
 1. PCA
 2. LASSO
 3. Ridge

Principal Component Analysis



Component 1: Condition Index = 1.00
Component 2: Condition Index = 1.86
Component 3: Condition Index = 2.99
Component 4: Condition Index = 4.07
Component 5: Condition Index = 3.90
Component 6: Condition Index = 5.22
Component 7: Condition Index = 13.63
Component 8: Condition Index = 16.35
Component 9: Condition Index = 18.40
Component 10: Condition Index = 22.27
Component 11: Condition Index = 37.10
Component 12: Condition Index = 99738880.56

First 5 PCA components explain over 95% variance, enabling efficient dimensionality reduction.
Condition indices <12 up to 6 components indicate no serious multicollinearity.

Regularized Linear Model

1. Ridge Regression

After performing 5-fold cross-validation on the Ridge regression model within the 50 values between 10^{-4} to 10^4 , the optimal L2 regularization rate value was selected.

```
Best alpha: {'alpha': 16.768329368110066}
```

2. Lasso Regression

After performing 5-fold cross-validation on the Lasso regression model within the 50 values between 10^{-4} to 10^4 , the optimal L1 regularization rate value was selected.

```
Best alpha: {'alpha': 0.3906939937054613}
```

Comparison of Models

1. Principal Component Analysis

After fitting the Linear Regression model on the first 6 Principal Components , the R 2 score on test data increased a little.

Test R² score: 0.7409606750406832

2. Ridge Regression

After fitting the Ridge model with the optiomal L2 regularization rate , the corresponding R2 score on test data is

Test R² score: 0.6014469391754558

3. Lasso Regression

After fitting the Lasso model with the optiomal L1 regularization rate , the corresponding R2 score on test data is

Test R² score: 0.5810748057817368

Comparison of Models



1. Principal Component Analysis

After fitting the Linear Regression model on the first 6 Principal Components , the R 2 score on test data increased a little.

Test R² score: 0.7409606750406832

2. Ridge Regression

After fitting the Ridge model with the optiomal L2 regularization rate , the corresponding R2 score on test data is

Test R² score: 0.6014469391754558

3. Lasso Regression

After fitting the Lasso model with the optiomal L1 regularization rate , the corresponding R2 score on test data is

Test R² score: 0.5810748057817368

Interpretation of Comparison

- PCA with 6 components gave the best test R^2 (0.74), effectively capturing key variance and reducing noise.
- Ridge ($R^2 = 0.60$) and Lasso ($R^2 = 0.58$) underperformed, likely due to multicollinearity and presence of noisy or irrelevant features.
- Regularization couldn't fully resolve these issues, while PCA transformed the data into uncorrelated components, enhancing model performance.

Conclusion

How Blinkit Can Increase its Revenue?!

- **Delivery Delays = Mood Killers**
 - Add more micro-warehouses (a.k.a. magic stockrooms).
 - Gamify fast deliveries: “Top Delivery Guy of the Week” gets rewards.
- **On-Time Delivery = Customer Satisfaction**
 - Promise “30-min or ₹50 back” — bold, spicy, unforgettable.
 - Track your heroes (delivery partners) and reward the punctual ones
- **More Products = More Reasons to Shop :**
 - Analyze what’s trending and load up (yes, even the weird snacks)
 - Localize inventory — what slaps in Delhi might flop in Chennai.
 - From instant noodles to ice cream — stock it all.

Conclusion

“If you stock more, deliver fast, show up on time, and make every cart count — you’ll not just earn revenue, you’ll earn loyalty, too.”

Thank You!
For your patient
hearing...

